



Published in final edited form as:

Econom J. 2021 September ; 24(3): 559–588. doi:10.1093/ectj/utab019.

Double/debiased machine learning for logistic partially linear model

MOLEI LIU[†], YI ZHANG[‡], DOUDOU ZHOU[§]

[†]Department of Biostatistics, Harvard T.H. Chan School of Public Health, 677 Huntington Avenue, Boston, MA 02115, USA.

[‡]Department of Statistics, Harvard University, One Oxford Street, Cambridge, MA 02138-2901, USA.

[§]Department of Statistics, University of California, Davis, One Shields Avenue, Davis, CA 95616, USA.

Summary:

We propose double/debiased machine learning approaches to infer a parametric component of a logistic partially linear model. Our framework is based on a Neyman orthogonal score equation consisting of two nuisance models for the nonparametric component of the logistic model and conditional mean of the exposure with the control group. To estimate the nuisance models, we separately consider the use of high dimensional (HD) sparse regression and (nonparametric) machine learning (ML) methods. In the HD case, we derive certain moment equations to calibrate the first order bias of the nuisance models, which preserves the model double robustness property. In the ML case, we handle the nonlinearity of the logit link through a novel and easy-to-implement ‘full model refitting’ procedure. We evaluate our methods through simulation and apply them in assessing the effect of the emergency contraceptive pill on early gestation and new births based on a 2008 policy reform in Chile.

Keywords

Logistic partially linear model; double machine learning; double robustness; regularized regression; calibration; C14

1. INTRODUCTION

Consider a logistic partially linear model. Let $\{(Y_i, A_i, \mathbf{X}_i): i = 1, 2, \dots, n\}$ be independent and identically distributed samples of $Y \in \{0, 1\}$, $A \in \mathbb{R}$, and $\mathbf{X} \in \mathbb{R}^p$. Assume that

molei_liu@g.harvard.edu .

SUPPORTING INFORMATION

Additional Supporting Information may be found in the online version of this article at the publisher’s website:

Replication Package

Co-editor Victor Chernozhukov handled this manuscript.

$$\mathbb{P}(Y = 1 \mid A, \mathbf{X}) = \text{expit}\{\beta_0 A + r_0(\mathbf{X})\}, \quad (1.1)$$

where $\text{expit}(\cdot) = \text{logit}^{-1}(\cdot)$, $\text{logit}(a) = \log\{a/(1-a)\}$, and $r_0(\cdot)$ are unknown nuisance functions of \mathbf{X} . In an experimental or observational study with A taken as the exposure variable, Y being the binary response of interest and \mathbf{X} representing the observed confounding variables, parameter β_0 is of particular interest as it measures the conditional effect of A on Y on the scale of the logarithmic odds ratio. As the most common and natural way to characterise the conditional model of a binary outcome against some exposure, model (1.1) has been extensively used in economics and policy science studies.

Our goal is to estimate and asymptotically infer β_0 at the rate $n^{-1/2}$. When \mathbf{X} is a scalar and $r_0(\cdot)$ is smooth, classic semiparametric kernel or sieve regression work well for this purpose (see Severini and Staniswalis, 1994; Lin and Carroll, 2006). When \mathbf{X} is of high dimensionality, these approaches can have poor performance due to the curse of dimensionality. Accordingly, it would be more desirable to estimate $r_0(\cdot)$ with modern high dimensional (HD) (parametric) or machine learning (ML) (nonparametric)¹ methods that are much more resistant to the growing dimensionality and complexity of \mathbf{X} . However, unlike the partially linear model scenario (see Chernozhukov, Chetverikov, et al., 2018; Dukes and Vansteelandt, 2020), robust and efficient inference of β_0 in (1.1) with HD or ML nuisance models has not yet been extensively studied.

In recent literature, Tan (2019a) proposed a simple and flexible doubly robust estimator to enhance the robustness to the potential misspecification of $r(\mathbf{x})$ specified as a fixed-dimensional parametric function: $r(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\gamma}$. The authors introduced a parametric model $m(\mathbf{x}) = g(\mathbf{x}^\top \boldsymbol{\alpha})$ with a (known) link function $g(\cdot)$ for the conditional mean model $m_0(\mathbf{x}) = \mathbb{E}(A \mid Y = 0, \mathbf{X} = \mathbf{x})$ and proposed a doubly robust estimating equation:

$$\frac{1}{n} \sum_{i=1}^n \hat{\phi}(\mathbf{X}_i) \left\{ Y_i e^{-\beta A_i} - \mathbf{X}_i^\top \hat{\boldsymbol{\gamma}} - (1 - Y_i) \right\} \{ A_i - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}}) \} = 0, \quad (1.2)$$

where $\hat{\phi}(\mathbf{x})$ is an estimation of some scalar nuisance function $\phi(\mathbf{x})$ affecting the asymptotic efficiency of the estimator, $\hat{\boldsymbol{\alpha}}$ and $\hat{\boldsymbol{\gamma}}$ are two fixed-dimensional nuisance model estimators. Estimator $\hat{\beta}$ solved from (1.2) is doubly robust in the sense that it is valid when either $r(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\gamma}$ is correctly specified for the logistic model nonparametric component, or when $m(\mathbf{x}) = g(\mathbf{x}^\top \boldsymbol{\alpha})$ is correctly specified for the conditional mean model $m_0(\mathbf{x}) = \mathbb{E}(A \mid Y = 0, \mathbf{X} = \mathbf{x})$. Prior to this, the doubly robust semiparametric estimation of the

¹Our HD setting refers to the HD parametric (linear or generalised linear) model and the ML setting refers to ML models of conditional mean estimation (prediction/classification) that is black box and usually nonparametric.

odds ratio was built on $p(A | X, Y = 0)$, the conditional density of A given X and $Y = 0$ (e.g., Chen, 2007; Tchetgen Tchetgen et al., 2010).²

Nevertheless, Tan (2019a) focused on fixed-dimensional parametric nuisance models that were still prone to model misspecification. Their proposed approach is not readily applicable to the recently developed and exploited HD ($p \gg n$ and the two nuisance components are specified as parametric models with sparse coefficients) or ML (the nuisance functions are estimated by arbitrary black box learning algorithm of condition mean) approaches. For the HD framework, this is because simply using regularised nuisance estimators in (1.2) would typically incur excessive bias and would not guarantee the parametric rate of convergence. In this paper, we realise bias reduction with respect to the regularisation errors by constructing certain Dantzig moment equations to estimate the nuisance parameters. With ultra-sparse nuisance parameters, i.e., sparsity level is $o(n^{1/2}/\log p)$, our proposed estimator preserves the model double robustness property such that it approaches β_0 at the rate $O_p(n^{-1/2})$ when either $r(\cdot)$ or $m(\cdot)$ is correctly specified. Under the ML framework, the nonlinearity and unextractability of the logit link makes it impossible to naturally estimator $r_0(\cdot)$ with a learning algorithm of conditional mean as the partially linear setting. We handle this challenge through an easy-to-implement ‘full model refitting’ (FMR) procedure that facilitates flexible implementation of arbitrary learning algorithms in our framework. Our double machine learning (DML) estimator for β_0 is rated doubly robust in the same sense as Chernozhukov, Chetverikov, et al. (2018), i.e., asymptotically normal at rate $n^{-1/2}$ when the two nuisance ML estimators are consistent for the true models and their root mean squared errors are controlled by $o_p(n^{-1/4})$.

In recent years, there has been a large body of literature developed for semiparametric inference with HD and ML nuisance models, which has garnered increased attention and applications in economics and policy sciences (e.g., Athey and Imbens, 2017; Knaus, 2018, 2020; Yang et al., 2020). As reflected in the organisation of our paper, there are two different frameworks in the literature, the HD and ML settings, both often referred to as ‘machine learning’ approaches. The main difference between them is that the HD setting imposes parametric assumptions on the nuisance models and may allow for their potential misspecification, while the ML setting uses nonparametric ML nuisance estimators that are supposed to approach the true functions at a certain geometric rate.

To estimate low dimensional parameters such as the average treatment effect (ATE) and conditional treatment effect in linear or log-linear models in the HD setting with potentially misspecified nuisance models, recent studies, including Smucler et al.(2019),Tan (2020a,2020b), Ning et al. (2020), and Dukes and Vansteelandt (2020), constructed ℓ_1 -regularised estimating equations with certain ℓ_∞ -constraints to simultaneously estimate the nuisance parameters and calibrate their first order bias. In comparison, Bradic et al. (2019) proposed a more sparsity-rate robust ATE estimator that requires substantially weaker sparsity assumptions but needs both HD parametric nuisance models to be correctly

²Chen (2007) and Tchetgen Tchetgen et al. (2010) estimate the conditional distribution of A given X while Tan (2019a) only needs to specify a conditional mean model of A .

specified. Our work is the first to consider the logistic partially linear model under a similar HD model robustness regime. Existing approaches including debiased (desparsified) LASSO (see Van de Geer et al., 2014; Janková and Van de Geer, 2016) and regularised Riesz representer (see Chernozhukov, Newey and Robins, 2018; Belloni et al., 2018) used the empirical inverse of the information matrix obtained with ℓ_1 -regularised regression to correct for the bias of the logistic LASSO estimator. They imposed a technical ultra-sparsity condition on the inverse of the information matrix, an approach which has been criticised as unreasonable and unverifiable (Xia et al., 2020). In contrast, our sparsity assumption is model-specific and therefore more reasonable and explainable.

We note that near the end of finalizing our paper, a parallel paper by Ghosh and Tan (2020) was published on [arXiv.org](https://arxiv.org). Our HD framework studies the same problem but uses a different doubly robust estimating equation and calibrated procedures for the nuisance models. While the two proposals have similar theoretical properties and numerical implementation strategies (see Sections 3.1, 4.1, and Appendix A4), our method under the ML setting introduced below has no overlap with their work. In parallel, Nekipelov et al. (2018) published a preprint (the latest version in October 2020) about DML estimation of the generalised partially linear model using an adjusted Neyman orthogonal estimating equation different from ours. As a comparison, their approach works for more general target models with arbitrary (smooth) link functions and types of outcome while our method is restricted to the logit link and binary outcome.³ However, their construction cannot achieve the model double robustness property as we can with either fixed or high dimensional parametric nuisance models. In particular, when the nuisance model $r_0(\mathbf{x})$ is misspecified, their estimator is invalid while ours is still valid as long as $m_0(\mathbf{x})$ is correctly specified. This is due to them using $\mathbb{E}(A | \mathbf{X} = \mathbf{x})$ as the nuisance model while we use $\mathbb{E}(A | Y = 0, \mathbf{X} = \mathbf{x})$. See Tan (2019a) for more discussion on this.

Under the nonparametric ML setting, Chernozhukov, Chetverikov, et al. (2018) established a DML framework utilising Neyman orthogonal scores and cross-fitting to construct a parametrically efficient ML-based casual estimator. Their framework has played a central role in semiparametric inference with ML. To complement their approach, recent work, including Chernozhukov, Newey, Robins and Singh (2018), Zimmert and Lechner (2019), and Colangelo and Lee (2020), localised the orthogonal score function to estimate conditional average treatment effect; Semenova and Chernozhukov (2020) constructed the best linear approximation of a structural function with ML; Farrell et al. (2018) used deep neural networks for DML estimation; Wager and Athey (2018) and Oprescu et al. (2019) proposed and studied the tree-based ML approaches for causal inference; and Cui and Tchetgen Tchetgen (2019) proposed a minimax data-driven model selection approach to choose the ML nuisance models with the lowest bias on the DML estimator. The above-mentioned work elaborated on specific inference problems including partially linear models, ATE, and heterogeneous treatment effects with nuisance models which can be directly estimated with arbitrary (supervised) ML algorithms. As mentioned above, $r_0(\cdot)$ in our case cannot be estimated with general ML algorithms due to the nonlinear structure of (1.1). To

³It is still an open problem to generalise our approach to general link functions in M-estimation.

the best of our knowledge, this paper is the first to solve this nontrivial technical problem through our proposed FMR procedure.

The rest of this paper is organised as follows. In Section 2, we define the Neyman orthogonal (doubly robust) score equation for logistic partially linear models. In Section 3, we introduce the realisation of debiased and DML inference for β_0 under the HD and ML settings respectively. In Section 4, we present and justify the asymptotic property of our HD and ML estimators. In Sections 5 and 6, we conduct simulations to study the empirical performance of our method and apply it to assess the effects of the emergency contraceptive (EC) pill on early gestation fetal and new births.

2. NEYMAN ORTHOGONAL SCORE

Before coming to the specific approaches in Section 3, we introduce a Neyman orthogonal (doubly robust) score function for logistic partially linear models and derive its first order bias, which plays a central role in motivating and guiding our methods and theoretical analysis. Let observation $\mathbf{D}_i = \{Y_i, A_i, \mathbf{X}_i\}$ for $i = 1, \dots, n$ and $\mathbf{D} = \{Y, A, \mathbf{X}\}$ be a realisation of \mathbf{D}_i . Motivated by equation (1.2) and Tan (2019a), we define the Neyman orthogonal score as

$$h(\mathbf{D}; \beta, \eta) = \psi(\mathbf{X}) \left\{ Y e^{-\beta A} - (1 - Y) e^{r(\mathbf{X})} \right\} \{A - m(\mathbf{X})\},$$

where $\eta = \{r(\cdot), m(\cdot), \psi(\cdot)\}$ represents the whole set of nuisance functions. Similar to (1.2), $r(\cdot)$ and $m(\cdot)$ correspond to the nonparametric component $r_0(\cdot)$ defined in (1.1) and $m_0(\mathbf{x}) = \mathbb{E}(A|Y = 0, \mathbf{X} = \mathbf{x})$, respectively. $\psi(\mathbf{x})$ is a nuisance function affecting the asymptotic variance of the estimator that may depend on $r(\mathbf{x})$ and $m(\mathbf{x})$ and actually corresponds to $\phi(\mathbf{x})e^{-r(\mathbf{x})}$ with $\phi(\mathbf{x})$ defined by (1.2).⁴

REMARK 2.1.

The score function $h(\mathbf{D}; \beta, \eta)$ is doubly robust in the sense that when $r(\cdot) = r_0(\cdot)$ or $m(\cdot) = m_0(\cdot)$, β_0 solves $\mathbb{E}h(\mathbf{D}; \beta, \eta) = 0$. We shortly demonstrate this as follows. When either $r(\cdot) = r_0(\cdot)$ or $m(\cdot) = m_0(\cdot)$ holds, we have

$$\begin{aligned} & \mathbb{E} \psi(\mathbf{X}) (1 - Y) \left\{ e^{r(\mathbf{X})} - e^{r_0(\mathbf{X})} \right\} \{A - m(\mathbf{X})\} \\ &= \mathbb{E} \left[\psi(\mathbf{X}) \left\{ e^{r(\mathbf{X})} - e^{r_0(\mathbf{X})} \right\} \{A - m(\mathbf{X})\} \mid Y = 0, \mathbf{X} \right] = 0, \end{aligned}$$

which combined with (1.1) leads to that

$$\begin{aligned} \mathbb{E} h(\mathbf{D}; \beta_0, \eta) &= \mathbb{E} \psi(\mathbf{X}) \left\{ Y e^{-\beta_0 A} - (1 - Y) e^{r(\mathbf{X})} \right\} \{A - m(\mathbf{X})\} \\ &= \mathbb{E} \psi(\mathbf{X}) e^{r_0(\mathbf{X})} \left\{ Y e^{-\beta_0 A} - r_0(\mathbf{X}) - (1 - Y) \right\} \{A - m(\mathbf{X})\} \\ &= \mathbb{E} \psi(\mathbf{X}) e^{r_0(\mathbf{X})} \{A - m(\mathbf{X})\} \frac{Y - \mathbb{P}(Y = 1 \mid A, \mathbf{X})}{\mathbb{P}(Y = 1 \mid A, \mathbf{X})} = 0. \end{aligned}$$

⁴We rewrite (1.2) with $\psi(\mathbf{x}) = \phi(\mathbf{x})e^{-r(\mathbf{x})}$ to form the score function so that one could find both its partial derivatives on r and ψ are Neyman orthogonal in a more explicit way.

Suppose the nuisance models $r_0(\mathbf{x})$ and $m_0(\mathbf{x})$ are estimated by $\hat{r}(\mathbf{x})$ and $\hat{m}(\mathbf{x})$ converging to $\bar{r}(\mathbf{x})$ and $\bar{m}(\mathbf{x})$ respectively, and let $\hat{\psi}(\mathbf{x})$ represent the estimator for $\psi(\mathbf{x})$ approaching $\bar{\psi}(\mathbf{x})$. Denote by $\bar{\eta} = \{\bar{r}(\cdot), \bar{m}(\cdot), \bar{\psi}(\cdot)\}$, and $\hat{\eta} = \{\hat{r}(\cdot), \hat{m}(\cdot), \hat{\psi}(\cdot)\}$. We then write the Gateaux (partial) derivative of the score function $h(\mathbf{D}; \beta_0, \bar{\eta})$ as

$$\begin{aligned} & \partial_{\eta} h(\mathbf{D}; \beta_0, \bar{\eta})[\eta - \bar{\eta}] \\ &= \partial_{\psi} h(\mathbf{D}; \beta_0, \bar{\eta})[\psi - \bar{\psi}] + \partial_r h(\mathbf{D}; \beta_0, \bar{\eta})[r - \bar{r}] + \partial_m h(\mathbf{D}; \beta_0, \bar{\eta})[m - \bar{m}] \\ &=: \left\{ Y e^{-\beta_0 A} - (1 - Y) e^{\bar{r}(\mathbf{X})} \right\} \left\{ A - \bar{m}(\mathbf{X}) \right\} \left\{ \psi(\mathbf{X}) - \bar{\psi}(\mathbf{X}) \right\} \\ &\quad - (1 - Y) \bar{\psi}(\mathbf{X}) e^{\bar{r}(\mathbf{X})} \left\{ A - \bar{m}(\mathbf{X}) \right\} \left\{ r(\mathbf{X}) - \bar{r}(\mathbf{X}) \right\} \\ &\quad - \bar{\psi}(\mathbf{X}) \left\{ Y e^{-\beta_0 A} - (1 - Y) e^{\bar{r}(\mathbf{X})} \right\} \left\{ m(\mathbf{X}) - \bar{m}(\mathbf{X}) \right\}. \end{aligned} \tag{2.1}$$

REMARK 2.2.

