# Nonparametric augmented probability weighting with sparsity

Xin He, Xiaojun Mao, and Zhonglei Wang [*][†]

September 29, 2022

## Abstract

Nonresponse frequently arises in practice, and simply ignoring it may lead to erroneous inference. Besides, the number of collected covariates may increase as the sample size in modern statistics, so parametric imputation or propensity score weighting usually leads to inefficiency without consideration of sparsity. In this paper, we propose a nonparametric imputation method with sparse learning by employing an efficient kernel-based learning gradient algorithm to identify truly informative covariates. Moreover, an augmented probability weighting framework is adopted to improve the estimation efficiency of the nonparametric imputation method and establish the limiting distribution of the corresponding estimator under regularity assumptions. The performance of the proposed method is also supported by several simulated examples and one real-life analysis.

*Keywords:* Central limit theorem, Reproducing kernel Hilbert space, Nonresponse, Sparse learning.

[*]X. He is with the School of Statistics and Management, Shanghai University of Finance and Economics, Shanghai 200433, China (e-mail:he.xin17@mail.shufe.edu.cn); X. Mao is with the School of Mathematical Sciences, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: maoxj@sjtu.edu.cn). Z. Wang is with Wang Yanan Institute for Studies in Economics and School of Economics, Xiamen University, Xiamen, Fujian 361005, China (e-mail: wangzl@xmu.edu.cn).

[†]X. He and X. Mao contributed equally to this work. Corresponding authors: Z. Wang.

# 1 Introduction

Nonresponse is a common problem in social science and other related fields, and simply ignoring it may lead to inefficiency or even erroneous inference due to confounding covariates (Rosenbaum and Rubin; 1983; Qu et al.; 2010; Abadie and Imbens; 2016; Lin et al.; 2018). Moreover, the number of collected covariates is relatively large in modern statistics, which makes learning nonresponse even more challenging (Yang et al.; 2020). How to deal with the nonresponse under the high-dimensional setup is still an open question.

Sparse learning bridges the gap between the high-dimensional data analysis and nonresponse. It is generally believed that among the numerous covariates, only a few of them contribute to the response of interest, known as truly informative ones, while others are noise. Thus, a variety of sparse learning methods have been proposed to identify those truly informative covariates under regularity assumptions. The linear response model assumption is popularly imposed, and various attempts have been made by designing sparse-induced regularization (Tibshirani; 1996; Fan and Li; 2001; Zou; 2006; Shen et al.; 2012, 2013), evaluating the marginal dependence (Fan and Lv; 2008; Wang and Leng; 2016), or checking variable robustness against added noise (Barber and Candes; 2019). Extended methods have also been developed for nonparametric models (Lin and Zhang; 2006; Huang et al.; 2010; Fan et al.; 2011). However, all these methods require explicit model assumptions that are difficult to validate in practice or suffer from heavy computational burden. To circumvent this difficulty, Belloni et al. (2013) proposed a feasible lasso method, which is similar to the adaptive lasso (Zou; 2006), for variable selection in a partial linear model. By using machine learning algorithms to handle high-dimensional nuisance parameters, Chernozhukov et al. (2018) proposed a double/debiased machine learning procedure to achieve parametric convergence rate for a low dimensional parameter. A valid double robust estimator using lasso-type penalty is discussed by Tan (2020). Recently, kernel-based sparse learning methods have been inspired

by the fact that the gradient functions provide an appropriate criterion to identify a general dependence structure in a model-free fashion. Specifically, Rosasco et al. (2013) proposed a novel learning-gradient method, which adds an empirical functional penalty on the gradients to a standard kernel ridge regression in a reproducing kernel Hilbert space (RKHS). Besides, Yang et al. (2016) employed pair-wise learning to estimate the gradient functions and considered a functional group lasso penalty to induce sparsity. It is worth pointing out that the lack of selection consistency (Rosasco et al.; 2013) and the high computational cost (Yang et al.; 2016) remain unsolved. To alleviate the difficulties, He et al. (2021) proposed a two-step sparse learning framework, which is computationally efficient in the sense that it only requires to fit the standard kernel ridge regression and the selection consistency is established under regularity assumptions. The method proposed by He et al. (2021) can be regarded as a nonparametric joint screening approach and achieves methodological flexibility, numerical efficiency and asymptotic consistency simultaneously.

Propensity score weighting is commonly used to handle nonresponse (Robins et al.; 1994; Wooldridge; 2007; Tan; 2010; Graham et al.; 2012; Zhao et al.; 2017), but conventional methods using all covariates may lead to numerical failure, including the lack of convergence and inefficiency, due to overfitting. Thus, sparse assumption is often imposed to estimate the propensity scores more efficiently (Shevade and Keerthi; 2003; Genkin et al.; 2007). A Bayesian variable selection method has been proposed by Chen et al. (1999) for logistic regression; also see Wainwright et al. (2007), Banerjee et al. (2008) and Ravikumar et al. (2010) for details about penalized logistic regression models. The group lasso (Yuan and Lin; 2006) was generalized to logistic regression model by Kim et al. (2006), and Meier et al. (2008) proposed a more efficient group lasso algorithm than that of Kim et al. (2006). Besides, Meier et al. (2008) also established the asymptotic consistency of the corresponding estimator. Ning et al. (2020) proposed a high-dimensional covariate balancing propensity score estimator, and they validated that their proposed estimator is of parametric convergence rate

and is asymptotically normal distributed. See Tang et al. (2014) and Bertsimas and King (2017) for a review of identifying informative covariates for logistic regression models.

In this paper, we propose a nonparametric Augmented Inverse Probability Weighting (AIPW) framework (Robins et al.; 1994) to handle the nonresponse under the assumption of sparsity. Inspired by the key observation that the gradient functions provide an appropriate information of the truly informative covariates in a model-free fashion, we employ a kernel-based sparse learning algorithm to efficiently impute the nonresponse. It only requires to fit the standard kernel ridge regression, which has an analytical solution, and the gradient functions can be directly computed by the derivative reproducing property. More importantly, the truly informative covariates can be exactly recovered with high probability. Even though the nonparametric imputation with sparse learning achieves consistency, its convergence rate is at most $O_{\mathbb{P}}(n^{-1/6}\log(n))$ under regularity assumptions, so it is hard to construct an interval estimator. To alleviate this difficulty, an AIPW framework is adopted to improve the convergence rate of the corresponding estimator, and a central limit theorem can be established. To achieve this goal, certain propensity score methods for analyzing sparse data suffice under regularity conditions; see Section 2.2 for details. The corresponding variance estimator is also discussed. The superior performance of the proposed nonparametric AIPW framework is also supported by the numerical comparisons against some state-of-the-art methods in several simulated examples and one real-life analysis.

The rest of this paper is organized as follows. Section 2 provides the background and introduces the proposed nonparametric AIPW framework. The theoretical properties of the corresponding estimator are established under regularity assumptions in Section 3. Section 4 reports the numerical experiments on the simulated and real-life examples. A brief summary is provided in Section 5.

# 2 Method

Consider

$$Y = f^*(\boldsymbol{x}) + \epsilon, \tag{1}$$

where $f^*(\boldsymbol{x}) = \mathbb{E}(Y \mid \boldsymbol{X} = \boldsymbol{x})$ is a continuous function of a covariate vector $\boldsymbol{x} = (x_1, ..., x_p)^\top$ taking values from a $p$ dimensional separable and compact metric space $\mathcal{X} \subset \mathbb{R}^p$, and $\epsilon$ denotes a random noise with conditional mean zero and bounded variance. We are interested in inferring $\theta^* = \mathbb{E}(Y)$ from a random sample $\{(\boldsymbol{x}_i, y_i) : i = 1, \ldots, n\}$ generated by (1).