We evaluate the Neyman orthogonal score on some limiting parameters $\bar{r}(\cdot)$ and $\bar{m}(\cdot)$ instead of on $r_0(\cdot)$ and $m_0(\cdot)$ as in Chernozhukov, Chetverikov, et al.(2018). This is because, different from their ML framework (and ours), assuming both nuisance estimators converge to the true models, i.e., $\bar{r}(\cdot) = r_0(\cdot)$ and $\bar{m}(\cdot) = m_0(\cdot)$, our HD realisation allows at most one nuisance model to be wrongly specified.

Inspired by our deduction in Remark 2.1, $\mathbb{E} \partial_{\psi} h(\mathbf{D}; \beta_0, \bar{\eta})[\psi - \bar{\psi}] = 0$ for any ψ whenever $\bar{r}(\cdot) = r_0(\cdot)$ or $\bar{m}(\cdot) = m_0(\cdot)$. Also, $\mathbb{E} \partial_r h(\mathbf{D}; \beta_0, \bar{\eta})[r - \bar{r}] = 0$ when $\bar{m}(\cdot) = m_0(\cdot)$ and $\mathbb{E} \partial_m h(\mathbf{D}; \beta_0, \bar{\eta})[m - \bar{m}] = 0$ when $\bar{r}(\cdot) = r_0(\cdot)$. Thus, under the ML setting, $h(\mathbf{D}; \beta_0, \bar{\eta})$ satisfies the Neyman orthogonality condition, $\partial_{\eta} h(\mathbf{D}; \beta_0, \bar{\eta})[\eta - \bar{\eta}] = 0$, as defined in Chernozhukov, Chetverikov, et al. (2018) and the first order (over-fitting) bias of $n^{-1} \sum_{i=1}^n h(\mathbf{D}_i; \beta, \hat{\eta})$ can be removed through cross-fitting (introduced in Section 3.2) and concentration. While under the HD parametric setting with $\bar{r}(\cdot) \neq r_0(\cdot)$ or $\bar{m}(\cdot) \neq m_0(\cdot)$, the moment equations for $\bar{r}(\cdot)$ and $\bar{m}(\cdot)$ must be carefully constructed to ensure the orthogonality conditions. Similar to existing literature, such as Chernozhukov, Chetverikov, et al. (2018) and Tan (2020a), removal of the second order (and beyond) bias relies on the assumption on quality of the nuisance estimators $\hat{r}(\cdot)$ and $\hat{m}(\cdot)$ (see Section 4).

3. METHOD

In this section, we separately present our specific construction procedures for HD parametric and ML nonparametric realisation of the debiased/DML estimator for β_0 , based on the derivation and discussion in Section 2.

3.1. High dimensional parametric model realisation

Consider the setting with $p \gg n$, where each \mathbf{X}_i has its first element being 1, and $r(\mathbf{x}) = \mathbf{x}^{\top} \boldsymbol{\gamma}$ and $m(\mathbf{x}) = g(\mathbf{x}^{\top} \boldsymbol{\alpha})$, where $g(\cdot)$ is a monotone and smooth link function with derivative $g'(\cdot)$. Inspired by Smucler et al. (2019), Tan (2020a), and Dukes and Vansteelandt (2020), we construct Dantzig moment equations to ensure the Neyman orthogonality empirically: $\partial_r h(\mathbf{D}; \beta_0, \bar{\eta})[r - \bar{r}] = 0$ and $\partial_m h(\mathbf{D}; \beta_0, \bar{\eta})[m - \bar{m}] = 0$, under potential misspecification of one the nuisance models.

First, we obtain $\tilde{\gamma}$ as some initial estimator for γ that converges to some limiting parameter γ^* equal to the true model parameter γ_0 when $r(\mathbf{x})$ is correctly specified as $r_0(\mathbf{x}) = \mathbf{x}^\top \gamma_0$. Let $\hat{\psi}(\mathbf{x})$ be some estimator of the nuisance function $\psi(\mathbf{x})$ depending on $\tilde{\gamma}$ with its limiting function being $\bar{\psi}(\mathbf{x})$, whose choice will be discussed in Section 3.3. According to (2.1), we obtain $\hat{\alpha}$ through the Dantzig moment equation:

$$\min_{\alpha \in \mathbb{R}^p} \|\alpha\|_1 \quad \text{s.t.} \quad \left\| n^{-1} \sum_{i=1}^n (1 - Y_i) \hat{\psi}(X_i) e^{X_i^\top \tilde{\gamma}} \{A_i - g(X_i^\top \alpha)\} X_i \right\|_\infty \leq \lambda_\alpha, \quad (3.1)$$

where λ_α is a tuning parameter controlling the regularisation bias. Finally, we obtain the nuisance estimator $\hat{\gamma}$ and the targeted HD estimator $\hat{\beta}_{\text{HD}}$ simultaneously from:

$$\begin{aligned} \min_{\beta \in \mathbb{R}, \gamma \in \mathbb{R}^p} \|\gamma\|_1 \quad \text{s.t.} \quad & \left\| n^{-1} \sum_{i=1}^n \hat{\psi}(X_i) \left\{ Y_i e^{-\beta A_i} - (1 - Y_i) e^{X_i^\top \gamma} \right\} g'(X_i^\top \hat{\alpha}) X_i \right\|_\infty \leq \lambda_\gamma; \\ & n^{-1} \sum_{i=1}^n \hat{\psi}(X_i) \left\{ Y_i e^{-\beta A_i} - (1 - Y_i) e^{X_i^\top \gamma} \right\} \{A_i - g(X_i^\top \hat{\alpha})\} = 0. \end{aligned} \quad (3.2)$$

Let the limits of $\{\hat{\alpha}, \hat{\gamma}\}$ be $\{\bar{\alpha}, \bar{\gamma}\}$ and $\bar{\eta} = \{\bar{r}(\cdot), \bar{m}(\cdot), \bar{\psi}(\cdot)\}$ where $\bar{r}(\mathbf{x}) = \mathbf{x}^\top \bar{\gamma}$, $\bar{m}(\mathbf{x}) = g(\mathbf{x}^\top \bar{\alpha})$ and $\bar{\psi}(\mathbf{x})$ as given in Section 3.3. We will comment on the orthogonality (moment) conditions of our proposal in Remark 3.1, compare our method with Dukes and Vansteelandt (2020) in Remark 3.2, and discuss its numerical implementation with a weighted LASSO formation in Remark 3.3.

REMARK 3.1.—Neglect the second order error terms for now. When $r(\mathbf{x})$ is correct (see Assumption HD1), i.e., $r_0(\mathbf{x}) = \mathbf{x}^\top \gamma_0$, it naturally holds that $\mathbb{E} \partial_m h(\mathbf{D}; \beta_0, \bar{\eta})[\hat{m} - m_0] = 0$ and $\gamma^* = \bar{\gamma} = \gamma_0$. Then the ℓ_∞ -constraint in (3.1) imposes that

$$\mathbb{E} \partial_m h(\mathbf{D}; \beta_0, \bar{\eta})[\hat{r} - r_0] \approx \mathbb{E} (1 - Y) \bar{\psi}(X) e^{X^\top \gamma_0} \{A - g(X^\top \bar{\alpha})\} X^\top (\hat{\gamma} - \gamma_0) = \mathbf{0}^\top (\hat{\gamma} - \gamma_0).$$

When $\bar{m}(\mathbf{x}) = m_0(\mathbf{x}) = g(\mathbf{x}^\top \alpha_0)$, we have $\mathbb{E} \partial_m h(\mathbf{D}; \beta_0, \bar{\eta})[\hat{r} - r_0]$ and $\bar{\alpha} = \alpha_0$ in turn. And the ℓ_∞ -constraint of (3.2) results in

$$\begin{aligned} \mathbb{E} \partial_m h(\mathbf{D}; \beta_0, \bar{\eta})[\hat{m} - m_0] &\approx \\ \mathbb{E} \bar{\psi}(X) \left\{ Y e^{-\beta_0 A} - (1 - Y) e^{X^\top \bar{\gamma}} \right\} g'(X^\top \alpha_0) X^\top (\hat{\alpha} - \alpha_0) & \\ = \mathbf{0}^\top (\hat{\alpha} - \alpha_0) & \end{aligned}$$

Thus, the Neyman orthogonality condition $\partial_\eta h(\mathbf{D}; \beta_0, \bar{\eta})[\eta - \bar{\eta}] = 0$ as introduced in Section 2 is satisfied under our construction when either $r(\cdot)$ or $m(\cdot)$ is correctly specified.

REMARK 3.2.—Similar to the HD partially linear (or log-linear) setting studied in Dukes and Vansteelandt (2020), estimating an equation for the nuisance parameter γ involves the unknown β . Unlike their construction procedure that plugs in β_0 as every $\beta \in \mathbb{R}$ to estimate γ and invert the resulted score-test p -values for interval estimation of β_0 , we solve for $\hat{\beta}$ and $\hat{\gamma}$ jointly from (3.2), the Dukes and Vansteelandt method, our approach is more friendly in terms of computation and with a doubly robust moment equation for $\hat{\beta}$, as demonstrated in Remark 2.1. Compared with implementation, directly provides a point estimator, and preserves a similar theoretical guarantee (see Section 4.1).

REMARK 3.3.—As is detailed in Appendix A4, one can construct LASSO problems with the same Karus–Kuhn–Tucker (KKT) conditions as the ℓ_∞ -norm constraints in (3.1) and (3.2) to obtain the estimators $\hat{\alpha}$ and $\hat{\gamma}$, which have equivalent theoretical properties as (3.2), numerical solution of the LASSO counterpart of (3.2) cannot be obtained with existing software, such as the **R** packages ‘glmnet’ (Friedman et al., 2010) and ‘RCAL’ (Tan, 2019b). A direct solution to this is programming an optimisation procedure such as the Fisher scoring descent algorithm used by Tan (2020b). We also found a convenient way to moderately modify the construction procedure to make the regularised estimating equations solvable with **R** package ‘RCAL’, and use it for numerical implementation. In Appendix A4, we outline this modification and demonstrate its theoretical guarantee.

3.2. Machine learning realisation

We turn now to a (nonparametric) ML setting under which any learning algorithms of conditional mean could potentially be applied to estimate the nuisance functions. Similar to Chernozhukov, Chetverikov, et al. (2018), we randomly split the n samples into K folds: $\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_K$ of equal size, to assist with removing the first order (over-fitting) bias through concentration. Then the cross-fitted estimating equation for β is constructed as

$$n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{F}_k} h(\mathbf{D}_i; \beta, \hat{\eta}^{[-k]}) = 0, \quad (3.3)$$

where $\hat{\eta}^{[-k]} = \{\hat{r}^{[-k]}(\cdot), m^{[-k]}(\cdot), \psi^{[-k]}(\cdot)\}$, representing ML estimators converging to $\bar{r}(\cdot) = r_0(\cdot)$, $\bar{m}(\cdot) = m_0(\cdot)$, and $\bar{\psi}(\cdot)$ obtained with the samples in $\mathcal{F}_{-k} = \{1, \dots, n\} \setminus \mathcal{F}_k$ and independent of the samples in \mathcal{F}_k . Now we present the specific construction procedures for $\hat{r}^{[-k]}(\cdot)$ and $m^{[-k]}(\cdot)$ with the choice of $\psi^{[-k]}(\cdot)$ again discussed in Section 3.3.

Suppose there is a black box learning algorithm $\mathcal{L}(R, C; \mathcal{F})$ that inputs samples $\mathcal{F} \subseteq \{1, 2, \dots, n\}$ with some response R_i and covariates C_i , and outputs an estimation of $E[R_i | C_i, i \in \mathcal{F}]$. We outline our approach utilising \mathcal{L} to estimate the nuisance functions. Corresponding to the definition of $m_0(\cdot)$, it can be estimated by: $\hat{m}^{[-k]}(\cdot) = \mathcal{L}(A, \mathbf{X}; \mathcal{F}_{-k} \cap \{i: Y_i = 0\})$. Compared to the partially linear setting in Chernozhukov, Chetverikov, et al. (2018), estimation of $r_0(\cdot)$ with \mathcal{L} is more sophisticated

as it is defined through a nonlinear form: $\mathbb{P}(Y = 1 | A, \mathbf{X}) = \text{expit}\{\beta_0 A + r_0(\mathbf{X})\}$. One could modify some ML approaches, e.g., neural network⁵ to accommodate this form. However, such modification is not readily available in general and typically requires additional human efforts for its implementation and validation. So, alternatively, we propose a ‘full model refitting’ (FMR) procedure that can use an arbitrary \mathcal{L} to estimate $r_0(\cdot)$. Our method is motivated by a simple proposition:

PROPOSITION 3.1.—*Let $M_0(A, \mathbf{X}) = \mathbb{P}(Y = 1 | A, \mathbf{X}) = \text{expit}\{\beta_0 A + r_0(\mathbf{X})\}$. We have:*

$$\beta_0 = \underset{\beta \in \mathbb{R}}{\text{argmin}} \mathbb{E}[\text{logit}\{M_0(A, \mathbf{X})\} - \beta(A - \mathbb{E}[A | \mathbf{X}])]^2.$$

Proof.: For any $\beta \in \mathbb{R}$, we have

$$\begin{aligned} \mathbb{E}[\text{logit}\{M_0(A, \mathbf{X})\} - \beta(A - \mathbb{E}[A | \mathbf{X}])]^2 &= \mathbb{E}\{\beta_0 A + r_0(\mathbf{X}) - \beta(A - \mathbb{E}[A | \mathbf{X}])\}^2 \\ &= \mathbb{E}\{(\beta_0 - \beta)(A - \mathbb{E}[A | \mathbf{X}]) + \tau(\mathbf{X})\}^2 = (\beta_0 - \beta)^2 \mathbb{E}(A - \mathbb{E}[A | \mathbf{X}])^2 + \mathbb{E}\{\tau(\mathbf{X})\}^2, \end{aligned}$$

where $\tau(\mathbf{X}) = r_0(\mathbf{X}) + \beta_0 \mathbb{E}[A | \mathbf{X}]$. Thus, β_0 minimises

$$\mathbb{E}[\text{logit}\{M_0(A, \mathbf{X})\} - \beta(A - \mathbb{E}[A | \mathbf{X}])]^2.$$

□

Further, randomly split \mathcal{F}_{-k} into K folds $\mathcal{F}_{-k,1}, \dots, \mathcal{F}_{-k,K}$, of equal size and denote by $\mathcal{F}_{-k,-j} = \mathcal{F}_{-k} \setminus \mathcal{F}_{-k,j}$. Motivated by Proposition 3.1, we first estimate the ‘full’ model $M_0(A, \mathbf{X})$ with $\mathcal{F}_{-k,-j}$ as:

$$\widehat{M}^{[-k, -j]}(\cdot) = \mathcal{L}\left(Y_i, (A_i, \mathbf{X}_i^\top)^\top; \mathcal{F}_{-k,-j}\right),$$

and learn $a_0(\mathbf{x}) = \mathbb{E}[A | \mathbf{X} = \mathbf{x}]$ by $\widehat{a}^{[-k, -j]}(\cdot) = \mathcal{L}(A_i, \mathbf{X}_i; \mathcal{F}_{-k,-j})$. Then we fit the (cross-fitted) least square regression to obtain:

$$\widehat{\beta}^{[-k]} = \underset{\beta \in \mathbb{R}}{\text{argmin}} \frac{1}{|\mathcal{F}_{-k}|} \sum_{j=1}^K \sum_{i \in \mathcal{F}_{-k,j}} \left[\text{logit}\left\{\widehat{M}^{[-k, -j]}(A_i, \mathbf{X}_i)\right\} - \beta\{A_i - \widehat{a}^{[-k, -j]}(\mathbf{X}_i)\} \right]^2, \quad (3.4)$$

as an estimator approaching β_0 at certain rate typically larger than $n^{-1/2}$. Then $r_0(\cdot)$ could be identified through $r_0(\mathbf{X}_i) = \text{logit}\{M_0(A_i, \mathbf{X}_i)\} - \beta_0 A_i$. Note that the empirically estimated version of $\text{logit}\{M_0(A_i, \mathbf{X}_i)\} - \beta_0 A_i$ typically involves A_i due to the discrepancy of β_0 and $M_0(\cdot)$ from their empirical estimation. This can essentially impede removal of the over-

⁵By setting the last layer of the neural network to be the combination of a complex network of \mathbf{X} and a linear function of A and linking it with the outcome through an expit link.

fitting bias given that $\partial_\beta h(\mathbf{D}; \beta_0, \bar{\eta})$ is not orthogonal to the error functions depending on A . So we further estimate the conditional mean of $\logit\{M_0(A, \mathbf{X})\} - \beta_0 A$ on \mathbf{X} to obtain the final estimator for $r_0(\cdot)$. Denote by $W_i = \logit\{\widehat{M}^{[-k, -j]}(A_i, \mathbf{X}_i)\}$ for each $i \in \mathcal{I}_{-k, j}$ and get $\hat{t}^{[-k]}(\cdot) = \mathcal{L}(W_i, \mathbf{X}_i; \mathcal{I}_{-k})$ to estimate $t_0(\mathbf{x}) = \mathbb{E}[\logit\{M_0(A, \mathbf{X})\} \mid \mathbf{X} = \mathbf{x}]$. Then the estimator of $r_0(\cdot)$ is given by:

$$\hat{r}^{[-k]}(\mathbf{x}) = \hat{t}^{[-k]}(\mathbf{x}) - \hat{\beta}^{[-k]} \hat{a}^{[-k]}(\mathbf{x}), \quad \text{where} \quad \hat{a}^{[-k]}(\mathbf{x}) = \frac{1}{K} \sum_{j=1}^K \hat{a}^{[-k, -j]}(\mathbf{x}). \quad (3.5)$$

Alternatively, one can estimate $r_0(\cdot)$ through

$$\hat{r}^{[-k]}(\cdot) = \log \left(\frac{\mathcal{L}\left(e^{-\hat{\beta}^{[-k]} A_i, \mathbf{X}_i; \mathcal{I}_{-k} \cap \{i: Y_i = 1\}}\right)}{\mathcal{L}(1 - Y_i, \mathbf{X}_i; \mathcal{I}_{-k})} \right).$$

motivated by the moment condition that is sufficient to identify $r_0(\cdot)$:

$$\mathbb{E}\left[Y e^{-\beta_0 A} - (1 - Y) e^{r_0(\mathbf{X})} \mid \mathbf{X}\right] = \mathbb{E}\left[e^{-\beta_0 A} \mid \mathbf{X}, Y = 1\right] - e^{r_0(\mathbf{X})} \mathbb{E}[(1 - Y) \mid \mathbf{X}] = 0.$$

We refer to the estimation step for $\hat{\beta}^{[-k]}$ and $\hat{r}^{[-k]}(\cdot)$ introduced above as ‘refitting’, and the whole procedure as FMR considering that we ‘refit’ the least square problem (3.4) and ML models \mathcal{L} to estimate $r_0(\cdot)$ with the initially estimated full model $\logit\{M_0(A, \mathbf{X}_i)\}$ as a pseudo-outcome. Finally, we solve (3.3) based on $\hat{\eta}^{[-k]}$ to obtain the DML estimator $\hat{\beta}_{\text{ML}}$.

REMARK 3.4.—We further use cross-fitting in FMR to avoid over-fitting of the models $\widehat{M}^{[-k, -j]}(\cdot)$ and $a^{[-k, -j]}(\cdot)$ when they are used to obtain the estimators $\hat{\beta}^{[-k]}$, $\hat{t}^{[-k]}(\mathbf{x})$ and $\hat{r}^{[-k]}(\mathbf{x})$. This is supposed to show empirical improvement of the FMR procedure.

REMARK 3.5.—The FMR implicitly assumes that \mathcal{L} should perform similarly well on different learning objects with the covariates set as either \mathbf{X} or $(A, \mathbf{X}^\top)^\top$. Classic nonparametric approaches like kernel smoothing or sieve may not satisfy this assumption whereas including one more covariate A in addition to the very low dimensional \mathbf{X} can have substantial impact on estimation performance. Thus, we recommend using more dimensionality-robust modern ML approaches, such as random forest and neural networks, in our ML framework. The classic ‘plug-in’ sieve or kernel method has been well-studied in existing literature (e.g., Severini and Staniswalis (1994); Lin and Carroll (2006)).

3.3. Efficiency considerations

The nuisance function $\psi(\cdot)$ in our framework is included and chosen in consideration of estimation efficiency. Tan (2019a) proposed and studied two options for $\phi(\cdot)$ used and defined in (1.2), with the corresponding function $\psi(\cdot)$ taken as:

$$\psi_{\text{opt}}(\mathbf{x}) = \frac{e^{-r(\mathbf{x})} \mathbb{E} \left[\{A - m(\mathbf{X})\}^2 \mid \mathbf{X} = \mathbf{x}, Y = 0 \right]}{\mathbb{E} \left[\{A - m(\mathbf{X})\}^2 / \text{expit} \{ \beta_0 A + r(\mathbf{X}) \} \mid \mathbf{X} = \mathbf{x}, Y = 0 \right]},$$

and $\psi_{\text{simp}}(\mathbf{x}) = \text{expit} \{ -r(\mathbf{x}) \}$.

REMARK 3.6.—It was shown in Tan (2019a) that when both nuisance models are correctly specified, the estimator solved with the weight $\psi_{\text{opt}}(\cdot)$ achieves the minimum asymptotic variance among all the doubly robust estimators obtained through our estimating equation. Given that the weighting function $\psi(\cdot)$ is independent of the exposure A , our estimating equations form a strict subset of all the score equations for the logistic partial model defined by (1.1). Thus, our estimator is generally less efficient than the semiparametric efficient estimator obtained through the maximum likelihood approach like Tchetgen Tchetgen et al. (2010).

Though $\psi_{\text{opt}}(\cdot)$ is the optimal choice in consideration of the efficiency under our construction, computation of $\psi_{\text{opt}}(\cdot)$ involves numerical integration with respect to \mathbf{X} given $Y = 0$, making it sometimes inconvenient to implement. So Tan (2019a) proposed a simplified but reasonable choice $\psi_{\text{simp}}(\mathbf{x})$ obtained by evaluating $\psi_{\text{opt}}(\mathbf{x})$ at $\beta_0 = 0$. In the following theoretical and numerical studies, we stick to $\psi(\mathbf{x}) = \psi_{\text{simp}}(\mathbf{x})$, $\hat{\psi}(\mathbf{x}) = \text{expit}(-\mathbf{x}^\top \hat{\gamma})$ and correspondingly $\bar{\psi}(\mathbf{x}) = \text{expit}(-\mathbf{x}^\top \gamma^*)$ under the HD setting, and $\hat{\psi}^{[-k]}(\mathbf{x}) = \text{expit} \{ -\hat{r}^{[-k]}(\mathbf{x}) \}$ and $\bar{\psi}(\mathbf{x}) = \text{expit} \{ -r_0(\mathbf{x}) \}$ under the ML setting. Our theoretical framework allows for other choices on $\psi(\cdot)$ as will be discussed in Section 4.2.

4. ASYMPTOTIC ANALYSIS

Let $o(\alpha_n)$, $O(\alpha_n)$, $\omega(\alpha_n)$, $\Omega(\alpha_n)$, and $\Theta(\alpha_n)$ represent the sequences growing at a smaller, equal/smaller, larger, equal/larger, and equal rate of α_n , respectively. Let o_p , O_p , ω_p , Ω_p , and Θ_p be the corresponding rates with probability approaching 1 as $n \rightarrow \infty$. Let $\mathcal{X} \subseteq \mathbb{R}^p$ be the domain of \mathbf{X} . First, we introduce the regularity condition for β and its estimating equation used under both HD and ML settings as Assumption REG, which is standard and can be commonly found in literature of the asymptotic analysis of M -estimator (see Van der Vaart, 2000, ch. 5). We will then study the asymptotic properties of $\hat{\beta}_{\text{HD}}$ and $\hat{\beta}_{\text{ML}}$ in Sections 4.1 and 4.2.

ASSUMPTION REG (REGULARITY OF ESTIMATING EQUATION).

Parameter β belongs to a compact set $\mathcal{B} \subseteq \mathbb{R}$ and there exists $\delta_n = \Omega(n^{-1/2} \log n)$ such that $(\beta_0 - \delta_n, \beta_0 + \delta_n) \subseteq \mathcal{B}$. Exposure A belongs to a compact set \mathcal{A} and $\sup_{x \in \mathcal{X}} |\mathbb{E}[A | \mathbf{X} = x, Y = y]| = O(1)$ for $y = 0, 1$. In addition,⁶

$$\mathbb{E} \bar{\psi}(\mathbf{X}) Y e^{-\beta_0 A} \{A - \bar{m}(\mathbf{X})\} = \Theta(1) \quad \text{and} \quad \mathbb{E} h^2(\mathbf{D}; \beta_0, \bar{\eta}) = \Theta(1).$$

4.1. High dimensional (parametric) setting

Let $\text{expit}'(\cdot)$ be the derivative function of $\text{expit}(\cdot)$, $\|\cdot\|_0$ represents the number of nonzero elements in a vector and $s = \max\{\gamma^*_0, \|\bar{\boldsymbol{\gamma}}\|_0, \|\bar{\boldsymbol{\alpha}}\|_0\}$. We introduce the following assumptions to regularise the covariates and nuisance estimators.

ASSUMPTION HD1 (MODEL DOUBLE ROBUSTNESS).—At least one of the following conditions hold: (a) there exists $\boldsymbol{\gamma}_0 \in \mathbb{R}^p$ such that $r_0(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\gamma}_0$ and $\boldsymbol{\gamma}^* = \bar{\boldsymbol{\gamma}} = \boldsymbol{\gamma}_0$; (b) there exists $\boldsymbol{\alpha}_0 \in \mathbb{R}^p$ such that $m_0(\mathbf{x}) = g(\mathbf{x}^\top \boldsymbol{\alpha}_0)$ and $\bar{\boldsymbol{\alpha}} = \boldsymbol{\alpha}_0$.

ASSUMPTION HD2 (CONCENTRATION RATE).—It holds that

$$\begin{aligned} \left\| n^{-1} \sum_{i=1}^n (1 - Y_i) \text{expit}(-\mathbf{X}_i^\top \boldsymbol{\gamma}^*) e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \{A_i - g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}})\} \mathbf{X}_i \right\|_\infty &= O_p\{(\log p/n)^{1/2}\}; \\ \left\| n^{-1} \sum_{i=1}^n \text{expit}(-\mathbf{X}_i^\top \boldsymbol{\gamma}^*) g'(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) \Psi_i \mathbf{X}_i \right\|_\infty &= O_p\{(\log p/n)^{1/2}\}; \\ \left\| n^{-1} \sum_{i=1}^n \text{expit}'(-\mathbf{X}_i^\top \boldsymbol{\gamma}^*) \{A_i - g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}})\} \Psi_i \mathbf{X}_i \right\|_\infty &= O_p\{(\log p/n)^{1/2}\}, \end{aligned}$$

where $\Psi_i = Y_i e^{-\beta_0 A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}}$.

ASSUMPTION HD3 (SMOOTH LINK FUNCTION).—There exists $L = \Theta(1)$ that for any $u, v \in \mathbb{R}$,

$$|g'(u) - g'(v)| \leq L|u - v|.$$

ASSUMPTION HD4 (RISK OF THE L1-REGULARISED ESTIMATORS).—**L1-REGULARISED ESTIMATORS.**—There exists tuning parameters $\lambda_n, \lambda_\gamma = \Theta\{(\log p/n)^{1/2}\}$ such that (3.1) and (3.2) have feasible solutions with probability approaching 1 and

$$\sup_{i \in \{1, \dots, n\}} |g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}})| = O_p(1); \quad \|\tilde{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*\|_1 + \|\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}}\|_1 + \|\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}}\|_1 = O_p\{s(\log p/n)^{1/2}\};$$

⁶To accommodate the notations of both HD and ML, we use $\bar{m}(\cdot)$ and $\bar{r}(\cdot)$ to represent the limiting models defined as $g(\mathbf{x}^\top \bar{\boldsymbol{\alpha}})$ and $\mathbf{x}^\top \bar{\boldsymbol{\gamma}}$ under the HD setting and just the true models $m_0(\cdot)$ and $r_0(\cdot)$ under ML.

$$n^{-1} \sum_{i=1}^n \left\{ 1 + e^{\mathbf{X}_i^\top \tilde{\gamma}} \left[\left\{ \mathbf{X}_i^\top (\hat{\gamma} - \tilde{\gamma}) \right\}^2 + \left\{ \mathbf{X}_i^\top (\tilde{\gamma} - \gamma^*) \right\}^2 \right] + (\hat{\beta}_{\text{HD}} - \beta_0)^2 \right\} = O_p(s \log p / n);$$

$$n^{-1} \sum_{i=1}^n \left\{ 1 + e^{\mathbf{X}_i^\top \tilde{\gamma}} \left[\left\{ \mathbf{X}_i^\top (\hat{\alpha} - \bar{\alpha}) \right\}^2 + \left\{ g(\mathbf{X}_i^\top \hat{\alpha}) - g(\mathbf{X}_i^\top \bar{\alpha}) \right\}^2 \right] \right\} = O_p(s \log p / n).$$

ASSUMPTION HD5 (ULTRA-SPARSITY).—It holds that $s = o(n^{1/2} / \log p)$.

REMARK 4.1.—Under Assumption HD1 and our constructions (3.1) and (3.2) (or the one introduced in Appendix A4), the expectations of the terms to be concentrated in Assumption HD2 are $\mathbf{0}$ by Remark 3.1. Then their maximum norms can be controlled by $O_p\{(\log p / n)^{1/2}\}$ as assumed in HD2, when the covariates \mathbf{X}_i are bounded, sub-exponential, or beyond (Kuchibhotla and Chakraborty, 2018), using the concentration results derived from the existing literature (Giné and Nickl, 2016).

REMARK 4.2.—Rates of the prediction and estimation risk of the nuisance estimators in Assumption HD4 can be derived following the general theoretical framework for ℓ_1 -regularised estimation introduced in Candès et al. (2007), Bickel et al. (2009), Bühlmann and Van de Geer (2011), and Negahban et al. (2012). The same rate properties have been used for analysing doubly robust estimators with HD nuisance models in existing literature (Smucler et al., 2019; Tan, 2020a; Dukes and Vansteelandt, 2020). Note that (3.1) and (3.2) involve the estimators $\tilde{\gamma}$ or $\hat{\alpha}$ being obtained beforehand. This will require some additional effort on removing the ‘plug-in’ errors of $\tilde{\gamma}$ or $\hat{\alpha}$ when deriving the risk rates for $\hat{\alpha}$ or $\tilde{\gamma}$ compared to the standard analysis procedures. See Tan (2020a) for a similar issue and the relevant technical details used to handle it. In addition, $\sup_{i \in \{1, \dots, n\}} |g(\mathbf{X}_i^\top \hat{\alpha})| = O_p(1)$ imposed in HD4 is not a standard assumption but is rather mild. This is because $\sup_{\mathbf{x} \in \mathcal{X}} |g(\mathbf{x}^\top \bar{\alpha})| = O(1)$ by Assumption REG and we only need $g(\mathbf{x}^\top \hat{\alpha}) - g(\mathbf{x}^\top \bar{\alpha})$ to be $O_p(1)$ uniformly.

REMARK 4.3.—The ultra-sparsity assumption HD5 was also imposed in existing literature, including Tan (2020a) and Dukes and Vansteelandt (2020), to control the rate of bias incurred by the HD estimators: $O_p(s \log p / n)$ below the parametric rate. For the linear nuisance model, existing work like Zhu et al. (2018) and Dukes and Vansteelandt (2020) suggested to add additional moment (KKT) constraints to relax the ultra-sparsity assumption. However, their approach has not yet been shown to be feasible for nonlinear models, so, while promising, remains unclear for our framework.

We present the asymptotic property of $\hat{\beta}_{\text{HD}}$ in Theorem 4.2 and its proof in Appendix A2.

THEOREM 4.1.—Denote by $\bar{I} = \mathbb{E} \bar{\psi}(\mathbf{X}) Y e^{-\beta_0 A} \{A - \bar{m}(\mathbf{X})\}$ and $\bar{\sigma}^2 = \bar{I}^{-2} \mathbb{E} h^2(\mathbf{D}; \beta_0, \bar{\eta})$. Under Assumptions REG and HD1–HD5, we have

$$\sqrt{n}\bar{\sigma}^{-1}(\hat{\beta}_{\text{HD}} - \beta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\bar{\sigma}^{-1})^{-1} h(\mathbf{D}_i; \beta_0, \bar{\eta}) + o_{\mathbb{P}}(1),$$

which weakly converge to $N(0, 1)$.

As an earlier proposed and commonly used approach, logistic debiased LASSO (Van de Geer et al., 2014; Jankova and Van de Geer, 2016) has been criticised because its sparse inverse information matrix condition is not explainable and justifiable, leading to a subpar performance theoretically and numerically (see Xia et al., 2020). Interestingly, we find the model sparsity assumption of our method is more reasonable than debiased LASSO and present a simple comparison of these two approaches in Remark 4.4.