If $y_1, \ldots, y_n$ were fully observed, the sample mean $\widehat{\theta}_n = n^{-1} \sum_{i=1}^n y_i$ would be an efficient estimator of $\theta^*$. However, it is generally not the case in practice, and the response of interest suffers from nonresponse. For $i = 1, \ldots, n$, denote $\delta_i$ to be the response indicator of $y_i$, where $\delta_i = 1$ if $y_i$ is observed and 0 otherwise. For simplicity, we assume missing at random (Rubin; 1976) for the response mechanism,

$$\Pr(\delta_i = 1 \mid \boldsymbol{x}_i, y_i) = \Pr(\delta_i = 1 \mid \boldsymbol{x}_i), \tag{2}$$

and denote $\pi^*(\boldsymbol{x}) = \Pr(\delta = 1 \mid \boldsymbol{x})$.

If consistent estimators $\widehat{f}_0(\boldsymbol{x})$ and $\widehat{\pi}(\boldsymbol{x})$ for $f^*(\boldsymbol{x})$ and $\pi^*(\boldsymbol{x})$ are available, then an AIPW estimator,

$$\widehat{\theta}_{AIPW} = \frac{1}{n} \sum_{i=1}^n \left[ \widehat{f}_0(\boldsymbol{x}_i) + \frac{\delta_i}{\widehat{\pi}(\boldsymbol{x}_i)} \{y_i - \widehat{f}_0(\boldsymbol{x}_i)\} \right], \tag{3}$$

can be applied to estimate $\theta^*$. More rigorously, the estimator (3) is not an AIPW estimator unless we replace $n^{-1}$ by $(\sum_{i=1}^n \delta_i \widehat{\pi}_i^{-1})^{-1}$. However, if the response model is correctly specified, we can show that $n^{-1} \sum_{i=1}^n \delta_i \widehat{\pi}_i^{-1} \to 1$ in probability under regularity conditions. When $p$ is small and $f^*(\boldsymbol{x})$ is a parametric model, standard statistical methods can be used

to obtain $\widehat{f}_0(\boldsymbol{x})$ and $\widehat{\pi}(\boldsymbol{x})$. As $p$ increases, however, it is not reasonable to include all covariates to estimate $f^*(\boldsymbol{x})$ and $\pi^*(\boldsymbol{x})$ due to the curse of dimensionality. Moreover, the model misspecification for $f^*(\boldsymbol{x})$ leads to erroneous inference.

## 2.1 Estimation of $f^*(x)$ via nonparametric sparse learning

To overcome those difficulties, we employ an efficient kernel-based sparse learning algorithm (He et al.; 2021) to estimate $f^*(\boldsymbol{x})$ in (1). Denote $\mathcal{H}_K$ to be an RKHS induced by a pre-specified kernel function $K(\cdot, \cdot)$, where $K(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is bounded, symmetric and positive semi-definite. It can be shown that $\mathcal{H}_K$ associated with the kernel $K(\cdot, \cdot)$ is the completion of the linear space spanned by $\{K_{\boldsymbol{x}}(\cdot) : \boldsymbol{x} \in \mathcal{X}\}$ with an inner product $\langle K_{\boldsymbol{x}}, K_{\boldsymbol{u}} \rangle_K = K(\boldsymbol{x}, \boldsymbol{u})$ for $\boldsymbol{x}, \boldsymbol{u} \in \mathcal{X}$, where $K_{\boldsymbol{x}}(\cdot) = K(\boldsymbol{x}, \cdot)$. Thus, $\mathcal{H}_K$ is uniquely determined by a kernel function $K(\cdot, \cdot)$ and the reproducing property, $\langle f, K_{\boldsymbol{x}} \rangle_K = f(\boldsymbol{x})$ for $f \in \mathcal{H}_K$ and $\boldsymbol{x} \in \mathcal{X}$. It is noteworthy that the RKHS induced by some universal kernel, such as the Gaussian kernel, is fairly large in the sense that any continuous function can be well approximated (Steinwart; 2005). Thus, by (2), we assume

$$f^*(\boldsymbol{x}) = \underset{f \in \mathcal{H}_K}{\operatorname{argmin}} \mathbb{E}[\delta\{y - f(\boldsymbol{x})\}]^2. \tag{4}$$

By (2), (4) implies that $f^*(\boldsymbol{x}) = \mathbb{E}(Y \mid \boldsymbol{X} = \boldsymbol{x}) \in \mathcal{H}_K$.

A covariate $x_l$ is non-informative for $f^*(\boldsymbol{x})$, if and only if $g_l^*(\boldsymbol{x}) = 0$ for $\boldsymbol{x} \in \mathcal{X}$ almost surely, where $x_l$ is the $l$th component of $\boldsymbol{x}$ for $l = 1, \ldots, p$, and $g_l^*(\boldsymbol{x}) = \partial f^*(\boldsymbol{x})/\partial x_l$. Thus, the usefulness of $x_l$ for estimating $f^*(\boldsymbol{x})$ can be measured via the $\mathcal{L}^2$-norm of $g_l^*(\boldsymbol{x})$, $\|g_l^*\|_2^2 = \int_{\mathcal{X}} \{g_l^*(\boldsymbol{x})\}^2 d\rho(\boldsymbol{x})$, where $\rho(\boldsymbol{x})$ is the marginal distribution of $\boldsymbol{X}$. Denote $\mathcal{A}^* = \{l : \|g_l^*\|_2^2 > 0\}$ to be the active set containing all informative covariates associated with $f^*(\boldsymbol{x})$. To estimate $g_l^*(\boldsymbol{x})$ efficiently, we consider the following derivative reproducing property (Zhou;

2007):

$$g_l^*(\boldsymbol{x}) = \langle f^*, \partial_l K_{\boldsymbol{x}} \rangle_K, \tag{5}$$

where $\partial_l K_{\boldsymbol{x}}(\cdot) = \partial K(\boldsymbol{x}, \cdot)/\partial x_l$. Specifically, by (4)–(5), once an initial estimator of $f^*(\boldsymbol{x})$ is available, say $\widehat{f}(\boldsymbol{x})$, its gradient function $g_l^*(\boldsymbol{x})$ can be estimated by (5). Without loss of generality, we assume that the first $m$ samples are fully observed, and the subsequent analysis is conditional on the realised sample, where $m = \sum_{i=1}^n \delta_i$. Under (2) and the high-dimensional setup, to obtain an initial estimator $\widehat{f}(\boldsymbol{x})$, we employ the standard kernel ridge regression by solving

$$\widehat{f} = \underset{f \in \mathcal{H}_K}{\arg\min} \Big[ \frac{1}{m} \sum_{i=1}^m \{y_i - f(\boldsymbol{x}_i)\}^2 + \lambda \|f\|_K^2 \Big], \tag{6}$$

where $\lambda > 0$ is a tuning parameter controlling the model complexity and typically goes to 0 as $m$ goes to infinity, and $\|\cdot\|_K^2$ denotes the RKHS-norm induced by the inner product $\langle \cdot, \cdot \rangle_K$; see Section 3 for details. By the representer theorem (Mercer; 1909), the solution of (6) is of the form

$$\widehat{f}(\boldsymbol{x}) = \sum_{i=1}^m \widehat{\alpha}_i K(\boldsymbol{x}_i, \boldsymbol{x}) = \widehat{\boldsymbol{\alpha}}^\top \boldsymbol{K}_m(\boldsymbol{x}), \tag{7}$$

where $\boldsymbol{K}_m(\boldsymbol{x}) = (K(\boldsymbol{x}_1, \boldsymbol{x}), ..., K(\boldsymbol{x}_m, \boldsymbol{x}))^\top$, and $\widehat{\boldsymbol{\alpha}} = (\widehat{\alpha}_1, ..., \widehat{\alpha}_m)^\top \subset \mathbb{R}^m$ are the estimated representer coefficients. That is, the representer theorem converts the original optimization task (6) in an infinite functional space $\mathcal{H}_K$ into a $m$-dimensional vector space. By (7), the optimization task (6) is equivalent to

$$\widehat{\boldsymbol{\alpha}} = \underset{\boldsymbol{\alpha}}{\text{argmin}} \Big[ \frac{1}{m} \sum_{i=1}^m \big\{y_i - \boldsymbol{\alpha}^\top \boldsymbol{K}_m(\boldsymbol{x}_i)\big\}^2 + \lambda \boldsymbol{\alpha}^T \boldsymbol{K} \boldsymbol{\alpha} \Big],$$

and its solution is $\widehat{\boldsymbol{\alpha}} = (\boldsymbol{K} + \lambda \boldsymbol{I})^{-1}\boldsymbol{y}$, where $\boldsymbol{K}$ is an $m \times m$ matrix with $(i, j)$th entry being $K(\boldsymbol{x}_i, \boldsymbol{x}_j)$ and $\boldsymbol{y} = (y_1, \ldots, y_m)^\top \in \mathbb{R}^m$.

Once $\widehat{\boldsymbol{\alpha}}$ is obtained, the gradient function in (5) can be estimated by

$$\widehat{g}_l(\boldsymbol{x}) = \frac{\partial \widehat{f}(\boldsymbol{x})}{\partial x_l} = \widehat{\boldsymbol{\alpha}}^\top \partial_l \boldsymbol{K}_m(\boldsymbol{x}), \quad (l = 1, ..., p),$$

where $\partial_l \boldsymbol{K}_m(\boldsymbol{x}) = (\partial_l K_{\boldsymbol{x}_1}(\boldsymbol{x}), \ldots, \partial_l K_{\boldsymbol{x}_m}(\boldsymbol{x}))^\top$. Since the marginal distribution $\rho(\boldsymbol{x})$ is seldom available, instead of the $\mathcal{L}^2$-norm, the empirical norm $\|\cdot\|_m$ is considered:

$$\|\widehat{g}_l\|_m^2 = \frac{1}{m} \sum_{i=1}^{m} \left\{ \widehat{g}_l(\boldsymbol{x}_i) \right\}^2 = \frac{1}{m} \sum_{i=1}^{m} \left\{ \widehat{\boldsymbol{\alpha}}^\top \partial_l \boldsymbol{K}_m(\boldsymbol{x}_i) \right\}^2,$$

and the estimated active set is $\widehat{\mathcal{A}}_{v_m} = \left\{ l : \|\widehat{g}_l\|_m^2 > v_m \right\}$, where $v_m$ is a thresholding value determined through a stability-based selection criterion (Sun et al.; 2013). Finally, we refit (6) with the selected covariates in $\widehat{\mathcal{A}}_{v_n}$ to obtain the nonparametric estimator $\widehat{f}_0(\boldsymbol{x})$.

It is worthy pointing out that the employed sparse learning algorithm was originally proposed by He et al. (2021), and they only focused on the purpose of variable selection and established the selection consistency without considering nonresponse. Yet, we generalized their method to handle incomplete samples in this paper and treat it as an valid intermediate estimator of our proposed estimator. More importantly, we further established a central limit theorem for the proposed nonparamtric estimator, which is rare and attractive in machine learning, and a variance estimation is also provided as well; see Section 3 for details.

## 2.2 Estimation of $\pi^*(x)$

There exist works to estimate $\pi^*(\boldsymbol{x})$ under the assumption of sparsity, and we consider the group lasso (Meier et al.; 2008) as an example by assuming

$$\text{logit}\{\pi^*(\boldsymbol{x}_i)\} = \beta_0^* + \boldsymbol{x}_i^\top \boldsymbol{\beta}_1^*,$$

where $\text{logit}(z) = \log(z) - \log(1 - z)$ for $z \in (0, 1)$. In addition, assume that the covariate vector can be rewritten as $\boldsymbol{x} = (\boldsymbol{x}_1^\top, \ldots, \boldsymbol{x}_G^\top)^\top$, where $\boldsymbol{x}_g \in \mathbb{R}^{\text{df}_g}$ contains the covariates of the $g$th group for $g = 1, \ldots, G$, and $\text{df}_g$ is the corresponding degrees of freedom. For example, $\text{df}_g = 3$ if $\boldsymbol{x}_g$ corresponds to a categorical covariate with four levels, and $\text{df}_g = 1$ if $x_g$ is continuous; see Meier et al. (2008) for details.

The log-likelihood estimator with a group lasso penalty is obtained by solving

$$\widehat{\boldsymbol{\beta}}_{\lambda_2} = \underset{\boldsymbol{\beta}}{\text{argmin}} \left\{ -l(\boldsymbol{\beta}) + \lambda_2 p(\boldsymbol{\beta}) \right\}, \tag{8}$$

where $l(\boldsymbol{\beta}) = \sum_{i=1}^n [\delta_i \log\{\pi(\boldsymbol{x}_i)\} + (1 - \delta_i) \log\{1 - \pi(\boldsymbol{x}_i)\}]$ is the log-likelihood of the response indicators, $\boldsymbol{\beta}^\top = (\beta_0, \boldsymbol{\beta}_1^\top) \in \mathbb{R}^{p+1}$ with $\beta_0 \in \mathbb{R}$ and $\boldsymbol{\beta}_1 \in \mathbb{R}^p$, $p(\boldsymbol{\beta}) = \sum_{g=1}^G \text{df}_g^{1/2} \|\boldsymbol{\beta}_g\|_2$ is the group lasso penalty, $\|\cdot\|_2$ is the Euclidean norm, and $\boldsymbol{\beta}_g$ corresponds to $\boldsymbol{x}_{i,g}$ for $g = 1, \ldots, G$. The block co-ordinate gradient descent algorithm is used to obtain $\widehat{\boldsymbol{\beta}}_{\lambda_2}$ in (8), and the detailed algorithm is adjourned to Appendix A.

**Remark 1.** *Since the estimated response probability is used to improve the convergence rate of the estimator in (3), the response model is assumed to be correctly specified; see Qin et al. (2017) for a similar assumption. In addition to the group lasso method (Meier et al.; 2008), other penalized logistic regression estimators (Fan et al.; 2014; Ning et al.; 2020) can be used to estimate the response probability. However, to guarantee the asymptotic central limit theorem in Theorem 2, the estimated response probability by other methods should satisfy*

*Lemma 3; see Section 3 for details.*

# 3   Theoretical Properties

In this section, we investigate the asymptotic consistency of $\widehat{f}_0(\boldsymbol{x})$ and establish the central limit theorem for the AIPW estimator in (3) under regularity assumptions.

Denote an integral operator $L_K : \mathcal{L}^2(\mathcal{X}, \rho) \to \mathcal{L}^2(\mathcal{X}, \rho)$ as $L_K(f)(\boldsymbol{x}) = \int_{\mathcal{X}} K(\boldsymbol{x}, \boldsymbol{u}) f(\boldsymbol{u}) d\rho(\boldsymbol{u})$, for $f \in \mathcal{L}^2(\mathcal{X}, \rho)$, where $\mathcal{L}^2(\mathcal{X}, \rho) = \{f : \int f^2(\boldsymbol{x}) d\rho(\boldsymbol{x}) < \infty\}$. If the RKHS $\mathcal{H}_K$ is separable, then by the spectral theorem (Fischer and Steinwart; 2020), we have $L_K f = \sum_{j \geq 1} \mu_j \langle f, e_j \rangle_2 e_j$, where $\{e_j : j = 1, 2, \ldots\}$ form an orthonormal basis of $\mathcal{L}^2(\mathcal{X}, \rho)$, $\{\mu_j : j = 1, 2, \ldots\}$ are the corresponding eigenvalues with respect to $L_K$, and $\langle f, g \rangle_2 = \int_{\mathcal{X}} f(\boldsymbol{x}) g(\boldsymbol{x}) d\rho(\boldsymbol{x})$ denotes the inner product of $f(\boldsymbol{x})$ and $g(\boldsymbol{x})$ in $\mathcal{L}^2(\mathcal{X}, \rho)$. By Mercer's theorem (Steinwart and Christmann; 2008), under regularity assumptions, the eigen-expansion of $K(\boldsymbol{x}, \boldsymbol{u})$ is $K(\boldsymbol{x}, \boldsymbol{u}) = \sum_{j \geq 1} \mu_j e_j(\boldsymbol{x}) e_j(\boldsymbol{u})$. Hence, the RKHS-norm of any $f \in \mathcal{H}_K$ can also be expressed as

$$\|f\|_K^2 = \sum_{j \geq 1} \frac{\langle f, e_j \rangle_2^2}{\mu_j}.$$

The above result implies the decay rate of $\mu_j$ fully characterizes the complexity of the RKHS.