REMARK 4.4.—Assume $\mathbb{P}(Y = 1 | A, \mathbf{X}) = \text{expit}\{\beta_0 A + \mathbf{X}^\top \boldsymbol{\gamma}_0\}$ is correctly specified. As is argued by Xia et al. (2020), assuming the information matrix of the logistic model has an ultra-sparse⁷ inverse, it is crucial to ensure the desirable properties of the debiased LASSO estimator for β_0 . However, this assumption is not explainable or convincing for the common Gaussian design with sparse precision matrix, due to the presence of the logistic canonical link. In comparison, we require that $\mathbb{E}(A | Y = 0, \mathbf{X} = \mathbf{x}) = g(\mathbf{X}^\top \boldsymbol{\alpha}_0)$ with $\|\boldsymbol{\alpha}_0\|_0 = o(n^{1/2}/\log p)$, which has two advantages over debiased LASSO. First, it accommodates nonlinear link function $g(\cdot)$ and can be more reasonable for a categorical A . Second, it is imposed on a conditional model directly and is thus more explainable. For example, consider a conditional Gaussian model: $(A, \mathbf{X}^\top)^\top | \{Y = j\} \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$ for $j = 0, 1$. Then we have $r_0(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\gamma}_0$ where $(\beta_0, \boldsymbol{\gamma}_0)^\top = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$, and $A | \mathbf{X}, Y = 0$ follows a Gaussian linear model with the coefficient $\boldsymbol{\alpha}_0$ determined by $\boldsymbol{\Sigma}^{-1}$. Therefore, our sparsity assumptions on $\boldsymbol{\alpha}_0$ and $\boldsymbol{\gamma}_0$ actually assume the data generation parameters $\boldsymbol{\Sigma}^{-1}$ and $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0$ to be sparse, which seems more explainable and verifiable in practice.

4.2. Machine learning (nonparametric) setting

Define $\|f(\cdot)\|_{\mathcal{Q}, q} = : \|f(U)\|_{\mathcal{Q}, q} = : \left\{ \int |f(u)|^q d\mathcal{Q}(u) \right\}^{1/q}$ for any real number $q > 0$, function $f(\cdot)$, random variables U and probability measure \mathcal{Q} . Let P denote the probability measure of the observed \mathbf{D} . We assume that $K = \Theta(1)$ and introduce the following assumption.

ASSUMPTION ML1 (QUALITY OF THE ML NUISANCE ESTIMATORS).—For each k ,

$$\begin{aligned} \sup_{\mathbf{x} \in \mathcal{X}} \left| \hat{r}^{[-k]}(\mathbf{x}) - r_0(\mathbf{x}) \right| + \left| \hat{m}^{[-k]}(\mathbf{x}) - m_0(\mathbf{x}) \right| &= o_{\mathbb{P}}(1); \\ \left\| \hat{r}^{[-k]}(\cdot) - r_0(\cdot) \right\|_{P, 2} + \left\| \hat{m}^{[-k]}(\cdot) - m_0(\cdot) \right\|_{P, 2} &= o_{\mathbb{P}}(n^{-1/4}). \end{aligned}$$

⁷Or approximately sparse (see recent work, e.g., Belloni et al., 2018; Ma et al., 2020; Liu et al., 2020).

REMARK 4.5.—Similar to Assumptions 3.2 and 3.4 of Chernozhukov, Chetverikov, et al. (2018), our Assumption ML1 requires that the ML estimators for $r_0(\cdot)$ and $m_0(\cdot)$ are uniformly consistent and their mean squared errors (MSE) achieve the rate $O_p(n^{-1/4})$. This assumption is also referred to as rate double robustness in Smucler et al. (2019) as it requires production of the MSEs of $\hat{r}^{[-k]}(\cdot)$ and $\hat{m}^{[-k]}(\cdot)$ to be $o_p(n^{-1/2})$. In Appendix A1, we provide justification for our proposed FMR procedure to derive that the resulted $\hat{r}^{[-k]}(\cdot)$ satisfies Assumption ML1 as long as the learning algorithm \mathcal{L} satisfies the same strong convergence properties as assumed in ML1 on all the learning tasks in FMR. Thus, FMR does not actually clip the wings of the ML algorithms being used in our framework.

We present the asymptotic property of $\hat{\beta}_{\text{ML}}$ in Theorem 4.2 with its proof found in Appendix A3.

THEOREM 4.2.—Denote by $I_0 = \mathbb{E}\bar{\psi}(\mathbf{X})Y e^{-\beta_0 A} \{A - m_0(\mathbf{X})\}$; $\sigma_0^2 = I_0^{-2} \mathbb{E}h^2(\mathbf{D}; \beta_0, \eta_0)$. Under Assumptions REG and ML1, we have

$$\sqrt{n}\sigma_0^{-1}(\hat{\beta}_{\text{ML}} - \beta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\sigma_0 I_0)^{-1} h(\mathbf{D}_i; \beta_0, \eta_0) + o_p(1),$$

which weakly converge to $N(0, 1)$.

REMARK 4.6.—Given that $\bar{\psi}(\mathbf{x}) = \text{expit}\{-r_0(\mathbf{x})\}$ and $\hat{\psi}^{[-k]}(\mathbf{x}) = \text{expit}\{-\hat{r}^{[-k]}(\mathbf{x})\}$ in our ML case, one could show that $\hat{\psi}^{[-k]}(\mathbf{x})$ achieves the same strong convergence and rate properties as $\hat{r}^{[-k]}(\cdot)$ under Assumption ML1. While generally speaking, uniform consistency of $\hat{\psi}^{[-k]}(\cdot)$ is sufficient for the desirable conclusion in Theorem 4.2 so our framework accommodates more flexible choices on $\psi(\mathbf{x})$, for example, $\psi_{\text{opt}}(\mathbf{x})$ as introduced in Section 3.3. We demonstrate this point during the proof of Theorem 4.2 in Appendix A3.

5. SIMULATION STUDY

We conduct simulation studies for our HD and ML settings separately in Sections 5.1 and 5.2, to study the point and interval estimation performance of our method.

5.1. High dimensional (parametric) setting

For the HD parametric setting, we design three data generation configurations introduced as follows to simulate different scenarios of model specification:

- i. First, generate Y following $P(Y = 1) = 1/2$. Then generate $(A, \mathbf{X}^T)^T \mid \{Y = j\} \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$ for $j = 0, 1$. Specification of $\boldsymbol{\mu}_j$ and $\boldsymbol{\Sigma}$ are presented in Appendix A5 such that $\beta_0 = 0.5$, $r_0(\mathbf{X}) = -0.22(X_1 + X_2) + 0.08(X_3 + X_4)$, and $m_0(\mathbf{X}) = -0.13(X_1 + X_2 + X_3 + X_4)$.

- ii. Generate Y following $P(Y = 1) = 1/2$ and $(A, \mathbf{X}^\top)^\top | \{Y = j\} \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$ for $j = 0, 1$. Specification of $\boldsymbol{\mu}_j$ and $\boldsymbol{\Sigma}_j$ are presented in Appendix A5 such that $\beta_0 = 0.5$, $r_0(\mathbf{X}) = -0.22(X_1 + X_2) + 0.08(X_3 + X_4) - 0.15(X_1X_2 + X_1X_3 + X_2X_3)$, and $m_0(\mathbf{X}) = -0.13(X_1 + X_2 + X_3 + X_4)$.
- iii. First generate $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ with $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$ given in Appendix A5. Then we generate A given \mathbf{X} from a Gaussian linear model with unit variance and conditional mean

$$E(A | \mathbf{X}) = 0.15(X_1 + X_2 + X_3 + X_4) + 0.075(X_1X_2 + X_1X_3 + X_2X_3).$$

Finally, generate Y by

$$\mathbb{P}(Y = 1 | A, \mathbf{X}) = \text{expit}(0.5A + 0.25X_1 + 0.25X_2 + 0.1X_3 + 0.1X_4).$$

We realise configurations (i)–(iii) with the sample size $n = 1,000, 1,500$ or $2,000$ separately and the dimension of \mathbf{X} fixed as $p = 200$. Under all these settings, we specify the nuisance models as: $r(\mathbf{X}) = \mathbf{X}^\top \boldsymbol{\gamma}$ and $m(\mathbf{X}) = \mathbf{X}^\top \boldsymbol{\alpha}$. Then both nuisance models are correctly specified under (i), only $m(\mathbf{X})$ is correctly specified under (ii), and only $r(\mathbf{X})$ is correct under (iii). Note that we cannot extract the explicit form of $m_0(\mathbf{x})$ under (iii) because A is generated conditional on \mathbf{X} without fixing $Y = 0$. But we still expect the linear model $m(\mathbf{X}) = \mathbf{X}^\top \boldsymbol{\alpha}$ to be misspecified under (iii) due to the nonlinear terms in $E(A | \mathbf{X})$. Implementing details of our HD approach are presented in Appendix A4. Specifically, all the tuning parameters in ℓ_1 -regularised regression are selected using cross-validation among the range $[0.2(\log p/n)^{1/2}, 2(\log p/n)^{1/2}]$. We conducted the regression in each setting with 300 repeated simulations.

Table 1 evaluates the performance of our estimator $\hat{\beta}_{\text{HD}}$ under all settings on its mean square error (MSE), absolute bias, and coverage probability (CP) of the 95% confidence interval (CI) estimated using bootstrap. Under all the settings, our method outputs low root-MSE and bias respectively being at most 18% and 7% of the magnitude of the true $\beta_0 (= 0.5)$ when $n = 1,000$, and at most 12% and 4% of the β_0 when $n = 2,000$. As the sample size n grows, one could see a trend of decaying on the MSEs and bias of our estimator. In addition, under all the settings, our interval estimation has proper CP locating in ± 0.03 range of the nominal level 0.95. Thus, our HD estimator performs steadily well under different model specification scenarios as long as at least one nuisance model is correctly specified.

5.2. Machine learning (nonparametric) setting

To study our proposed method under the ML setting, we let $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$ with $\Sigma_{ii} = 1$, $\Sigma_{ij} = 0.2$ for $i \neq j$, and truncate the randomly generated $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ vectors by $(-2, 2)$ to obtain the covariates \mathbf{X} . We then generate A from the Gaussian model given \mathbf{X} with unit variance and conditional mean $a_0(\mathbf{X}) = \boldsymbol{\zeta}_a^\top f_a(\mathbf{X})$ where $f_a(\mathbf{x})$ is a nonlinear basis function of \mathbf{x} including various types of effects (interaction, indicator, and trigonometric function, etc.), as defined

in Appendix A5 and ζ_0 represents its loading coefficients also given in Appendix A5. Then $A = a_0(x) + e$ where e is a standard normal variable truncated into $(-2, 2)$. Finally, we set $\beta_0 = 1$, $r_0(\mathbf{X}) = \zeta_0^\top f_r(\mathbf{X})$ with the nonlinear basis $f_r(\mathbf{X})$ (see Appendix A5) and generate Y following model (1.1). We fix $p = 20$ and set $n = 1,000$ and $2,000$ separately.

To estimate the nuisance function $r_0(x)$, we use the FMR procedure with its last step being (3.5). The number of fold for cross-fitting is set as $K = 5$. For choice of the learning algorithms \mathcal{L} , we consider four ML methods and a hybrid method of the ML estimators introduced as follows.

- a. Gradient boosted machines (GBM): an ensemble approach of classification and regressiontree (CART) using gradient boosting. Implemented by **R** package ‘gbm’ (Greenwell et al., 2020).
- b. Random forest (RF): ensemble of CART with bagging. Implemented with **R** package ‘RandomForest’ (Liaw and Wiener, 2002).
- c. Support vector machine (SVM): with linear kernel and implemented using **R** package ‘e1071’ (Dimitriadou et al., 2004).
- d. Neural network (NN): single hidden layer neural network implemented with **R** package ‘nnet’ (Ripley and Venables, 2016).
- e. Best nuisance models (Best): similar to Chernozhukov, Chetverikov, et al. (2018), for each nuisance component, we use a simple hybrid method choosing the ML estimator among (a)–(d) as the one with best prediction performance evaluated by the cross-validated sum-squared loss.

All the above-mentioned ML algorithms have been commonly used in recent years and considered in the literature of DML; for example, see Chernozhukov, Chetverikov, et al. (2018) and Cui and Tchetgen Tchetgen (2019). Tuning parameters of the ML models including the number of trees of GBM and RF, the margin of SVM, and the number of units and the weight decay of NN are selected using the resampling approach of **R** package ‘caret’ (Kuhn et al., 2020). We conducted the regression in each setting with 300 repeated simulations.

Table 2 presents the resulted average MSE, absolute bias and CP of 95% CI of $\hat{\beta}_{ML}$ obtained with the five ML modelling strategies for $n = 1,000$ and $n = 2,000$ separately. The five approaches have relatively consistent performance in terms of MSE, bias, and CP under both settings, with the variation of their MSEs smaller than 0.015. This demonstrates that performance of our framework is robust to the choice of ML algorithms. While, to a certain degree, ‘NN’ has the best performance (with the lowest bias and MSE) when $n = 1,000$ and ‘GBM’ and ‘Best’ have the best performance when $n = 2,000$. Also, interval estimations of all the approaches achieve proper coverage rates, all of which are between 0.90 and 0.95.

6. A REAL EXAMPLE: EFFECT OF EC PILL ON EARLY FETAL GESTATION

In this section, we implement our proposed HD and ML methods to study the effect of the emergency contraceptive (EC) pill on the rate of new birth and early gestation fetal death (abortion), by revisiting and exploring the data of a quasi-experimental study based on the policy reform on the EC pill in Chile (Bentancor and Clarke, 2017). In the original study, the authors collected all records of birth and fetal deaths in Chile, as well as a number of municipality-level features (education, salary, and healthcare, etc.) of women at the reproductive age (15–34), in the years around 2008, during which the country was experiencing a reform relating to the legislation of EC pills. As a consequence of this, about half of the municipalities in Chile started to provide EC pill freely in 2009, while in the remaining half, the EC pill was not available or extremely restricted in use during that period. This policy was mostly dependent on the political, economic, and public health factors characterised by a total of 16 features (denoted as Z) such as education spending, public health spending, condom use, and political conservativeness. Thus, the treatment of EC pill ($A = 1$ for EC pill accessible; $A = 0$ for EC pill not accessible) can be regarded as exogenous for the individuals.

Let $Y^{(1)}$ denote the indicator for the status of early gestation fetal death of each individual record and $Y^{(2)}$ indicate new births (pregnant and did not incur fetal death). Assume that

$$\mathbb{P}(Y^{(1)} = 1 \mid A, Z) = \text{expit}\{\beta_0^{(1)}A + r_0^{(1)}(Z)\};$$

$$\mathbb{P}(Y^{(2)} = 1 \mid A, Z) = \text{expit}\{\beta_0^{(2)}A + r_0^{(2)}(Z)\},$$

where $r_0^{(1)}(\cdot)$ and $r_0^{(2)}(\cdot)$ are two unknown functions. We are interested in inferring the two parameters $\beta_0^{(1)}$ and $\beta_0^{(2)}$ characterising the log odd ratios (log-OR) of abortion (among the pregnant individuals) and birth (among all individuals) to the treatment of the EC pill respectively. To investigate $\beta_0^{(1)}$ as the effect of the EC pill on abortion, we follow a similar strategy as Bentancor and Clarke (2017) that focuses on the individual records at the stage between 15 and 25, on which early gestation fetal death can be viewed as a reasonable proxy for illegal abortion. Note that the prevalence of $Y^{(1)}$ and $Y^{(2)}$ in their corresponding populations are both less than 5%, which could cause the logistic model to be unstable to fit. We randomly downsample the zeros in both analysis to make the prevalence of $Y^{(1)}$ and $Y^{(2)}$ 0.25 and 0.4, respectively. This procedure only changes intercepts of the logistic models and does not affect the target parameters. We find our results are not sensitive to the prevalence set for $Y^{(1)}$ and $Y^{(2)}$ as long as they are in a proper range, say 0.2 to 0.5. The resulting data set for analysing $\beta_0^{(1)}$ (abortion) has $n^{(1)} = 5,824$ samples. We take a subset with $n^{(2)} = 10,000$ samples for $\beta_0^{(2)}$ so that our algorithms will not require excessive computation time.

For our HD approach, we let X be the $p = 175$ dimensional basis joining Z , all the interaction terms of Z and the three-dimensional natural splines of all the continuous

variables in \mathbf{Z} . We specify the nuisance functions as $m(\mathbf{x}) = \text{expit}(\mathbf{x}^\top \boldsymbol{\alpha}^{(\ell)})$, and $r(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\gamma}^{(\ell)}$ for $\ell = 1, 2$. For the ML approach, we take $\mathbf{X} = \mathbf{Z}$ as the input covariates of the nuisance models. The choice and implementation of the ML algorithms are the same as in Section 5.2, except that we also use dropout⁸ to avoid over-fitting due to the fact that most covariates are at the municipality level while the records are at the individual level.

Tables 3 and 4 present the point estimation, 95% CI and (two-side) p -values of our approaches $\beta_0^{(1)}$ and $\beta_0^{(2)}$, respectively. For $\beta_0^{(1)}$, point estimations of all methods are negative and around -0.18 ± 0.04 . Their interval estimations are also internally consistent except that SVM outputs a relatively narrow CI and NN includes 0 near its CI upper bound. Correspondingly, all methods but NN reject the null ' $\beta_0^{(1)} = 0$ ' at level 0.05. We note that NN outputs a slightly worse prediction model for $A \mid \mathbf{X}, Y = 0$, which causes it to produce a relatively wider CI. The result of our hybrid method 'Best' is very consistent with HD, indicating that our methods under both settings lead to basically the same conclusion. A similar situation occurs to the estimators of $\beta_0^{(2)}$. All methods reject ' $\beta_0^{(2)} = 0$ ' at level 0.05 and their estimations are all negative values, and internally consistent on the magnitudes and CIs (SVM shows a moderate variation from other methods).