The following technical assumptions are made to investigate the theoretical properties of the AIPW in (3).

**Assumption 1**: There exists a positive constant $r \in (1/2, 1]$, such that $f^*(\boldsymbol{x})$ is in the range of the $r$th power of $L_K$, denoted as $L_K^r$. Besides, the distribution of $\epsilon$ has a $q$-exponential tail with some function $q(\cdot)$; that is, there exists a constant $c_1 > 0$, such that $\Pr(|\epsilon| > t) \leq c_1 \exp\{-q(t)\}$ for any $t > 0$.

**Assumption 2**: There exist positive values $\kappa_{1,p}$ and $\kappa_{2,p}$, which may depend on $p$, such that $\sup_{\boldsymbol{x} \in \mathcal{X}} \|K_{\boldsymbol{x}}\|_K \leq \kappa_{1,p}$ and $\sup_{\boldsymbol{x} \in \mathcal{X}} \|\partial_l K_{\boldsymbol{x}}\|_K \leq \kappa_{2,p}$ for $l = 1, ..., p$.

**Assumption 3**: There exists a positive constant $\xi_1 > 1$ such that

$$\min_{l \in \mathcal{A}^*} \|g_l^*\|_2^2 >$$
$$c_m \max \left\{ \kappa_{1,p} \|f^*\|_K, q^{-1} \left( \log \frac{4c_1 m}{\delta_m} \right) \right\} m^{-\frac{2r-1}{2(2r+1)}} (\log p)^{\xi_1},$$

where $c_m$ is provided in Lemma 1.

**Assumption 4**: There exists $\kappa \in (0, 1/2)$ such that $\kappa < \pi(\boldsymbol{x}) < 1 - \kappa$ for all $\boldsymbol{x} \in \mathcal{X}$.

**Assumption 5**: $\mathbb{E}(\boldsymbol{X}\boldsymbol{X}^\top)$ is invertible, and its smallest eigenvalue is bounded away from zero by a fixed positive constant $c_{min}$, and recall that $\boldsymbol{X}$ is the random vector associated with $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n$.

**Assumption 6**: Let $\boldsymbol{X}_g$ be the random vector associated with the $g$th group, and we normalize $\boldsymbol{X}_g$ such that $\mathbb{E}(\boldsymbol{X}_g^\top \boldsymbol{X}_g)$ is a $\mathrm{df}_g \times \mathrm{df}_g$ identity matrix. Then, there exists $L_n$ such that $\max_{\boldsymbol{x}} \max_g (\boldsymbol{x}_g^\top \boldsymbol{x}_g) \leq nL_n^2$, where $\boldsymbol{x}_g$ corresponds to the normalized $\boldsymbol{X}_g$.

**Assumption 7**: $\max_{g=1,\ldots,G} \mathrm{df}_g = O(1)$, there exists a constant number $\zeta > 0$ such that $\log(G) = o(n^{1/3-2\zeta})$ and $G \gg \log(n)$, $N_0 = O(1)$, $\lambda_2 \asymp \log(G)$, i.e. $\lambda_2$ is of the order $\log(G)$, and $L_n^2 = O\{1/\log(G)\}$, where $N_0$ is the number of non-zero group effects.

Assumptions 1–3 are proposed for the kernel-based sparse learning algorithm, and Assumptions 4–7 are required by the group lasso logistic regression. In Assumption 1, the integral operator $L_K$ is self-adjoint and semi-positive definite, so its fractional operator $L_K^r$ is well-defined, and its range is contained in $\mathcal{H}_K$ as long as $r \geq 1/2$; see Smale and Zhou (2007) and Mendelson and Neeman (2010) for details. This implies that for some function $h \in \mathcal{L}^2(\mathcal{X}, \rho)$, it holds $L_K^r f^* = \sum_{j \geq 1} \mu_j^r \langle h, e_j \rangle_2 e_j \in \mathcal{H}_K$, so ensures strong estimation consistency under the RKHS-norm. The second part of Assumption 1 characterizes the tail behavior of the error distribution, and it relaxes the commonly-used bounded response assumption in the machine learning literature (Smale and Zhou; 2007; Rosasco et al.; 2013; Lv et al.; 2018). Besides, the assumption on the error distribution is general and can be

satisfied by a variety of distributions (Wang and Leng; 2016; Zhang et al.; 2016). For example, if the error distribution is sub-Gaussian or bounded, then $q(t) = O(t^2)$ suffices; if the distribution of $\epsilon$ is sub-exponential, $q(t) = O(\min\{t/C, t^2/C^2\})$ suffices with $C > 0$. Assumption 2 bounds the RKHS-norms associated with the kernel function and its gradient functions, and it is satisfied by popular kernels, including the Gaussian kernel, linear kernel and the Sobolev kernel (Smale and Zhou; 2007; Rosasco et al.; 2013; Yang et al.; 2016). For example, $\kappa_{1,p} = \kappa_{2,p} = 1$ for the Gaussian kernel, $K(\boldsymbol{x}, \boldsymbol{u}) = \exp\{-\|\boldsymbol{x} - \boldsymbol{u}\|_2^2/(2\sigma^2)\}$, and $\kappa_{1,p} = Cp$ and $\kappa_{2,p} = C$ for the linear kernel, $K(\boldsymbol{x}, \boldsymbol{u}) = \boldsymbol{x}^\top \boldsymbol{u}$, for some positive constant $C$. Assumption 3 requires that the gradient functions contain sufficient information about the truly informative covariates. It is worthy pointing out that we measure the significance of each gradient function to distinguish informative and uninformative covariates without any explicit model specification. The minimal signal strength in Assumption 3 is much tighter than those in other nonparametric sparse learning methods (Huang et al.; 2010; Yang et al.; 2016), which often require the signal strength to be bounded below by some positive constant. Assumption 4 bounds the response probability, and it is commonly used to avoid inefficient estimators. To obtain the desired convergence rate, Assumption 4 guarantees that $m \asymp n$ in probability, where $a_n \asymp b_n$ is equivalent to $a_n = O(b_n)$ and $b_n = O(a_n)$. The smallest eigenvalue of $\mathbb{E}(\boldsymbol{X}\boldsymbol{X}^\top)$ is bounded by a fixed positive constant in Assumption 5, and it is a special case of assumption (b) of Meier et al. (2008). In Assumption 6, the convergence rate of $L_n$ is related with that of the estimated response probability. Assumption 7 is used to guarantee that

$$\mathbb{E}\{|\eta_{\widehat{\boldsymbol{\beta}}_{\lambda_2}}(\boldsymbol{X}) - \eta_{\boldsymbol{\beta}^*}(\boldsymbol{X})|^2\} = O_{\mathbb{P}}(n^{-2/3-2\zeta}), \tag{9}$$

where $\eta_{\boldsymbol{\beta}}(\boldsymbol{x}) = \beta_0 + \boldsymbol{x}^\top \boldsymbol{\beta}_1$, $\boldsymbol{\beta}^* = (\beta_0^*, \boldsymbol{\beta}_1^*)$ and the expectation is taken with respect to $\boldsymbol{X}$ conditional on $\widehat{\boldsymbol{\beta}}_{\lambda_2}$ or the observations. Specifically, the value $\zeta$ is used to show $n^{-\zeta/2} \log(n) \to 0$ as $n \to 0$, so it can be chosen arbitrarily small; see the proof of Theorem 2 for details.