Our results reveal that distribution of the EC pill could significantly reduce the rate of illegal abortion (in the age group 15–25) and new births. This is consistent with the results of Bentancor and Clarke (2017) obtained through their municipality-level analysis. Although the estimated effect sizes are at different scales⁹ and thus incomparable between the two studies, our p -values appear to show more significance in that nearly all of them are below 0.05 while their estimated p -values are between 0.05 and 0.1. This is because we use more complex and robust nuisance models to adjust for the confounding effects of \mathbf{Z} and perform our analysis at the individual level, which enables us to have larger sample sizes.

ACKNOWLEDGEMENTS

The authors thank their adviser Tianxi Cai for helpful discussion and comments on this paper.

APPENDIX A1.: JUSTIFICATION OF THE FMR PROCEDURE

In this section, we derive error rates for the ML estimator $\hat{r}^{[k]}(\cdot)$ obtained with the FMR procedure introduced in Section 3.2. Assume that the learning algorithm \mathcal{L} attains the same strong convergence and rate properties as those for $\hat{m}^{[k]}(\cdot)$ in Assumption ML1, i.e., for each $j \in \{1, 2, \dots, K\}$ and $k \in \{1, 2, \dots, K\}$:

⁸A common and flexible technique in ML research used for regularisation and avoiding overfitting. Here we randomly and independently set each entry of the training covariates matrix as $N(0, 1)$ variable with probability 0.4 and 0.3, for the first and second study, respectively.

⁹Their effect is defined in a partially linear model of the abortion/birth rate against the treatment and control variables. While we are measuring the effect of the EC pill in a logistic model at the individual level.

$$\begin{aligned} \sup_{\mathbf{x} \in \mathcal{X}, a \in \mathcal{A}} \left| \widehat{M}^{[-k, -j]}(a, \mathbf{x}) - M_0(a, \mathbf{x}) \right| &= o_p(1); \quad \left\| \widehat{M}^{[-k, -j]}(\cdot) - M_0(\cdot) \right\|_{p,2} = o_p(n^{-1/4}); \\ \sup_{\mathbf{x} \in \mathcal{X}} \left| \widehat{a}^{[-k, -j]}(\mathbf{x}) - a_0(\mathbf{x}) \right| &= o_p(1); \quad \left\| \widehat{a}^{[-k, -j]}(\cdot) - a_0(\cdot) \right\|_{p,2} = o_p(n^{-1/4}); \\ \sup_{\mathbf{x} \in \mathcal{X}} \left| \widehat{t}^{[-k]}(\mathbf{x}) - t_0^{[-k]}(\mathbf{x}) \right| &= o_p(1); \quad \left\| \widehat{t}^{[-k]}(\cdot) - t_0^{[-k]}(\cdot) \right\|_{p,2} = o_p(n^{-1/4}). \end{aligned}$$

where $t_0^{[-k, -j]}(\mathbf{x}) = : \mathbb{E} \left[\text{logit} \left\{ \widehat{M}^{[-k, -j]}(a, \mathbf{x}) \right\} \mid \mathbf{X} = \mathbf{x}, \widehat{M}^{[-k, -j]}(\cdot) \right]$. We justify as follows that

$$\sup_{\mathbf{x} \in \mathcal{X}} \left| \widehat{r}^{[-k]}(\mathbf{x}) - r_0(\mathbf{x}) \right| = o_p(1); \quad \left\| \widehat{r}^{[-k]}(\cdot) - r_0(\cdot) \right\|_{p,2} = o_p(n^{-1/4}).$$

First, given that logit is a smooth function, it is not hard to show that

$$\begin{aligned} \sup_{\mathbf{x} \in \mathcal{X}, a \in \mathcal{A}} \left| \text{logit} \left\{ \widehat{M}^{[-k, -j]}(a, \mathbf{x}) \right\} - \text{logit} \left\{ M_0(a, \mathbf{x}) \right\} \right| &= o_p(1); \\ \left\| \text{logit} \left\{ \widehat{M}^{[-k, -j]}(\cdot) \right\} - \text{logit} \left\{ M_0(\cdot) \right\} \right\|_{p,2} &= o_p(n^{-1/4}), \end{aligned}$$

under some mild regularity conditions. Then derive the error rate of $\widehat{\beta}^{[-k]}$ as follows.

$$\begin{aligned} & |\mathcal{J}_{-k}|^{-1} \sum_{j=1}^K \sum_{i \in \mathcal{J}_{-k,j}} \text{logit} \left\{ \widehat{M}^{[-k, -j]}(A_i, \mathbf{X}_i) \right\} \{A_i - \widehat{a}^{[-k, -j]}(\mathbf{X}_i)\} \\ &= |\mathcal{J}_{-k}|^{-1} \sum_{j=1}^K \sum_{i \in \mathcal{J}_{-k,j}} (\text{logit} \{M_0(A_i, \mathbf{X}_i)\} \{A_i - a_0(\mathbf{X}_i)\} \\ &\quad + \left[\text{logit} \left\{ \widehat{M}^{[-k, -j]}(A_i, \mathbf{X}_i) \right\} - \text{logit} \{M_0(A_i, \mathbf{X}_i)\} \right] \{A_i - a_0(\mathbf{X}_i)\} \\ &\quad + \text{logit} \{M_0(A_i, \mathbf{X}_i)\} \{a_0(\mathbf{X}_i) - \widehat{a}^{[-k, -j]}(\mathbf{X}_i)\} \\ &\quad + \left[\text{logit} \left\{ \widehat{M}^{[-k, -j]}(A_i, \mathbf{X}_i) \right\} - \text{logit} \{M_0(A_i, \mathbf{X}_i)\} \right] \{a_0(\mathbf{X}_i) - \widehat{a}^{[-k, -j]}(\mathbf{X}_i)\}) \\ &= |\mathcal{J}_{-k}|^{-1} \sum_{j=1}^K \sum_{i \in \mathcal{J}_{-k,j}} \text{logit} \{M_0(A_i, \mathbf{X}_i)\} \{A_i - a_0(\mathbf{X}_i)\} \\ &\quad + \left\| \widehat{a}^{[-k, -j]}(\cdot) - a_0(\cdot) \right\|_{p,2} + \left\| \text{logit} \left\{ \widehat{M}^{[-k, -j]}(\cdot) \right\} - \text{logit} \{M_0(\cdot)\} \right\|_{p,2} + O_p(n^{-1/2}) \\ &= |\mathcal{J}_{-k}|^{-1} \sum_{j=1}^K \sum_{i \in \mathcal{J}_{-k,j}} \text{logit} \{M_0(A_i, \mathbf{X}_i)\} \{A_i - a_0(\mathbf{X}_i)\} + o_p(n^{-1/4}) + O_p(n^{-1/2}), \end{aligned}$$

under some mild regularity conditions. Similarly, we have

$$\begin{aligned} & |\mathcal{J}_{-k}|^{-1} \sum_{j=1}^K \sum_{i \in \mathcal{J}_{-k,j}} \left\{ A_i - \widehat{a}^{[-k, -j]}(\mathbf{X}_i) \right\}^2 \\ &= |\mathcal{J}_{-k}|^{-1} \sum_{j=1}^K \sum_{i \in \mathcal{J}_{-k,j}} \{A_i - a_0(\mathbf{X}_i)\}^2 + o_p(n^{-1/4}) + O_p(n^{-1/2}). \end{aligned}$$

And consequently, by Proposition 3.1 and

$$\check{\beta}^{[-k]} = \frac{\sum_{i \in \mathcal{I}_{-k}} \logit\{M_0(A_i, \mathbf{X}_i)\} \{A_i - a_0(\mathbf{X}_i)\}}{\sum_{i \in \mathcal{I}_{-k}} \{A_i - a_0(\mathbf{X}_i)\}^2} + o_p(n^{-1/4}) = \beta_0 + o_p(n^{-1/4}).$$

Then by Assumption REG that β_0 and $|a_0(\mathbf{x})|$ are bounded and recall that

$$\hat{a}^{[-k]}(\mathbf{x}) = K^{-1} \sum_{j=1}^K \hat{a}^{[-k, -j]}(\mathbf{x}),$$

the estimator $\hat{r}^{[-k]}(\cdot)$ given by equation (3.5) satisfies that:

$$\begin{aligned} & \sup_{\mathbf{x} \in \mathcal{X}} |\hat{r}^{[-k]}(\mathbf{x}) - r_0(\mathbf{x})| \\ & \leq \sup_{\mathbf{x} \in \mathcal{X}} |\hat{r}^{[-k]}(\mathbf{x}) - t_0(\mathbf{x})| + |\check{\beta}^{[-k]} - \beta_0| |a_0(\mathbf{x})| + |\check{\beta}^{[-k]}| |\hat{a}^{[-k]}(\mathbf{x}) - a_0(\mathbf{x})| \\ & \leq \sup_{\mathbf{x} \in \mathcal{X}, j} |\hat{r}^{[-k]}(\mathbf{x}) - t_j^{[-k, -j]}(\mathbf{x})| + |t_j^{[-k, -j]}(\mathbf{x}) - t_0(\mathbf{x})| + o_p(1) = o_p(1), \end{aligned}$$

where $\sup_{\mathbf{x} \in \mathcal{X}, j} |t_j^{[-k, -j]}(\mathbf{x}) - t_0(\mathbf{x})| = o_p(1)$ is a consequence of

$$\sup_{\mathbf{x} \in \mathcal{X}, a \in \mathcal{A}} \left| \logit\{\hat{M}^{[-k, -j]}(a, \mathbf{x})\} - \logit\{M_0(a, \mathbf{x})\} \right| = o_p(1).$$

And then

$$\begin{aligned} & \|\hat{r}^{[-k]}(\cdot) - r_0(\cdot)\|_{p,2} \\ & \leq \|\hat{r}^{[-k]}(\cdot) - t_0(\cdot)\|_{p,2} + |\check{\beta}^{[-k]} - \beta_0| \|a_0(\mathbf{x})\|_{p,2} + |\check{\beta}^{[-k]}| \|\hat{a}^{[-k]}(\mathbf{x}) - a_0(\mathbf{x})\|_{p,2} \\ & = \max_{j \in \{1, \dots, K\}} \|t_j^{[-k, -j]}(\cdot) - t_0(\cdot)\|_{p,2} + o_p(n^{-1/4}) \\ & = \max_{j \in \{1, \dots, K\}} \left\| \mathbb{E} \left[\logit\{\hat{M}^{[-k, -j]}(A, \mathbf{X})\} \mid \mathbf{X} = \mathbf{x} \right] - \mathbb{E} \left[\logit\{M_0(A, \mathbf{X})\} \mid \mathbf{X} = \mathbf{x} \right] \right\|_{p,2} \\ & \quad + o_p(n^{-1/4}) \\ & \leq \max_{j \in \{1, \dots, K\}} \left\| \logit\{\hat{M}^{[-k, -j]}(\cdot)\} - \logit\{M_0(\cdot)\} \right\|_{p,2} + o_p(n^{-1/4}) = o_p(n^{-1/4}). \end{aligned}$$

Thus, $\hat{r}^{[-k]}(\mathbf{x})$ satisfies Assumption ML1.

APPENDIX A2.: PROOF OF THEOREM 4.1

Proof.

By (3.2), we have

$$n^{-1} \sum_{i=1}^n h(\mathbf{D}_i; \hat{\beta}_{\text{HD}}, \hat{\eta}) = n^{-1} \sum_{i=1}^n \hat{\psi}(\mathbf{X}_i) \left\{ Y_i e^{-\hat{\beta}_{\text{HD}} A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \hat{\gamma}} \right\} \{A_i - g(\mathbf{X}_i^\top \hat{\alpha})\} = 0.$$

Our main involvement is to remove the approximation error $n^{-1} \sum_{i=1}^n h(\mathbf{D}_i; \hat{\beta}_{\text{HD}}, \hat{\eta}) - h(\mathbf{D}_i; \hat{\beta}_{\text{HD}}, \bar{\eta})$, asymptotically. Note that

$$\begin{aligned}
& n^{-1} \sum_{i=1}^n h(\mathbf{D}_i; \hat{\beta}_{\text{HD}}, \hat{\eta}) - h(\mathbf{D}_i; \hat{\beta}_{\text{HD}}, \bar{\eta}) \\
&= n^{-1} \sum_{i=1}^n \hat{\psi}(\mathbf{X}_i)(1 - Y_i) \left\{ e^{\mathbf{X}_i^\top \hat{\gamma}} - e^{\mathbf{X}_i^\top \bar{\gamma}} \right\} \{A_i - g(\mathbf{X}_i^\top \hat{\alpha})\} \\
&+ n^{-1} \sum_{i=1}^n \hat{\psi}(\mathbf{X}_i) \left\{ Y_i e^{-\hat{\beta}_{\text{HD}} A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \bar{\gamma}} \right\} \{g(\mathbf{X}_i^\top \bar{\alpha}) - g(\mathbf{X}_i^\top \hat{\alpha})\} \\
&+ n^{-1} \sum_{i=1}^n \left\{ \text{expit}(-\mathbf{X}_i^\top \hat{\gamma}) - \text{expit}(-\mathbf{X}_i^\top \gamma^*) \right\} \\
&\times \left\{ Y_i e^{-\hat{\beta}_{\text{HD}} A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \bar{\gamma}} \right\} \{A_i - g(\mathbf{X}_i^\top \bar{\alpha})\} \\
&=: \Delta_1 + \Delta_2 + \Delta_3.
\end{aligned}$$

We handle the terms Δ_1 , Δ_2 and Δ_3 separately as follows. First, we have

$$\begin{aligned}
\Delta_1 &= n^{-1} \sum_{i=1}^n (1 - Y_i) \{ \hat{\psi}(\mathbf{X}_i) - \bar{\psi}(\mathbf{X}_i) \} e^{\mathbf{X}_i^\top \bar{\gamma}} \left\{ 1 - e^{\mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma})} \right\} \{A_i - g(\mathbf{X}_i^\top \bar{\alpha})\} \\
&+ n^{-1} \sum_{i=1}^n (1 - Y_i) \hat{\psi}(\mathbf{X}_i) e^{\mathbf{X}_i^\top \bar{\gamma}} \left\{ 1 - e^{\mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma})} \right\} \{g(\mathbf{X}_i^\top \bar{\alpha}) - g(\mathbf{X}_i^\top \hat{\alpha})\} \\
&+ n^{-1} \sum_{i=1}^n (1 - Y_i) \bar{\psi}(\mathbf{X}_i) e^{\mathbf{X}_i^\top \bar{\gamma}} \left\{ 1 - e^{\mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma})} - \mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma}) \right\} \{A_i - g(\mathbf{X}_i^\top \hat{\alpha})\} \\
&+ n^{-1} \sum_{i=1}^n (1 - Y_i) \bar{\psi}(\mathbf{X}_i) e^{\mathbf{X}_i^\top \bar{\gamma}} \{g(\mathbf{X}_i^\top \bar{\alpha}) - g(\mathbf{X}_i^\top \hat{\alpha})\} \mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma}) \\
&+ n^{-1} \sum_{i=1}^n (1 - Y_i) \bar{\psi}(\mathbf{X}_i) e^{\mathbf{X}_i^\top \bar{\gamma}} \{A_i - g(\mathbf{X}_i^\top \bar{\alpha})\} \mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma}) \\
&=: \Delta_{11} + \Delta_{12} + \Delta_{13} + \Delta_{14} + \Delta_{15}
\end{aligned}$$

As $\sup_{i \in \{1, \dots, n\}} |\mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma})| = O_p(1)$ by Assumption HD4, there exists $M_1 = O(1)$ such that with probability approaching 1,

$$\left| 1 - e^{\mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma})} \right| \leq M_1 |\mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma})|, \quad \left| 1 - e^{\mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma})} - \mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma}) \right| \leq M_1 \{ \mathbf{X}_i^\top (\hat{\gamma} - \bar{\gamma}) \}^2; \tag{A1}$$

$$|\hat{\psi}(\mathbf{X}_i) - \bar{\psi}(\mathbf{X}_i)| = \frac{e^{\mathbf{X}_i^\top \gamma^*} \left| 1 - e^{\mathbf{X}_i^\top (\bar{\gamma} - \gamma^*)} \right|}{(1 + e^{\mathbf{X}_i^\top \gamma^*})(1 + e^{\mathbf{X}_i^\top \bar{\gamma}})} \leq M_1 |\mathbf{X}_i^\top (\bar{\gamma} - \gamma^*)|. \tag{A2}$$