**Lemma 1.** *Suppose Assumptions 1–2 are satisfied, and $\lambda = m^{-1/(2r+1)}$. Then, for any $\delta_m \in (0,1)$, with probability at least $1 - \delta_m$, there holds*

$$\max_{1 \leq l \leq p} \left| \|\widehat{g}_l\|_m^2 - \|g_l^*\|_2^2 \right| \leq c_m c_{p,q} \log\left(\frac{8p}{\delta_m}\right) m^{-\Theta},$$

*where $c_m$ is a constant depending only on $\kappa_{1,p}, \kappa_{2,p}$ and $\|f^*\|_K^2$, $c_{p,q} = \max\left\{ \kappa_{1,p} \|f^*\|_K, q^{-1}\left(\log \frac{4c_1 m}{\delta_m}\right) \right\}$ with $q^{-1}(\cdot)$ denoting the inverse function of $q(\cdot)$, and $\Theta = \frac{2r-1}{2(2r+1)}$.*

Lemma 1 guarantees that $\|\widehat{g}_l\|_m^2$ converges to $\|g_l^*\|_2^2$ with high probability, and it is crucial to establish the selection consistency of the employed sparse learning algorithm. The convergence result in Lemma 1 still holds even when the dimension diverges with the sample size, and the quantities $\|f^*\|_K^2$ and $\|L_K^{-r} f^*\|_2$, which may depend on the number of truly informative covariates of $f^*(\boldsymbol{x})$, may also diverge as the sample size increases. For instance, if $f^*(\boldsymbol{x}) = \boldsymbol{x}^T \boldsymbol{\beta}^*$, then $\|f^*\|_K^2 = \|\boldsymbol{\beta}^*\|_2^2$, which clearly depends on the number of truly informative covariates. However, such dependency is difficult to quantify explicitly in a fully general case (Fukumizu and Leng; 2014).

The following lemma establishes the asymptotic selection consistency of the proposed sparse learning method.

**Lemma 2.** *Suppose that the assumptions of Lemma 1 and Assumption 3 are satisfied. If $v_m = 0.5 c_m c_{p,q} m^{-\Theta} (\log p)^{\xi_1}$, then $\Pr\left(\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\right) \to 1$, as $m \to \infty$.*

Lemma 2 shows that the selected covariates can exactly recover the truly informative ones with probability tending to 1. This result is particularly general in that it is established without any model specification. The following theorem shows that $\widehat{f}_0(\boldsymbol{x})$ achieves a fast convergence rate in term of the infinity norm, where $\widehat{f}_0(\boldsymbol{x})$ is obtained by the standard kernel ridge regression (6) based on the selected covariates in $\widehat{\mathcal{A}}_{v_n}$.

**Theorem 1.** *Suppose the assumptions of Lemma 2 are satisfied and denote the probability* $\Pr(\widehat{\mathcal{A}}_{v_n} \neq \mathcal{A}^*) = \Delta_m$. *If* $\lambda = m^{-\frac{1}{2r+1}}$, *then with probability at least* $1 - \delta_m - \Delta_m$, *there holds*

$$\|\widehat{f}_0 - f^*\|_K \;\; \leq \;\; c_{m,2} c_{p_0,q} \log\left(\frac{4}{\delta_m}\right) m^{-\Theta},$$

*where* $p_0 = |\mathcal{A}^*|$,

$$\begin{aligned}
c_{m,2} \;\; &= \;\; 4 \max\{\kappa_{2,p_0}^2, \kappa_{2,p_0}^2 \|f^*\|_K, \|f^*\|_K^2\} \\
&\quad \times \max\{3\kappa_{1,p_0}, 2\sqrt{2}\kappa_{2,p_0}^2, \|L_K^{-r} f^*\|_2\},
\end{aligned}$$

*and* $c_{p_0,q} = \max\left\{\kappa_{1,p_0}\|f^*\|_K, q^{-1}\left(\log\frac{4c_1 m}{\delta_m}\right)\right\}$.

*Additionally, if we take* $r = 1$, *and assume that* $p_0 = O(1)$, $\epsilon$ *has sub-Gaussian or subexponential tail and by Assumption 4, we have* $\left\|\widehat{f}_0 - f^*\right\|_\infty = O_\mathbb{P}(n^{-\frac{1}{6}} \log n)$.

Theorem 1 establishes the uniform convergence rate of the refitted estimator $\widehat{f}_0(\boldsymbol{x})$, and it plays a crucial role to establish the central limit theory of the AIPW estimator in (3). The required tail behavior of $\epsilon$ in Theorem 1 is to quantify $q^{-1}(\cdot)$ explicitly for simplicity, and it can be extended to any error distribution satisfying Assumption 1.

**Lemma 3.** *Suppose Assumptions 4–7 are satisfied. Then, given* $\widehat{\boldsymbol{\beta}}_{\lambda_2}$, *there holds* $\mathbb{E}[|\widehat{\pi}(\boldsymbol{X}) - \pi^*(\boldsymbol{X})|^2] = O_\mathbb{P}(n^{-2/3-2\zeta})$.

Lemma 3 establishes the convergence rate of the estimated response probability using group lasso logistic regression. By Lemma 3, we essentially require that the estimated response probability should be at least consistent. A similar requirement is also discussed by Tan (2020). Specifically, Tan (2020) assumed a correctly specified response model in order to achieve valid interval estimator. If other penalized logistic regression estimators are considered, Assumptions 4–7 should be replaced in order to guarantee Lemma 3.

By Theorem 1 and Lemma 3, we can validate the following central limit theorem for the AIPW estimator in (3).

**Theorem 2.** *Suppose all the assumptions in Theorem 1 and Lemma 3 are satisfied. If* $\mathbb{E}\{|f^*(\boldsymbol{X})+\delta\pi^*(\boldsymbol{X})^{-1}\{Y-f^*(\boldsymbol{X})\}|^2\} < \infty$, *then* $\sqrt{n}(\widehat{\theta}_{AIPW}-\theta^*) \to \mathrm{N}(0,\sigma^2)$, *in distribution, where* $\pi^*(\boldsymbol{X}) = \mathrm{Pr}(\delta=1 \mid \boldsymbol{X})$ *and* $\sigma^2 = \mathrm{var}\left[f^*(\boldsymbol{X}) + \delta\pi^*(\boldsymbol{X})^{-1}\{Y - f^*(\boldsymbol{X})\}\right]$.

It is worthy pointing out that the derived result is particularly attractive given the fact that the central limit theorem is built by nonparametric estimation of $f^*(\boldsymbol{x})$ with diverging dimension, and to our knowledge, such a result is novel in literature. More importantly, the variance term $\sigma^2$ can be estimated by the sample variance of $\{\widehat{f}_0(\boldsymbol{x}_i)+\delta_i\widehat{\pi}(\boldsymbol{x}_i)^{-1}\{y_i-\widehat{f}_0(\boldsymbol{x}_i)\} :$ $i = 1, \ldots, n\}$:

$$\widehat{\sigma}^2_{AIPW} = \frac{1}{n-1}\sum_{i=1}^{n}(\widehat{y}_i - \widehat{\theta}_{AIPW})^2,$$

where $\widehat{y}_i = \widehat{f}_0(\boldsymbol{x}_i) + \delta_i\widehat{\pi}(\boldsymbol{x}_i)^{-1}\{y_i - \widehat{f}_0(\boldsymbol{x}_i)\}$. Thus, based on Theorem 2 and the estimated variance $\widehat{\sigma}^2_{AIPW}$, we can also obtain the interval estimators of $\theta^*$.