And by Assumptions REG and HD4, there exists $M_2 = \Theta(1)$ that

$\sup_{i \in \{1, \dots, n\}} |A_i - g(\mathbf{X}_i^\top \bar{\alpha})| + |A_i - g(\mathbf{X}_i^\top \hat{\alpha})| \leq M_2$. Consequently, by Assumptions HD4 and boundness of $\psi(\cdot)$, we have

$$\begin{aligned}
|\Delta_{11}| &\leq n^{-1} \sum_{i=1}^n M_1^2 M_2 e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \|\mathbf{X}_i^\top (\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*)\| \|\mathbf{X}_i^\top (\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}})\| \\
&\leq M_1^2 M_2 \left[n^{-2} \sum_{i=1}^n e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \{\mathbf{X}_i^\top (\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*)\}^2 \sum_{i=1}^n e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \{\mathbf{X}_i^\top (\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}})\}^2 \right]^{1/2} \\
&= O_p\left(\frac{s \log p}{n}\right); \\
|\Delta_{12}| &\leq n^{-1} \sum_{i=1}^n M_1 e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \|\mathbf{X}_i^\top (\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}})\| |g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}})| \\
&\leq M_1 \left[n^{-2} \sum_{i=1}^n e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \{\mathbf{X}_i^\top (\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}})\}^2 \sum_{i=1}^n e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \{g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}})\}^2 \right]^{1/2} \\
&= O_p\left(\frac{s \log p}{n}\right); \\
|\Delta_{13}| &\leq n^{-1} \sum_{i=1}^n M_1 M_2 e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \{\mathbf{X}_i^\top (\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}})\}^2 = O_p\left(\frac{s \log p}{n}\right); \\
|\Delta_{14}| &\leq n^{-1} \sum_{i=1}^n e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \|\mathbf{X}_i^\top (\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}})\| |g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}})| = O_p\left(\frac{s \log p}{n}\right), \text{ similar to } |\Delta_{12}|.
\end{aligned}$$

By Assumptions HD2 and HD4,

$$|\Delta_{15}| \leq \left\| n^{-1} \sum_{i=1}^n (1 - Y_i) \bar{\boldsymbol{\psi}}(\mathbf{X}_i) e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \{A_i - g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}})\} \mathbf{X}_i \right\|_{\infty} \cdot \|\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}}\|_1 = O_p\left(\frac{s \log p}{n}\right).$$

Thus, we have $|\Delta_1| = O_p(s \log p/n)$, For Δ_2 , we have

$$\begin{aligned}
\Delta_2 &= n^{-1} \sum_{i=1}^n \hat{\boldsymbol{\psi}}(\mathbf{X}_i) Y_i \left(e^{-\hat{\beta}_{\text{HD}} A_i} - e^{-\beta_0 A_i} \right) \{g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}})\} \\
&\quad + n^{-1} \sum_{i=1}^n \{\hat{\boldsymbol{\psi}}(\mathbf{X}_i) - \bar{\boldsymbol{\psi}}(\mathbf{X}_i)\} \left\{ Y_i e^{-\beta_0 A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \right\} \{g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}})\} \\
&\quad + n^{-1} \sum_{i=1}^n \bar{\boldsymbol{\psi}}(\mathbf{X}_i) \left\{ Y_i e^{-\beta_0 A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \right\} \\
&\quad \quad \times \{g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}}) - g'(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) \mathbf{X}_i^\top (\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}})\} \\
&\quad + n^{-1} \sum_{i=1}^n \bar{\boldsymbol{\psi}}(\mathbf{X}_i) \left\{ Y_i e^{-\beta_0 A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \right\} g'(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) \mathbf{X}_i^\top (\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}}) \\
&=: \Delta_{21} + \Delta_{22} + \Delta_{23} + \Delta_{24}.
\end{aligned}$$

Again using Assumptions REG and HD4, there exists $M_3 = \Theta(1)$ such that

$$\left| e^{-\hat{\beta}_{\text{HD}} A_i} - e^{-\beta_0 A_i} \right| \leq M_3 |\hat{\beta}_{\text{HD}} - \beta_0|.$$

(A3)

And by Assumption HD3 and the mean value theorem, for each i , there exists t_i lying between $\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}}$ and $\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}$ such that

$$\begin{aligned}
& |g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}}) - g'(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) \mathbf{X}_i^\top (\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}})| \\
& \leq |g'(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g'(t_i)| \|\mathbf{X}_i^\top (\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}})\| \leq L \|\mathbf{X}_i^\top (\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}})\|^2.
\end{aligned} \tag{A4}$$

These combined with (A1), (A2) and Assumptions REG and HD4 lead to that

$$\begin{aligned}
|\Delta_{21}| &= O\left(\left|\hat{\beta}_{\text{HD}} - \beta_0\right| \left[n^{-1} \sum_{i=1}^n \{g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}})\}^2\right]^{1/2}\right) = O_p\left(\frac{s \log p}{n}\right); \\
|\Delta_{22}| &= O\left(\left[n^{-1} \sum_{i=1}^n \{1 + e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}}\} \{\hat{\psi}(\mathbf{X}_i) - \bar{\psi}(\mathbf{X}_i)\}^2\right.\right. \\
&\quad \left.\left. \times n^{-1} \sum_{i=1}^n \{1 + e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}}\} \{g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) - g(\mathbf{X}_i^\top \hat{\boldsymbol{\alpha}})\}^2\right]^{1/2}\right) = O_p\left(\frac{s \log p}{n}\right); \\
|\Delta_{23}| &= O\left(n^{-1} \sum_{i=1}^n \{1 + e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}}\} \|\mathbf{X}_i^\top (\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}})\|^2\right) = O_p\left(\frac{s \log p}{n}\right).
\end{aligned}$$

And by Assumptions HD2 and HD4,

$$\begin{aligned}
|\Delta_{24}| &\leq \left\| n^{-1} \sum_{i=1}^n \bar{\psi}(\mathbf{X}_i) \left\{ Y_i e^{-\beta_0 A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \right\} g'(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}}) \mathbf{X}_i \right\|_\infty \cdot \|\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}}\|, \\
&= O_p\left(\frac{s \log p}{n}\right).
\end{aligned}$$

So we also have $|\Delta_2| = O_p(s \log p/n)$. For Δ_3 , we have

$$\begin{aligned}
\Delta_3 &= n^{-1} \sum_{i=1}^n \left\{ \text{expit}(-\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}) - \text{expit}(-\mathbf{X}_i^\top \boldsymbol{\gamma}^*) \right\} Y_i \left(e^{-\hat{\beta}_{\text{HD}} A_i} - e^{-\beta_0 A_i} \right) \{A_i - g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}})\} \\
&\quad + n^{-1} \sum_{i=1}^n \left[\left\{ \text{expit}(-\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}) - \text{expit}(-\mathbf{X}_i^\top \boldsymbol{\gamma}^*) - \text{expit}'(-\mathbf{X}_i^\top \boldsymbol{\gamma}^*) \mathbf{X}_i^\top (\boldsymbol{\gamma}^* - \bar{\boldsymbol{\gamma}}) \right\} \right. \\
&\quad \left. \times \left\{ Y_i e^{-\beta_0 A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \right\} \{A_i - g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}})\} \right] \\
&\quad + n^{-1} \sum_{i=1}^n \text{expit}'(-\mathbf{X}_i^\top \boldsymbol{\gamma}^*) \mathbf{X}_i^\top (\boldsymbol{\gamma}^* - \bar{\boldsymbol{\gamma}}) \left\{ Y_i e^{-\beta_0 A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}} \right\} \{A_i - g(\mathbf{X}_i^\top \bar{\boldsymbol{\alpha}})\} \\
&= : \Delta_{31} + \Delta_{32} + \Delta_{33}.
\end{aligned}$$

Again using the mean value theorem and the fact that $|\text{expit}'(u) - \text{expit}'(v)| \leq |u - v|$ for any $u, v \in \mathbb{R}$ (by $|\text{expit}''(\cdot)| \leq 1$), we have

$$\left| \text{expit}(-\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}) - \text{expit}(-\mathbf{X}_i^\top \boldsymbol{\gamma}^*) - \text{expit}'(-\mathbf{X}_i^\top \boldsymbol{\gamma}^*) \mathbf{X}_i^\top (\boldsymbol{\gamma}^* - \bar{\boldsymbol{\gamma}}) \right| \leq \|\mathbf{X}_i^\top (\boldsymbol{\gamma}^* - \bar{\boldsymbol{\gamma}})\|^2.$$

Then by (A2), (A3), and Assumptions REG and HD4, we have

$$|\Delta_{31}| = O\left(\left|\hat{\beta}_{\text{HD}} - \beta_0\right| \left[n^{-1} \sum_{i=1}^n \{X_i^\top(\tilde{y} - \gamma^*)\}^2\right]^{1/2}\right) = O_{\mathbb{P}}\left(\frac{s \log p}{n}\right);$$

$$|\Delta_{32}| = O\left(n^{-1} \sum_{i=1}^n \{1 + e^{X_i \tilde{y}}\} \{X_i^\top(\tilde{y} - \gamma^*)\}^2\right) = O_{\mathbb{P}}\left(\frac{s \log p}{n}\right).$$

And by Assumptions HD2 and HD4,

$$|\Delta_{33}| \leq \left\| n^{-1} \sum_{i=1}^n \expit'(-X_i^\top \gamma^*) \left\{ Y_i e^{-\beta_0 A_i} - (1 - Y_i) e^{X_i^\top \tilde{y}} \right\} \{A_i - g(X_i^\top \bar{\alpha})\} X_i \right\|_{\infty} \\ \times \|\gamma^* - \tilde{\gamma}\|_1 = O_{\mathbb{P}}\left(\frac{s \log p}{n}\right).$$

Thus, we have $|\Delta_3| = O_{\mathbb{P}}(s \log p/n)$, and by Assumption HD5,

$$n^{-1} \sum_{i=1}^n h(\mathbf{D}_i; \hat{\beta}_{\text{HD}}, \hat{\eta}) - h(\mathbf{D}_i; \hat{\beta}_{\text{HD}}, \bar{\eta}) = O_{\mathbb{P}}\left(\frac{s \log p}{n}\right) = o_{\mathbb{P}}\left(\frac{1}{\sqrt{n}}\right),$$

which leads to that

$$n^{-1} \sum_{i=1}^n h(\mathbf{D}_i; \hat{\beta}_{\text{HD}}, \bar{\eta}) + o_{\mathbb{P}}\left(\frac{1}{\sqrt{n}}\right) = 0.$$

This combined with Remark 2.1 that $\mathbb{E}h(\mathbf{D}; \beta_0, \bar{\eta}) = 0$ under Assumption HD1, the regularity Assumption REG, and Theorem 5.21 of Van der Vaart (2000), comes to the conclusion of Theorem 4.1. \square

APPENDIX A3.: PROOF OF THEOREM 4.2

Following the general results of the DML estimator with nonlinear Neyman orthogonal score presented in Section 3.3 and Theorem 3.3 of Chernozhukov, Chetverikov, et al. (2018), we only need to verify their Assumptions 3.3 and 3.4 on our score function $h(\mathbf{D}; \beta, \eta)$, specifically, Assumptions A1 and A2 presented as follows.

ASSUMPTION A1 (MOMENT CONDITION WITH NEYMAN ORTHOGONALITY).

It holds: (a) $\mathbb{E}h(\mathbf{D}; \beta_0, \eta_0) = 0$ and \mathcal{B} contains an interval of length $\Theta(n^{-1/2} \log n)$ centred at β_0 ; (b) the map $(\beta, \eta) \rightarrow \mathbb{E}h(\mathbf{D}; \beta, \eta_0)$ is twice continuously Gateaux-differentiable; (c) $|\mathbb{E}h(\mathbf{D}; \beta, \eta_0)| \geq \min\{ |J_0(\beta - \beta_0)|, c_0 \}$ where the parameters $\eta_0 = \{r_0(\cdot), m_0(\cdot), \bar{\psi}\}$, $c_0 = \Theta(1)$ and $J_0 = \partial_{\beta} \mathbb{E}h(\mathbf{D}; \beta, \eta_0)|_{\beta = \beta_0} = \Theta(1)$; (d) $h(\mathbf{D}; \beta, \eta_0)$ obeys Neyman orthogonality, i.e., $\partial_{\eta} \mathbb{E}h(\mathbf{D}; \beta, \eta_0)[\eta - \eta_0] = 0$ for all $\eta \in \mathcal{E}$ where the parameter space of η : $\mathcal{E} \subseteq \{\eta: \mathbb{E}|h(\mathbf{D}; \beta_0, \eta_0)[\eta - \eta_0]| < \infty\}$.

ASSUMPTION A2 (QUALITY OF THE NUISANCE ESTIMATORS).

It holds: (a) $\hat{\eta}^{[k]}$ belongs to the realisation set \mathcal{T}_n for each $k \in \{1, 2, \dots, K\}$, with probability approaching 1 where \mathcal{T}_n satisfies $\eta_0 \in \mathcal{T}_n$ and conditions given as follows; (b) The space of β , \mathcal{B} is bounded and for each $\eta \in \mathcal{T}_n$, the functional space $\mathcal{F}_\eta = \{h(\cdot; \beta, \eta); \beta \in \mathcal{B}\}$ is measurable and its uniform covering number satisfies: there exists positive constant $R = \Theta(1)$ and $v = \Theta(1)$ such that

$$\sup_{\mathcal{Q}} \log \mathcal{N}(\epsilon \|F_\eta(\cdot)\|_{Q,2}, \mathcal{F}_\eta, \|\cdot\|_{Q,2}) \leq v \log(R/\epsilon), \quad \forall \epsilon \in (0, 1],$$

where $F_\eta(\cdot)$ is a measurable envelope function for \mathcal{F}_η ; $F_\eta(\mathbf{D}) \geq |h(\mathbf{D}; \beta, \eta)|$ for all \mathbf{D} and $\beta \in \mathcal{B}$, and there exists $q > 2$ such that $\|F_\eta(\cdot)\|_{p,q} = O(1)$; (c) there exists sequence τ_n :

$$\begin{aligned} \eta = \{r, m, \psi\} \in \mathcal{T}_n, \beta \in \mathcal{B} & \quad |\mathbb{E}h(\mathbf{D}; \beta, \eta) - \mathbb{E}h(\mathbf{D}; \beta, \{r_0, m_0, \psi\})| = o(\tau_n), \\ \eta \in \mathcal{T}_n, |\beta - \beta_0| \leq \tau_n & \quad \mathbb{E}[h(\mathbf{D}; \beta, \eta) - h(\mathbf{D}; \beta_0, \{r_0, m_0, \psi\})]^2 \\ & \quad + \mathbb{E}[h(\mathbf{D}; \beta_0, \{r_0, m_0, \psi\}) - h(\mathbf{D}; \beta_0, \eta_0)]^2 = o(1), \\ r \in (0, 1), \eta \in \mathcal{T}_n, |\beta - \beta_0| \leq \tau_n & \quad \partial_r^2 \mathbb{E}h\{\mathbf{D}; \beta_0 + r(\beta - \beta_0), \eta_0 + r(\eta - \{r_0(\cdot), m_0(\cdot), \psi\})\} = o(n^{-1/2}); \end{aligned}$$

and (d) $\mathbb{E}h^2(\mathbf{D}; \beta_0, \eta_0) = \Theta(1)$.

We simplified and adapted the original assumptions in Chernozhukov, Chetverikov, et al. (2018) to form Assumptions A1–A2, according to our own setting. The only nontrivial change made here is that in Assumption A2 (c), we require

$$\begin{aligned} \eta = \{r, m, \psi\} \in \mathcal{T}_n, \beta \in \mathcal{B} & \quad |\mathbb{E}h(\mathbf{D}; \beta, \eta) - \mathbb{E}h(\mathbf{D}; \beta, \{r_0, m_0, \psi\})| = o(\tau_n), \\ r \in (0, 1), \eta \in \mathcal{T}_n, |\beta - \beta_0| \leq \tau_n & \quad \partial_r^2 \mathbb{E}h\{\mathbf{D}; \beta_0 + r(\beta - \beta_0), \eta_0 + r(\eta - \{r_0(\cdot), m_0(\cdot), \psi\})\} = o(n^{-1/2}); \end{aligned} \tag{A5}$$

instead of:

$$\begin{aligned} \eta \in \mathcal{T}_n, \beta \in \mathcal{B} & \quad |\mathbb{E}h(\mathbf{D}; \beta, \eta) - \mathbb{E}h(\mathbf{D}; \beta, \eta_0)| = o(\tau_n), \\ r \in (0, 1), \eta \in \mathcal{T}_n, |\beta - \beta_0| \leq \tau_n & \quad |\partial_r^2 \mathbb{E}h\{\mathbf{D}; \beta_0 + r(\beta - \beta_0), \eta_0 + r(\eta - \eta_0)\}| = o(n^{-1/2}), \end{aligned}$$

as used in Assumption 3.4 (c) of Chernozhukov, Chetverikov, et al. (2018). The first inequality of (A5) is used by Chernozhukov, Chetverikov, et al. (2018) to derive a preliminary rate for the DML estimator: $|\hat{\beta}_{\text{ML}} - \beta_0| = o_p(\tau_n)$ (see their Step 1 of the proof of Lemma 6.3), and the second inequality of (A5) is used in their Step 3 the proof of Lemma 6.3 to process the second order error of

$$n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{J}_k} h(\mathbf{D}_i; \beta_0, \eta_0) - h(\mathbf{D}_i; \beta, \hat{\eta}^{[-k]})$$

uniformly for β satisfying $|\beta - \beta_0| \leq \tau_n$.

Note that our modified two assumptions are still sufficient for deriving these results. Given that $\mathbb{E}h(\mathbf{D}; \beta_0, \{r_0, m_0, \psi\}) = 0$ holds for all ψ (see Remark 2.1), there is actually no need to consider $\mathbb{E}h(\mathbf{D}; \beta, \{r_0, m_0, \psi\}) - \mathbb{E}h(\mathbf{D}; \beta, \{r_0, m_0, \psi_0\})$ when deriving $|\hat{\beta}_{\text{ML}} - \beta_0| = o_{\mathbb{P}}(\tau_n)$. While for the Step 3 of Chernozhukov, Chetverikov, et al. (2018), one can instead handle the second order error of

$$n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{J}_k} h(\mathbf{D}_i; \beta_0, \{r_0, m_0, \hat{\psi}^{[-k]}\}) - h(\mathbf{D}_i; \beta, \hat{\eta}^{[-k]}),$$

again using that $\mathbb{E}h(\mathbf{D}; \beta_0, \{r_0, m_0, \psi\}) = 0$ holds for all ψ , and then remove the remaining error:

$$n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{J}_k} h(\mathbf{D}_i; \beta_0, \{r_0, m_0, \psi_0\}) - h(\mathbf{D}_i; \beta_0, \{r_0, m_0, \hat{\psi}^{[-k]}\}),$$

through concentration based on $\partial_{\psi} \mathbb{E}h(\mathbf{D}; \beta_0, \eta_0)[\psi - \psi_0] = 0$ and the second inequality of Assumption A2 (c): $\sup_{\eta \in \mathcal{J}_n, |\eta - \beta_0| \leq \tau_n} \mathbb{E}[h(\mathbf{D}; \beta_0, \{r_0, m_0, \psi\}) - h(\mathbf{D}; \beta_0, \eta_0)]^2 = o(1)$.

Alternatively, this modification essentially reduces our requirement on the quality of $\hat{\psi}^{[-k]}(\cdot)$, as mentioned in Remark 4.6. As will be seen from the proof below, to fulfil the modified Assumption A2 (c), we do not require $\|\hat{\psi}^{[-k]}(\cdot) - \bar{\psi}(\cdot)\|_{P,2} = o_{\mathbb{P}}(\tau_n)$ as will happen when following the original version of Chernozhukov, Chetverikov, et al. (2018), but only needs $\hat{\psi}^{[-k]}(\cdot)$ to be uniformly consistent (though the former one may still be justifiable in our case because we take $\bar{\psi}(\mathbf{x}) = \text{expit}\{-\bar{r}(\mathbf{x})\}$). We now verify Assumptions A1 and A2 based on our Assumptions REG and ML1.

Proof.

Assumption A1 (a) is directly given by our logistic partial model assumption 1.1. Assumption A1 (b) is naturally satisfied as $h(\mathbf{D}; \beta, \eta_0 + r(\eta - \eta_0))$ and is a twice continuously differentiable in (β, r) . Assumption A1 (c) is directly given by Assumption REG. And Assumption A1 (d) holds by equation (2.1), combined with the model assumptions (1.1) and $\mathbb{E}[A | \mathbf{X} = \mathbf{x}, Y = 0] = m_0(\mathbf{x})$.

By Assumption ML1 and $\bar{\psi}(\mathbf{x}) = \text{expit}\{-\bar{r}(\mathbf{x})\}$, there exists $\zeta_{1,n} = o(1)$ and $\zeta_{2,n} = o(n^{-1/4})$ such that $\hat{\eta}^{[-k]}$ belongs to

$$\mathcal{T}_n = : \left\{ \eta = (r, m, \psi) : \sup_{\mathbf{x} \in \mathcal{X}} \left| \psi(\mathbf{x}) - \bar{\psi}(\mathbf{x}) \right| + |r(\mathbf{x}) - r_0(\mathbf{x})| + |m(\mathbf{x}) - m_0(\mathbf{x})| \leq \zeta_{1,n}, \right. \\ \left. \text{and } \|r(\cdot) - r_0(\cdot)\|_{p,2} + \|m(\cdot) - m_0(\cdot)\|_{p,2} \leq \zeta_{2,n} \right\},$$

with probability approaching 1, for each $k \in \{1, \dots, K\}$. We define \mathcal{T}_n of Assumption A2 in this way such that A2 (a) is satisfied. Now we validate Assumption A2 (b). By Assumption REG that A and β belong to compact sets, and $|m(\mathbf{x})| \leq m_0(\mathbf{x}) + o(1)$ is uniformly bounded for $\eta = \{r, m, \psi\} \in \mathcal{T}_n$, there exists positive $C_1 = \Theta(1)$ (1) such that for $\eta = \{r, m, \psi\} \in \mathcal{T}_n$,

$$h(\mathbf{D}; \beta, \eta) = \psi(\mathbf{X}) \left\{ Y e^{-\beta A} - (1 - Y) e^{r(\mathbf{X})} \right\} \{A - m(\mathbf{X})\} \\ \leq \left| \psi(\mathbf{X}) Y e^{-\beta A} \{A - m(\mathbf{X})\} \right| + \left| \psi(\mathbf{X}) e^{r(\mathbf{X})} (1 - Y) \{A - m(\mathbf{X})\} \right| \\ \leq C_1 \left\{ \psi(\mathbf{X}) + \psi(\mathbf{X}) e^{r(\mathbf{X})} \right\} = C_1 \left\{ \text{expit}(-\psi(\mathbf{X})) + \text{expit}(\psi(\mathbf{X})) \right\} \\ \leq C_1 + 1; \\ \partial_\beta h(\mathbf{D}; \beta, \eta)_{p,2}^2 = \mathbb{E} \left[\psi(\mathbf{X}) Y e^{-\beta A} A \{A - m(\mathbf{X})\} \right]^2 \leq C_1.$$

Then by Example 19.7 of Van der Vaart (2000), Assumption A2 (b) holds with $v = 1$ and R being the diameter of \mathcal{B} . Note that Assumption A2 (d) is again directly given by Assumption REG. It remains to verify Assumption A2 (c). For each $\eta = \{r, m, \psi\} \in \mathcal{T}_n$ and $\beta \in \mathcal{B}$, using the boundness of β , A , $m_0(\mathbf{x})\psi(\mathbf{x})$ and $\psi(\mathbf{x})e^{r(\mathbf{x})}$, there exists $C_2 = \Theta(1)$ such that

$$\left| \mathbb{E}h(\mathbf{D}; \beta, \eta) - \mathbb{E}h(\mathbf{D}; \beta, \{r_0, m_0, \psi\}) \right| \\ \leq \left| \mathbb{E} \psi(\mathbf{X}) Y e^{-\beta A} \{m_0(\mathbf{X}) - m(\mathbf{X})\} \right| + \left| \mathbb{E} \psi(\mathbf{X}) (1 - Y) e^{r(\mathbf{X})} \{m_0(\mathbf{X}) - m(\mathbf{X})\} \right| \\ + \left| \mathbb{E} \psi(\mathbf{X}) e^{r(\mathbf{X})} (1 - Y) \{1 - e^{r_0(\mathbf{X})} - r(\mathbf{X})\} \{A - m_0(\mathbf{X})\} \right| \\ \leq C_2 \left(\|m_0(\mathbf{X}) - m(\mathbf{X})\|_{p,2} + \|r_0(\mathbf{X}) - r(\mathbf{X})\|_{p,2} + \|r_0(\mathbf{X}) - r(\mathbf{X})\|_{p,2}^2 \right) \leq 3C_2 \zeta_{2,n}.$$

So we take $\tau_n = n^{-1/4}$ and by $\zeta_{2,n} = o(n^{-1/4})$, the first inequality of Assumption A2 (c) is satisfied. Again by the boundness of β , A , $m_0(\mathbf{x})\psi(\mathbf{x})$ and $\psi(\mathbf{x})e^{r(\mathbf{x})}$; and $\zeta_{1,n}$, $\tau_n = o(1)$, there exists $C_3 = \Theta(1)$ such that

$$\mathbb{E} [h(\mathbf{D}; \beta, \eta) - h(\mathbf{D}; \beta_0, \{r_0, m_0, \psi\})]^2 + \mathbb{E} [h(\mathbf{D}; \beta_0, \{r_0, m_0, \psi\}) - h(\mathbf{D}; \beta_0, \eta_0)]^2 \\ \leq \mathbb{E} [h(\mathbf{D}; \beta, \eta) - h(\mathbf{D}; \beta_0, \eta)]^2 + \mathbb{E} [h(\mathbf{D}; \beta_0, \eta) - h(\mathbf{D}; \beta_0, \{r_0, m_0, \psi\})]^2 \\ + \mathbb{E} [h(\mathbf{D}; \beta_0, \{r_0, m_0, \psi\}) - h(\mathbf{D}; \beta_0, \eta_0)]^2 \\ \leq \mathbb{E} \left[\psi(\mathbf{X}) e^{-\beta_0 A} \{e^{(\beta_0 - \beta)A} - 1\} \{A - m(\mathbf{X})\} \right]^2 \\ + \mathbb{E} \left[\psi(\mathbf{X}) e^{-\beta_0 A} \{m_0(\mathbf{X}) - m(\mathbf{X})\} \right]^2 + \mathbb{E} \left[\psi(\mathbf{X}) e^{r(\mathbf{X})} \{m_0(\mathbf{X}) - m(\mathbf{X})\} \right]^2 \\ + \mathbb{E} \left[\psi(\mathbf{X}) e^{r(\mathbf{X})} (1 - Y) \{1 - e^{r_0(\mathbf{X})} - r(\mathbf{X})\} \{A - m_0(\mathbf{X})\} \right]^2 \\ + \mathbb{E} \left[\left| \bar{\psi}(\mathbf{X}) - \psi(\mathbf{X}) \right| \left\{ e^{-\beta_0 A} + e^{r_0(\mathbf{X})} \right\} |A - m_0(\mathbf{X})| \right]^2 \\ \leq C_3 \sup_{a \in \mathcal{A}} \left| e^{(\beta_0 - \beta)a} - 1 \right| + C_3 \sup_{\mathbf{x} \in \mathcal{X}} \left\{ \left| (1 + e^{r(\mathbf{x})}) [\bar{\psi}(\mathbf{x}) - \psi(\mathbf{x})] \right| \right. \\ \left. + |m(\mathbf{x}) - m_0(\mathbf{x})| + |e^{r(\mathbf{x})} - r_0(\mathbf{x}) - 1| \right\} \\ \leq C_3 \sup_{a \in \mathcal{A}} \left| e^{(\beta_0 - \beta)a} - 1 \right| + C_3 \sup_{\mathbf{x} \in \mathcal{X}} \left\{ \left| \bar{\psi}(\mathbf{x}) - \psi(\mathbf{x}) \right| \right. \\ \left. + |m(\mathbf{x}) - m_0(\mathbf{x})| + \{ \text{expit}(r_0(\mathbf{x})) + 1 \} |e^{r(\mathbf{x})} - r_0(\mathbf{x}) - 1| \right\} \\ \leq C_3 \left\{ (e^{\tau_n C_3} - 1) + 2\zeta_{1,n} + 2|e^{\zeta_{1,n}} - 1| \right\} = o(1),$$

which validates the second inequality of Assumption A2 (c). At last, for each $r \in (0, 1)$, denote by $\beta^* = \beta_0 + r(\beta - \beta_0)$, $\eta^* = \{r^*, m^*, \psi\} = \{r_0(\cdot), m_0(\cdot), \psi\} + r(\eta - \{r_0(\cdot), m_0(\cdot), \psi\})$. Similar as the above deduction, we have that there exists $C_4 = \Theta(1)$,

$$\begin{aligned} & \partial_r^2 \mathbb{E} h\{\mathbf{D}; \beta_0 + r(\beta - \beta_0), \eta_0 + r(\eta - \{r_0(\cdot), m_0(\cdot), \psi\})\} \\ &= \mathbb{E} \psi(\mathbf{X}) Y e^{-\beta^* A} A (\beta - \beta_0) \{m_0(\mathbf{X}) - m(\mathbf{X})\} \\ & \quad + \mathbb{E} \psi(\mathbf{X}) (1 - Y) e^{r^* \mathbf{X}} \{r_0(\mathbf{X}) - r(\mathbf{X})\} \{m_0(\mathbf{X}) - m(\mathbf{X})\} \\ & \leq C_4 |\beta - \beta_0| \cdot \mathbb{E} |m_0(\mathbf{X}) - m(\mathbf{X})| + C_4 \mathbb{E} |r_0(\mathbf{X}) - r(\mathbf{X})| \cdot \mathbb{E} |m_0(\mathbf{X}) - m(\mathbf{X})| \\ & = O(\|m(\cdot) - m_0(\cdot)\|_{p,2}^2) + O\left\{(\beta - \beta_0)^2\right\} + O(\|r(\cdot) - r_0(\cdot)\|_{p,2}^2) \\ & = O(\xi_{2,n}^2) + o(\tau_n^2) = o(n^{-1/4}). \end{aligned}$$

□

Using the verified Assumptions A1 and A2, one can follow nearly the same proof procedures as those of Theorem 3.3 and Lemma 6.3 in Chernozhukov, Chetverikov, et al. (2018) to prove our Theorem 4.2. The only minor difference concerning the processing of $\psi^{[-k]}$ has been presented as above. As we point out, one can handle this smoothly by first considering $n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{J}_k} h(\mathbf{D}_i; \beta_0, \{r_0, m_0, \hat{\psi}^{[-k]}\})$ when deriving the as initial rate and asymptotic expansion of $\hat{\beta}_{\text{ML}}$ as $\mathbb{E}[h(\mathbf{D}_i; \beta_0, \{r_0, m_0, \hat{\psi}^{[-k]}\}) | \hat{\psi}^{[-k]}] = 0$, and finally concentrate $n^{-1} \sum_{k=1}^K \sum_{i \in \mathcal{J}_k} h(\mathbf{D}_i; \beta_0, \{r_0, m_0, \psi_0\}) - h(\mathbf{D}_i; \beta_0, \{r_0, m_0, \hat{\psi}^{[-k]}\})$ using that $\partial_\psi \mathbb{E} h(\mathbf{D}; \beta_0, \eta_0) [\psi - \psi_0] = 0$ and $\hat{\psi}^{[-k]}(\cdot)$ is uniformly consistent.

APPENDIX A4.: NUMERICAL IMPLEMENTATION OF THE HD APPROACH

We present and demonstrate the implementation procedure of our HD approach mentioned in Remark 3.3 that uses LASSO instead of the Dantzig equation and modifies the construction procedures to make it solvable using the **R** packages *glmnet* and *RCAL*. Let $G(u) = \int g(u) du$, and recall that $\tilde{\gamma}$ is some initial estimator obtained through ℓ_1 -regularised logistic regression for Y versus $\{A, \mathbf{X}\}$ and $\hat{\psi}(x) = \text{expit}(-\mathbf{x}^\top \tilde{\gamma})$. Then we fit

$$\min_{\alpha \in \mathbb{R}^p} n^{-1} \sum_{i=1}^n (1 - Y_i) \hat{\psi}(\mathbf{X}_i) e^{\mathbf{X}_i^\top \tilde{\gamma}} \{-A_i \mathbf{X}_i^\top \alpha + G(\mathbf{X}_i^\top \alpha)\} + \lambda_\alpha \|\alpha\|_1, \quad (\text{A6})$$

to obtain $\hat{\alpha}$. It is not hard to see that the KKT (or subgradient) condition of (A6) is equivalent to the ℓ_∞ -constraint in (3.1). When the link function of $g(\cdot)$ is identity (linear model) or $\text{expit}(\cdot)$ (logistic model), can be solved using the **R** package *glmnet* with proper specification of the sample weights, i.e., $(1 - Y_i) \hat{\psi}(\mathbf{X}_i) e^{\mathbf{X}_i^\top \tilde{\gamma}}$. Then we solve

$$n^{-1} \sum_{i=1}^n \hat{\psi}(\mathbf{X}_i) \left\{ Y_i e^{-\beta A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \tilde{\gamma}} \right\} \{A_i - g(\mathbf{X}_i^\top \hat{\alpha})\} = 0,$$

to obtain a preliminary estimator $\tilde{\beta}$. It can be shown that when either $r(\mathbf{x})$ or $m(\mathbf{x})$ is correctly specified, the estimator $\tilde{\beta}$ should approach β_0 at the rate $O_p\left\{\left(\frac{\log p}{n}\right)^{1/2}\right\}$, i.e., the ℓ_2 errors of $\tilde{\gamma}$ and $\hat{\alpha}$. So $\tilde{\beta}$ provides a good enough approximation of β_0 that can be used for the ℓ_1 -regularised (weighted) calibration regression (Tan, 2020b) to estimate γ :

$$\min_{\beta \in \mathbb{R}, \eta \in \mathbb{R}^{pn}} n^{-1} \sum_{i=1}^n \hat{\psi}(\mathbf{X}_i) e^{\mathbf{X}_i^\top \tilde{\gamma} g'(\mathbf{X}_i^\top \hat{\alpha})} \left\{ Y_i e^{-\tilde{\beta} A_i} - \mathbf{X}_i^\top \gamma + (1 - Y_i)(\tilde{\beta} A_i + \mathbf{X}_i^\top \gamma) \right\} + \lambda_\gamma \|\gamma\|_1. \quad (\text{A7})$$