# 4 Numerical analysis

In this section, we compared the numerical performance of the proposed AIPW estimator, denoted as Prop, against several state-of-the-art competitors under two simulated experiments and a real-data application. For Prop, in all the scenarios, we applied a Gaussian kernel, $K(\boldsymbol{x},\boldsymbol{u}) = \exp\left(-\|\boldsymbol{x}-\boldsymbol{u}\|_2^2/(2\sigma_n^2)\right)$ with $\sigma_n$ being the median of all the pairwise distances among the covariates (Jaakkola et al.; 1999). As suggested by He et al. (2021), we also applied the stability-based selection criterion (Sun et al.; 2013) to determine the thresholding value $v_n$ and set the ridge parameter $\lambda_n = 0.001$ for the employed sparse learning algorithm.

## 4.1 Simulated experiments

In this section, we considered $n \in \{800, 1\,000\}$ and $p \in \{400, 2\,000\}$, and covariates were generated by $x_{il} \sim U(-0.5, 0.5)$ for $i = 1, \ldots, n$ and $l = 1, \ldots, p$, where $x_{il}$ denoted the $l$th element of $\boldsymbol{x}_i$, and $U(-0.5, 0.5)$ denoted a uniform distribution over $[-0.5, 0.5]$. The following regression models were applied to generate the response of interest:

M1. Linear regression model: $y_i = 5x_{i1} + 6x_{i2} + 4x_{i3} + 4x_{i4} + \epsilon_i$ with $\epsilon_i \sim N(0, 1)$.

M2. Nonlinear regression model: $y_i = 6x_{i1} + 4(2x_{i2}+1)(2x_{i3}-1) + 6h(x_{i4}) + 5\sin(x_{i5}\pi)/\{2 - \sin(x_{i5}\pi)\} + \epsilon_i$, where $h(x) = 0.1\sin(x_{i4}\pi) + 0.2\cos(x_{i4}\pi) + 0.3\sin(x_{i4}\pi)^2 + 0.4\cos(x_{i4}\pi)^3 + 0.5\sin(x_{i4}\pi)^3$, and $\epsilon_i \sim N(0, 1)$.

For $i = 1, \ldots, n$, the response indicator $\delta_i$ was generated by a Bernoulli distribution with success probability $\pi^*(\boldsymbol{x}_i)$, which was obtained by the following models:

R1. Logistic response model: $\text{logit}\{\pi^*(\boldsymbol{x}_i)\} = -0.1 + 2x_{i1} + 2x_{i3}$.

R2. Multi-modal response model $\pi^*(\boldsymbol{x}_i) = \sin(6x_{i2} + 8x_{i4})/3 + 0.5$.

The linear regression model M1 is commonly assumed in practice (Fan and Li; 2001). The nonlinear regression model M2, however, is more complex, and the interaction effect is also taken into consideration. The logistic response model R1 is widely used in practice. However, the response model R2 violates (2.2), so it is used to test the robustness of the proposed AIPW estimator.

The primary interest was to estimate $\theta^* = \mathbb{E}(Y)$. For the regression model M1, we had $\theta^* = 0$. However, instead of deriving $\theta^*$ analytically, we used $\tilde{\theta} = L^{-1}\sum_{l=1}^{L} y_l$ as the "true value" for the regression model M2, where $\{y_l : l = 1, \ldots, L\}$ was a random sample of size $L = 1\,000\,000$. The following competitors were considered:

CC. The sample mean of the complete cases, $\widehat{\theta}_{cc} = m^{-1}\sum_{i=1}^{n} \delta_i y_i$, where $m = \sum_{i=1}^{n} \delta_i$.

PS. Conventional propensity score estimator $\widehat{\theta}_{ps} = n^{-1}\sum_{i=1}^{n}\delta_i\pi^{-1}(\boldsymbol{x}_i;\widehat{\boldsymbol{\beta}})y_i$, where $\widehat{\boldsymbol{\beta}}^{\top} = (\widehat{\beta}_0,\widehat{\boldsymbol{\beta}}_1^{\top})$ solves $\sum_{i=1}^{n}\{\delta_i - \pi(\boldsymbol{x}_i;\widehat{\boldsymbol{\beta}})\}(1,\boldsymbol{x}_i^{\top}) = 0$ without consideration of the sparsity.

DI. Deterministic imputation using kernel ridge regression (Wang and Kim; 2021) $\widehat{\theta}_{di} = n^{-1}\sum_{i=1}^{n}\{\delta_i y_i + (1 - \delta_i)\widehat{f}(\boldsymbol{x}_i)\}$, where $\widehat{f}(\boldsymbol{x})$ is the fitted kernel ridge regression model based on the fully observed data $\{(\boldsymbol{x}_i, y_i) : \delta_i = 1\}$ without employing sparse learning.

NAIPW. Naive AIPW estimator $\widehat{\theta}_{AIPW1} = n^{-1}\sum_{i=1}^{n}\{\widehat{f}(\boldsymbol{x}_i) + \delta_i\pi^{-1}(\boldsymbol{x}_i;\widehat{\boldsymbol{\beta}})\{y_i - \widehat{f}(\boldsymbol{x}_i)\}\}$, where $\widehat{\boldsymbol{\beta}}$ is the same as that in the PS estimator without consideration of the sparsity, and $\widehat{f}(\boldsymbol{x})$ is the same as that in the DI estimator without employing sparse learning.

The CC estimator completely ignore the unobserved data, leading to a biased estimator if $\mathbb{E}(\delta_i \mid \boldsymbol{x}_i)$ involves covariates used in the regression model. PS estimator is widely used in causal inference (Rosenbaum and Rubin; 1983) and missing data analysis (Wooldridge; 2007). The imputation methods are commonly used to provide a complete dataset, especially in survey sampling; see Kim and Shao (2013, Chapter 4) for details. Recently, Wang and Kim (2021) has proposed a kernel-based deterministic imputation method, and we consider their method for comparison as well. Except for the proposed AIPW estimator, we also considered the naive AIPW estimator based on the conventional propensity score estimator and the deterministic imputation estimator.

We conducted $M = 500$ Monte Carlo simulations for each estimator under different model setups. First, we compared different estimators in terms of the Monte Carlo bias and the Monte Carlo standard error:

$$\mathrm{Bias} = \bar{\theta}_n^{(M)} - \theta,$$

$$\mathrm{SE} = \left\{\frac{1}{M-1}\sum_{m=1}^{M}(\widehat{\theta}_n^{(m)} - \bar{\theta}_n^{(M)})^2\right\}^{1/2},$$

where $\bar{\theta}_n^{(M)} = M^{-1} \sum_{m=1}^{M} \widehat{\theta}_n^{(m)}$, $\widehat{\theta}_n^{(m)}$ was a specific estimator of $\theta$ for the $m$th Monte Carlo simulation. Simulation results were summarized in Table 1. The CC estimator is biased since the response of interest $y_i$ is correlated with the response index $\delta_i$. Even though a logistic model is correctly specified for the response model R1, the PS estimator is still biased or even unrealistic due to the curse of dimensionality. Since the NAIPW estimator was obtained using the same response model as the PS estimator, it is also questionable, especially when the sample size is small and the number of useless covariates is large. Since the response probability was not used by the DI estimator, it does not suffer the same problem as the PS estimator. However, even under the linear regression model, the bias of the DI estimator may not be negligible compared with its standard error. Compared with its competitors, the proposed AIPW estimator performs the best since it has the smallest bias under most model setups, and its standard error is reasonably small.