Similarly, the KKT condition of (A7) corresponds to the ℓ_∞ -constraints in (3.2), though they are not always imposing the same moment conditions: when the nuisance model $r(\mathbf{x})$ is misspecified, $\tilde{\gamma}$ and $\hat{\gamma}$ typically have different limits. We use **R** package *RCAL* to solve (A7) with the response taken as Y_i , regressors as \mathbf{X}_i , sample weight $\hat{\psi}(\mathbf{X}_i) e^{\mathbf{X}_i^\top \tilde{\gamma} g'(\mathbf{X}_i^\top \hat{\alpha})}$ and offset $\tilde{\beta} A_i$ for each i . Denoting the solution of (A7) as $\hat{\gamma}$, we finally obtain the estimator $\hat{\beta}_{\text{HD}}$ by solving

$$n^{-1} \sum_{i=1}^n \hat{\psi}(\mathbf{X}_i) e^{\mathbf{X}_i^\top (\tilde{\gamma} - \hat{\gamma})} \left\{ Y_i e^{-\beta A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \hat{\gamma}} \right\} \{A_i - g(\mathbf{X}_i^\top \hat{\alpha})\} = 0. \quad (\text{A8})$$

Here the final estimating equation is asymptotically equivalent to the second row of (3.2) only when $r(\mathbf{x})$ is correctly specified $\tilde{\gamma}$ and $\hat{\gamma}$ have the same limiting values). When $r(\mathbf{x})$ is misspecified, the orthogonal score function used in (A8), denoted by $h'(\mathbf{D}, \beta_0, \eta)$, is not the same as the $h(\mathbf{D}, \beta_0, \eta)$ used in the main text. We will point out that this does not hurt the Neyman orthogonality of $h'(\mathbf{D}, \beta_0, \eta)$. It is because when $r(\mathbf{x})$ is misspecified but $m(\mathbf{x})$ is correct (by our model Assumption HD1, at least one nuisance model need to be correct), $\partial_r h'(\mathbf{D}; \beta_0, \tilde{\eta})[r - \bar{r}]$ is naturally satisfied due to the correctness of $m(\mathbf{x})$, and $\partial_m h'(\mathbf{D}; \beta_0, \tilde{\eta})[m - \bar{m}]$ is satisfied according to the KKT (moment) condition of (A7). When $m(\mathbf{x})$ is misspecified but $r(\mathbf{x})$ is correct, $\partial_m h'(\mathbf{D}; \beta_0, \tilde{\eta})[m - \bar{m}]$ is naturally satisfied and (A8) is asymptotically equivalent with

$$n^{-1} \sum_{i=1}^n \hat{\psi}(\mathbf{X}_i) \left\{ Y_i e^{-\beta A_i} - (1 - Y_i) e^{\mathbf{X}_i^\top \hat{\gamma}} \right\} \{A_i - g(\mathbf{X}_i^\top \hat{\alpha})\} = 0,$$

as $\tilde{\gamma}$ and $\hat{\gamma}$ approach the true γ_0 and the second order errors is asymptotically negligible. So $\partial_r h'(\mathbf{D}; \beta_0, \tilde{\eta})[r - \bar{r}]$ is satisfied by (A7). Thus, our modified construction procedure does not break the theoretical guarantee of $\hat{\beta}_{\text{HD}}$.

APPENDIX A5.: ADDITIONAL DETAILS OF NUMERICAL EXPERIMENTS

First, we present the mean vector and covariance matrix used to generate A and X in Section 5.1:

- i. Take $\boldsymbol{\mu}_1 = (0.4, -0.25, -0.25, 0, \dots, 0)$, $\boldsymbol{\mu}_0 = \mathbf{0}$, and

$$(\boldsymbol{\Sigma}^{-1})_{ij} = \begin{cases} 1.5 & i = j = 1 \\ 1.2 & i = j \geq 2 \\ 0.2 & i = 1, 2 \leq j \leq 5 \text{ or } 2 \leq i \leq 5, j = 1 \\ 0 & \text{else} \end{cases}$$

- ii. Take $\boldsymbol{\mu}_1 = (0.4, -0.25, -0.25, 0, \dots, 0)$, $\boldsymbol{\mu}_0 = \mathbf{0}$,

$$(\boldsymbol{\Sigma}^{-1})_{ij} = \begin{cases} 1.5 & i = j = 1 \\ 1.2 & i = j \geq 2 \\ 0.2 & i = 1, 2 \leq j \leq 5 \text{ or } 2 \leq i \leq 5, j = 1 \\ 0 & \text{else} \end{cases}$$

and

$$(\boldsymbol{\Sigma}_0^{-1})_{ij} = \begin{cases} 1.5 & i = j = 1 \\ 1.2 & i = j \geq 2 \\ 0.2 & i = 1, 2 \leq j \leq 5 \text{ or } 2 \leq i \leq 5, j = 1 \\ 0.075 & i = 3, 4, j = 2 \text{ or } i = 2, j = 3, 4 \text{ or } i = 3, j = 4 \text{ or } i = 4, j = 3 \\ 0 & \text{else} \end{cases}$$

- iii. Take the covariance of \mathbf{X} as

$$\boldsymbol{\Sigma}_{ij} = \begin{cases} 0.5 & i = j \\ 0.15 & i \leq 4, j \leq 4, i \neq j \\ 0 & \text{else} \end{cases}$$

Then we present the specific choice on the basis functions for data generation in Section 5.2. In specific, we take

$$f_a(\mathbf{x}) = \left\{ \frac{1}{1 + e^{x_1}}, \frac{1}{1 + e^{x_2}}, \sin(x_3), \cos(x_4), I(x_5 > 0), I(x_6 > 0), x_7 x_8, x_9 x_{10} \right\}^T;$$

$$\boldsymbol{\zeta}_a = (1, -1, 0.5, 0.5, 0.25, -0.25, 0.1, 0.1)^T,$$

for $a_0(\mathbf{x})$, the conditional mean of A given $\mathbf{X} = \mathbf{x}$. And that

$$f_r(\mathbf{x}) = \left\{ x_1 x_2 x_3, x_4 x_5, x_6^3, \sin^2(x_7), \cos(x_8), \frac{1}{1 + x_9^2}, \frac{1}{1 + e^{x_{10}}}, I(x_{11} > 0), I(x_{12} > 0) \right\}^T;$$

$$\boldsymbol{\zeta}_r = (0.1, 0.1, 0.1, -0.5, 0.5, 1, -1, 0.25, -0.25)^T,$$

to specify $r_0(\mathbf{x})$.

REFERENCES

- Athey S. and Imbens GW (2017). The state of applied econometrics: Causality and policy evaluation. *Journal of Economic Perspectives* 31(2), 3–32. [PubMed: 29465214]
- Belloni A, Chernozhukov V, Chetverikov D, Hansen C. and Kato K. (2018). High-dimensional econometrics and regularized GMM. Paper 1806.01888, [arXiv.org](https://arxiv.org/abs/1806.01888).
- Bentancor A. and Clarke D. (2017). Assessing plan B: The effect of the morning after pill on children and women. *Economic Journal* 127(607), 2525–52.
- Bickel PJ., Ritov Y and Tsybakov AB (2009). Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics* 37(4), 1705–32.
- Bradic J, Wager S. and Zhu Y. (2019). Sparsity double robust inference of average treatment effects. Paper 1905.00744, [arXiv.org](https://arxiv.org/abs/1905.00744).
- Bühlmann P. and van de Geer S. (2011). *Statistics for High-Dimensional Data: Methods, Theory and Applications*, Berlin: Springer Science & Business Media.
- Candes E. and Tao T. (2007). The Dantzig selector: Statistical estimation when p is much larger than n . *Annals of Statistics* 35(6), 2313–51.
- Chen HY (2007). A semiparametric odds ratio model for measuring association. *Biometrics* 63(2), 413–21. [PubMed: 17688494]
- Chernozhukov V, Chetverikov D, Demirer M, Duflo E, Hansen C, Newey W. and Robins J. (2018). Double/debiased machine learning for treatment and structural parameters. *Econometrics Journal* 21(1), C1–C68.
- Chernozhukov V, Newey W. and Robins J. (2018). Double/debiased machine learning using regularized Riesz representers. CEMMAP Working paper 15/18, Institute of Fiscal Studies, University College London, UK.
- Chernozhukov V, Newey W, Robins J. and Singh R. (2018). Double/debiased machine learning of global and local parameters using regularized Riesz representers. Paper 1802.08667, [arXiv.org](https://arxiv.org/abs/1802.08667).
- Colangelo K. and Lee Y-Y (2020). Double debiased machine learning nonparametric inference with continuous treatments. Paper 2004.03036, [arXiv.org](https://arxiv.org/abs/2004.03036).
- Cui Y. and Tchetgen Tchetgen E. (2019). Bias-aware model selection for machine learning of doubly robust functionals. Paper 1911.02029, [arXiv.org](https://arxiv.org/abs/1911.02029).
- Dimitriadou E, Hornik K, Leisch F, Meyer D. and Weingessel A. (2004). R package e1071: Misc functions of the Department of Statistics. Software package, Technical University of Vienna, Austria.
- Dukes O. and Vansteelandt S. (2020). Inference for treatment effect parameters in potentially misspecified high-dimensional models. *Biometrika*, 108, 321–34.
- Farrell MH, Liang T. and Misra S. (2018). Deep neural networks for estimation and inference: Application to causal effects and other semiparametric estimands. Paper 1809.09953, [arXiv.org](https://arxiv.org/abs/1809.09953).
- Friedman J, Hastie T. and Tibshirani R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* 33(1), 1–22. [PubMed: 20808728]
- Ghosh Sand Tan Z(2020). Doublyrobustsemiparametricinferenceusingregularizedcalibratedestimation with high-dimensional data. Paper 2009.12033, [arXiv.org](https://arxiv.org/abs/2009.12033).
- Giné E. and Nickl R. (2016). *Mathematical Foundations of Infinite-Dimensional Statistical Models*, vol. 40. New York: Cambridge University Press.
- Greenwell B, Boehmke B. and Cunningham J, GBM Developers(2020). Package ‘gbm’, R package version, vol. 2.
- Janková J. and van de Geer S. (2016). Confidence regions for high-dimensional generalized linear model under sparsity. Paper 1610.01353, [arXiv.org](https://arxiv.org/abs/1610.01353).
- Knaus MC (2018). A double machine learning approach to estimate the effects of musical practice on student’s skills. Paper 1805.10300, [arXiv.org](https://arxiv.org/abs/1805.10300).
- Knaus MC (2020). Double machine learning based program evaluation under unconfoundedness. Paper 2003.03191, [arXiv.org](https://arxiv.org/abs/2003.03191).

- Kuchibhotla AK and Chakraborty A. (2018). Moving beyond sub-Gaussianity in high-dimensional statistics: Applications in covariance estimation and linear regression. Paper 1804.02605, [arXiv.org](https://arxiv.org/abs/1804.02605).
- Kuhn M, Wing J, Weston S, Williams A, Keefer C, Engelhardt A, Cooper T, Mayer Z. and Kenkel B, R Core Team, et al. (2020). Package caret, CRAN.
- Liaw A. and Wiener M. (2002). Classification and regression by RandomForest. R News 2(3), 18–22.
- Lin X. and Carroll RJ (2006). Semiparametric estimation in general repeated measures problems. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 68(1), 69–88.
- Liu M, Xia Y, Cho K. and Cai T. (2020). Integrative high dimensional multiple testing with heterogeneity under data sharing constraints. Paper 2004.00816, [arXiv.org](https://arxiv.org/abs/2004.00816).
- Ma R, Cai TT and Li H. (2020). Global and simultaneous hypothesis testing for high-dimensional logistic regression models. Journal of the American Statistical Association, 116(534), 984–98. [PubMed: 34421157]
- Negahban SN, Ravikumar P, Wainwright MJ and Yu B. (2012). A unified framework for high-dimensional analysis of M -estimators with decomposable regularizers. Statistical Science 27(4), 538–57.
- Nekipelov D, Semenova V. and Syrgkanis V (2018). Regularized orthogonal machine learning for nonlinear semiparametric models. Paper 1806.04823, [arXiv.org](https://arxiv.org/abs/1806.04823).
- Ning Y, Sida P. and Imai K. (2020). Robust estimation of causal effects via a high-dimensional covariate balancing propensity score. Biometrika 107(3), 533–54.
- Oprescu M, Syrgkanis V. and Wu ZS (2019). Orthogonal random forest for causal inference. Proceedings of Machine Learning Research 97, 4932–41.
- Ripley B. and Venables W. (2016). Package nnet. R package version 7.3–12, CRAN.
- Semenova V. and Chernozhukov V. Debiased machine learning of conditional average treatment effects and other causal functions. Econometrics Journal, Published ahead of print, 29 August 2020. 10.1093/ectj/utaa027.
- Severini TA and Staniswalis JG (1994). Quasi-likelihood estimation in semiparametric models. Journal of the American statistical Association 89(426), 501–11.
- Smucler E, Rotnitzky A. and Robins JM (2019). A unifying approach for doubly-robust ℓ_1 -regularized estimation of causal contrasts. Paper 1904.03737, [arXiv.org](https://arxiv.org/abs/1904.03737).
- Tan Z. (2019a). On doubly robust estimation for logistic partially linear models. Statistics & Probability Letters 155, 108577.
- Tan Z. (2019b). RCAL: Regularized calibrated estimation. R package version 1.0, CRAN.
- Tan Z. (2020a). Model-assisted inference for treatment effects using regularized calibrated estimation with high-dimensional data. Annals of Statistics 48(2), 811–37.
- Tan Z. (2020b). Regularized calibrated estimation of propensity scores with model misspecification and high-dimensional data. Biometrika 107(1), 137–58.
- Tchetgen Tchetgen EJ, Robins JM and Rotnitzky A. (2010). On doubly robust estimation in a semiparametric odds ratio model. Biometrika 97(1), 171–80. [PubMed: 23049119]
- Van de Geer S, Bühlmann P, Ritov Y. and Dezeure R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. Annals of Statistics 42(3), 1166–202.
- Van der Vaart AW (2000). Asymptotic statistics, vol. 3. Cambridge: Cambridge University Press.
- Wager S. and Athey S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. Journal of the American Statistical Association 113(523), 1228–42.
- Xia L, Li B. Nanand Y (2020). A revisit to debiased Lasso for generalized linear models. Paper 2006.12778, [arXiv.org](https://arxiv.org/abs/2006.12778).
- Yang J-C., Chuang H-C and Kuan C-M (2020). Double machine learning with gradient boosting and its application to the big n audit quality effect. Annals issue in honor of George Tiao: Statistical learning for dependent data. Journal of Econometrics 216(1), 268–83.
- Zhu Y. and Bradic J. (2018). Significance testing in non-sparse high-dimensional linear models. Electronic Journal of Statistics 12(2), 3312–64.
- Zimmert M. and Lechner M. (2019). Nonparametric estimation of causal heterogeneity under high-dimensional confounding. Paper 1908.08779, [arXiv.org](https://arxiv.org/abs/1908.08779).

Table 1.

Average mean square error (MSE), average absolute bias (Bias), and average coverage probability (CP) of 95% CI of our HD estimator with the sample size set as 1,000, 1,500 and 2,000, under configurations (i)–(iii) described in Section 5.1. Number of repetition for each setting is 300.

	Configuration (i)			Configuration (ii)			Configuration (iii)		
n	1,000	1,500	2,000	1,000	1,500	2,000	1,000	1,500	2,000
MSE	0.008	0.006	0.004	0.007	0.005	0.004	0.007	0.004	0.003
Bias	0.021	0.021	0.018	0.020	0.022	0.018	0.030	0.016	0.013
CP	0.91	0.92	0.96	0.94	0.92	0.94	0.94	0.93	0.96

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 2.

Average mean square error (MSE), average absolute bias (Bias), and average coverage probability (CP) of 95% CI of our ML estimator with sample sizes set as 1,000 and 2,000, which the nuisance models estimated using the four ML algorithms as well as the ‘Best’ approach described in Section 5.2. The number of repetitions for each setting is 300.

	<i>n</i> = 1,000					<i>n</i> = 2,000				
	GBM	RF	SVM	NN	Best	GBM	RF	SVM	NN	Best
MSE	0.013	0.014	0.013	0.012	0.013	0.007	0.008	0.008	0.007	0.007
Bias	0.037	0.046	0.048	0.015	0.037	0.039	0.051	0.049	0.042	0.039
CP	0.92	0.95	0.93	0.94	0.93	0.90	0.92	0.90	0.91	0.91

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 3.

Point estimations, 95% CI lower/upper bounds (LB/UB) and two-side p -values for $\beta_0^{(1)}$ (log odds ratio of early gestation fetal death to the treatment of the EC pill) of our HD and ML (with the five different realisations described in Section 5.2) approaches.

Method	HD	GBM	SVM	RF	NN	Best
β_0	-0.171	-0.220	-0.153	-0.215	-0.177	-0.190
CI LB	-0.310	-0.355	-0.279	-0.378	-0.391	-0.324
CI UB	-0.028	-0.085	-0.027	-0.052	0.038	-0.057
p -value	0.018	0.000	0.016	0.007	0.088	0.003

Table 4.

Point estimations, 95% CI lower/upper bounds (LB/UB) and two-side p -values for $\beta_0^{(2)}$ (log odds ratio of new birth to the treatment of the EC pill) of our HD and ML (with the five different realisation described in Section 5.2) approaches.

Method	HD	GBM	SVM	RF	NN	Best
β_0	-0.186	-0.135	-0.124	-0.112	-0.148	-0.125
CI LB	-0.287	-0.235	-0.217	-0.224	-0.259	-0.222
CI UB	-0.089	-0.035	-0.030	0.000	-0.036	-0.029
p -value	0.000	0.007	0.009	0.033	0.010	0.009