Next, the proposed AIPW estimator was evaluated by the relative bias of the variance estimator and its coverage rate of a 95% confidence interval:

$$
\mathrm{RB} = \frac{\bar{\sigma}_n^{2(M)} - \mathrm{SE}^2}{\mathrm{SE}^2},
$$

$$
\mathrm{CR} = \frac{1}{M} \sum_{m=1}^{M} \mathbb{I}(\widehat{\theta}_n^{(m)} - 1.96\widehat{\sigma}_n^{(m)} \leq \theta \leq \widehat{\theta}_n^{(m)} + 1.96\widehat{\sigma}_n^{(m)}),
$$

where $\bar{\sigma}_n^{2(M)} = M^{-1} \sum_{i=1}^{M} \widehat{\sigma}_n^{2(m)}$, $\widehat{\sigma}_n^{2(m)}$ is the variance estimator for the $m$th Monte Carlo simulation, $\widehat{\sigma}_n^{(m)}$ is the square root of $\widehat{\sigma}_n^{2(m)}$, and $\mathbb{I}(a \leq x \leq b)$ is an indicator function of $u$ for given $a \leq b$, and it takes value 1 if $u \in [a, b]$ and 0 otherwise. We conducted 500 Monte Carlo simulations, and Table 2 summarized the corresponding results. Generally, as the sample size increases from $n = 800$ to $n = 1\,000$, the relative bias of the variance estimator decreases, and it is negligible if sample size $n = 1\,000$. Thus, the proposed variance estimator performs well, especially when the sample size is large. The coverage rates are close to their nominal

Table 1: Summary of the Monte Carlo bias (Bias) and standard error (SE) corresponding to the five estimators under different model setups, and the unit is 0.1. For "Model", "C1–C4" represent (M1, R1),(M2, R1), (M2, R2) and (M2, R2), respectively. For "Size", "I–IV" corresponds to $(n, p) = (800, 400), (1\,000, 400), (800, 2\,000)$, and $(1\,000, 2\,000)$, respectively. Notation "-" is used when the absolute value of either bias or standard error is greater than 100.

| Model | Size | CC | PS | DI | NAIPW | Prop |
|-------|------|------|------|------|-------|------|
| C1 | I | 6.9 (1.5) | - (-) | 0.9 (1.1) | - (-) | 0.3 (1.2) |
| | II | 7.0 (1.3) | - (-) | 0.9 (1.0) | - (-) | 0.1 (1.0) |
| | III | 7.0 (1.5) | 3.4 (0.7) | 1.1 (1.2) | 1.1 (1.2) | 0.4 (1.3) |
| | IV | 6.9 (1.3) | 3.3 (0.6) | 1.0 (1.0) | 1.0 (1.0) | 0.2 (1.1) |
| C2 | I | -0.7 (2.8) | - (-) | 0.1 (2.0) | - (-) | 0.0 (2.0) |
| | II | -0.8 (2.4) | - (-) | -0.1 (1.7) | - (-) | -0.1 (1.6) |
| | III | -0.7 (2.6) | 2.5 (1.3) | 0.3 (2.1) | 0.3 (2.1) | 0.0 (2.0) |
| | IV | -0.6 (2.3) | 2.6 (1.1) | 0.0 (1.7) | 0.0 (1.7) | -0.1 (1.7) |
| C3 | I | -1.2 (1.5) | - (-) | -0.3 (1.1) | - (-) | -0.2 (1.1) |
| | II | -1.3 (1.3) | 0.3 (57.1) | -0.3 (1.0) | 0.2 (8.4) | -0.2 (1.0) |
| | III | -1.2 (1.5) | -0.6 (0.7) | -0.2 (1.1) | -0.2 (1.1) | -0.1 (1.1) |
| | IV | -1.2 (1.3) | -0.6 (0.7) | -0.2 (1.0) | -0.2 (1.0) | -0.2 (1.0) |
| C4 | I | 1.7 (2.6) | - (-) | 0.2 (1.9) | - (-) | 0.0 (1.9) |
| | II | 1.7 (83.6) | -3.0 (2.2) | 0.2 (1.6) | -0.2 (17.9) | 0.0 (1.7) |
| | III | 1.8 (2.6) | 3.6 (1.3) | 0.3 (2.0) | 0.3 (2.0) | 0.1 (2.0) |
| | IV | 1.9 (2.2) | 3.6 (1.1) | 0.3 (1.6) | 0.3 (1.6) | 0.0 (1.6) |

truth 0.95 when sample size is large. Since the variance of the proposed AIPW estimator is under-estimated for the setup with regression model $M1$ and response model $R1$, the corresponding coverage rate is much lower than 0.95. As the sample size increases, however, the coverage rate gets closer to its nominal truth. For the two setups with response model R2, a logistic regression model is wrongly specified for the response indicator. However, the absolute values of the relative bias of the variance estimator are generally less than 0.05 and the corresponding coverage rates are close to the nominal truth 0.95, specially when the sample size is large. Thus, the proposed AIPW estimator is indeed robust against a wrongly specified response model.

Table 2: Relative bias (RB) of the variance estimator and coverage rate (CR) of the 95% confidence interval of $\theta$ for the proposed method. For "Model", "C1–C4" represent (M1, R1),(M2, R1), (M2, R2) and (M2, R2), respectively. "I–IV" corresponds to $(n, p) = (800, 400)$, $(1\,000, 400)$, $(800, 2\,000)$, and $(1\,000, 2\,000)$, respectively.

| Model | RB | | | | CR | | | |
|---|---|---|---|---|---|---|---|---|
| | I | II | III | IV | I | II | III | IV |
| C1 | -0.10 | -0.07 | -0.33 | -0.20 | 0.92 | 0.94 | 0.88 | 0.92 |
| C2 | -0.13 | 0.01 | -0.14 | -0.01 | 0.94 | 0.94 | 0.93 | 0.95 |
| C3 | 0.04 | -0.02 | -0.03 | 0.05 | 0.95 | 0.95 | 0.95 | 0.96 |
| C4 | -0.10 | -0.01 | -0.14 | 0.05 | 0.94 | 0.94 | 0.93 | 0.96 |

## 4.2 Application to a supermarket dataset

In this section, the proposed AIPW estimator and its competitors were applied to a supermarket dataset (Wang; 2009), which was collected from a major supermarket located in northern China, consisting of daily sale records of $p = 6\,398$ products on $n = 464$ days. This data included almost all kinds of daily necessities and the response of interest was the number of customers on each day, and the covariates are the daily sale volumes of each product. For simplicity, denote $y_i$ and $\boldsymbol{x}_i = (x_{i1}, \ldots, x_{ip})^\top$ to the be the response of interest and the corresponding covariate for the $i$th day. In this section, we were interested in estimating the average number of customers visiting the supermarket. This dataset was fully observed, and it was studentized before analyzing. Thus, the sample mean $\widehat{\theta} = n^{-1} \sum_{i=1}^{n} y_i = 0$ served as a benchmark.

To compare the performance of the proposed AIPW estimator with other competitors, we considered the following missing mechanism for $y_i$:

$$\text{logit}\{\pi(\boldsymbol{x}_i)\} = 1 - 0.6x_{i5} - x_{i6} + 0.5x_{i10}, \tag{10}$$

and $y_i$ was treated as observed if and only if $\delta_i = 1$, where $\delta_i = 1$ with probability $\pi_i$. The mechanism in (10) was MAR, and the corresponding covariates were identified as informative

for estimating the response of interest by He et al. (2021). Then, instead of observing the whole data, we assumed that all the covariates and only $\{y_i : \delta_i = 1\}$ were available, and the resulting response rate was about 0.70.

We generated 500 incomplete datasets using (10) and compare the estimators in Section 4.1. Table 3 summarized the average of the estimators and the corresponding standard error. The CC estimator is highly biased since it ignores the missing mechanism. The performance of the PS estimator is also questionable in that the response model using all covariates results in overfitting. The DI, NAIPW and Prop estimators outperform the CC and PS estimators since their estimates are much closer to 0. However, the standard error of the Prop estimator is much smaller than the other two, illustrating the superior of the proposed AIPW estimator.

Table 3: The average and standard error of 500 incomplete datasets for estimating the number of customers visiting the supermarket. Since studentized is applied, the sample mean 0 serves as the benchmark.

|                | CC    | PS    | DI    | NAIPW | Prop  |
|----------------|-------|-------|-------|-------|-------|
| Estimate       | -0.20 | -0.13 | -0.04 | -0.04 | -0.04 |
| Standard error | 0.03  | 0.02  | 0.03  | 0.03  | 0.01  |

# 5  Conclusion

In this paper, we propose a novel AIPW estimator to infer the population mean, which incorporates an efficient nonparametric imputation with sparse structure and a penalized propensity score estimator under the assumption of missing at random. The proposed method is computationally efficient and allows the dimension diverging. More importantly, the estimation consistency as well as the corresponding central limit theorem are established under regularity assumptions. Its superior is also supported by several simulated examples and one application to a supermarket dataset.

# Acknowledgement

# A    Block co-ordinate gradient descent algorithm

For $g$th group of $\boldsymbol{\beta}_1$, consider a vector $\boldsymbol{d}$ such that $\boldsymbol{d}_k = 0$ for $k \neq g$, and assume that the $\mathrm{df}_g \times \mathrm{df}_g$ submatrix is of the form $H_{gg}^{(t)} = h_g^{(t)} I_{\mathrm{df}_g}$ for some scalar $h_g^{(t)}$, where $I_m$ is an $m \times m$ identity matrix.

If $\|\nabla l(\widehat{\boldsymbol{\beta}}^{(t)})_g - h_g^{(t)} \widehat{\boldsymbol{\beta}}_g^{(t)}\|_2 \leq \lambda_2 \mathrm{df}_g^{1/2}$, let $\boldsymbol{d}_g^{(t)} = -\widehat{\boldsymbol{\beta}}_g^{(t)}$. Otherwise,

$$\boldsymbol{d}_g^{(t)} = -\frac{1}{h_g^{(t)}} \left\{ \nabla l(\widehat{\boldsymbol{\beta}}^{(t)})_g - \lambda_2 \mathrm{df}_g^{1/2} \frac{\nabla l(\widehat{\boldsymbol{\beta}}^{(t)})_g - h_g^{(t)} \widehat{\boldsymbol{\beta}}_g^{(t)}}{\|\nabla l(\widehat{\boldsymbol{\beta}}^{(t)})_g - h_g^{(t)} \widehat{\boldsymbol{\beta}}_g^{(t)}\|_2} \right\}.$$

If $\boldsymbol{d}^{(t)} \neq 0$, let $\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + \alpha^{(t)} \boldsymbol{d}^{(t)}$, where $\alpha^{(t)}$ is the largest value among $\{\alpha_0 \delta^l : l \geq 0\}$ such that

$$S_{\lambda_2}(\boldsymbol{\beta}^{(t)} + \alpha^{(t)} \boldsymbol{d}^{(t)}) - S_{\lambda_2}(\boldsymbol{\beta}^{(t)}) \leq \alpha^{(t)} \sigma \Delta^{(t)},$$

$\delta \in (0, 1)$, $\sigma \in (0, 1)$, $\alpha_0 > 0$, and

$$\Delta^{(t)} = -\left(\boldsymbol{d}^{(t)}\right)^\top \nabla l(\widehat{\boldsymbol{\beta}}^{(t)}) + \lambda_2 \mathrm{df}_g^{1/2} \|\widehat{\boldsymbol{\beta}}_g^{(t)} + \boldsymbol{d}_g^{(t)}\|_2 - \lambda_2 \mathrm{df}_g^{1/2} \|\widehat{\boldsymbol{\beta}}_g^{(t)}\|_2.$$

See Meier et al. (2008) for details.

# B Proofs

**Proposition 1.** *Suppose Assumptions 1–2 are satisfied. Then, with probability at least $1 - \delta_n/2$, there holds*

$$\|\widehat{f} - f^*\|_K \leq 2\log\left(\frac{8}{\delta_n}\right)\left[\frac{3\kappa_{1,p}}{n^{1/2}\lambda_n}\left\{\kappa_{1,p}\|f^*\|_K + q^{-1}\left(\log\frac{4c_1 n}{\delta_n}\right)\right\}\right.$$
$$\left. + \lambda_n^{r-1/2}\|L_K^{-r}f^*\|_2\right].$$

The proof of Proposition 1 is similar as that in He et al. (2021) and thus we omit it here.

*Proof of Lemma 1.* The proof of Lemma 1 is similar as that in He et al. (2021) by using Proposition 1, the property of Hilbert-Schmidt operators and the concentration inequalities in Hilbert-Schmidt operator space. Thus we omit the detail here. $\square$

*Proof of Lemma 2.* The proof of Lemma 2 is similar as that in He et al. (2021), and thus we omit the detail here. $\square$

*Proof of Theorem 1.* Define the event that

$$\mathcal{C}_1 = \left\{\|\widehat{f}_0 - f^*\|_\infty > c_{m,2}\max\left\{\kappa_{1,p_0}\|f^*\|_K, q^{-1}\left(\log\frac{2c_1 m}{\delta_m}\right)\right\}\right.$$
$$\left. \times \log\left(\frac{4}{\delta_m}\right)m^{-\frac{2r-1}{2(2r+1)}}\right\}. \tag{11}$$

Then, the probability $\Pr(\mathcal{C}_1)$ can be decomposed as

$$\Pr\left(\mathcal{C}_1\right) = \Pr\left(\mathcal{C}_1, \{\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\}\right) + \Pr\left(\mathcal{C}_1, \{\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\}\right)$$
$$= \Pr\left(\mathcal{C}_1 \mid \{\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\}\right)\Pr\left(\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\right)$$
$$+ \Pr\left(\mathcal{C}_1 \mid \{\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\}\right)\Pr\left(\widehat{\mathcal{A}}_{v_m} \neq \mathcal{A}^*\right)$$
$$\leq \Pr\left(\mathcal{C}_1 \mid \{\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\}\right)(1 - \Delta_m) + \Delta_m.$$

By Lemma 2, we have $\Delta_m \to 0$ and $(1 - \Delta_m) \to 1$. For $\Pr(\mathcal{C}_1 \mid \{\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\})$, by applying the proof in Proposition 1 conditioning on $\{\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\}$, with probability at least $1 - \delta_m$, there holds

$$
\begin{aligned}
\left\|\widehat{f} - f^*\right\|_K \leq{} & c_{m,2} \max\left\{\kappa_{1,p_0}\|f^*\|_K, q^{-1}(\log \tfrac{2c_1 m}{\delta_m})\right\} \\
& \times \log\left(\tfrac{4}{\delta_m}\right) m^{-\frac{2r-1}{2(2r+1)}},
\end{aligned}
$$

which implies $\Pr(\mathcal{C}_1 \mid \{\widehat{\mathcal{A}}_{v_m} = \mathcal{A}^*\}) \leq \delta_m$. Combining the above results, we have $\Pr(\mathcal{C}_1) \leq \delta_m + \Delta_m$. This completes the proof of the first part in Theorem 1.

Additionally, by Assumption 4, we have $m = O(n)$, and if we take $r = 1$ and assume that $p_0 = O(1)$ and $\epsilon$ has sub-Gaussian or sub-exponential tail, there holds

$$
\left\|\widehat{f} - f^*\right\|_K = O_\mathbb{P}(n^{-\frac{1}{6}} \log n).
$$

Note that

$$
\begin{aligned}
\left\|\widehat{f} - f^*\right\|_\infty &= \sup_{\boldsymbol{x}} |\widehat{f}(\boldsymbol{x}) - f^*(\boldsymbol{x})| \\
&= \sup_{\boldsymbol{x}} |\langle \widehat{f} - f^*, K_{\boldsymbol{x}}\rangle_K| \leq \kappa_{1,p_0}\left\|\widehat{f} - f^*\right\|_K,
\end{aligned}
$$

which completes the proof. $\qquad\square$

*Proof of Lemma 3.* By Assumptions 4–7, Meier et al. (2008) showed (9). Denote $g(x) = \{1 + \exp(-x)\}^{-1}$, and we can show that $g'(x) = \mathrm{d}g(x)/\mathrm{d}x = g(x)\{1 - g(x)\}$. That is,

$$
|g'(x)| \leq 1, \tag{12}
$$

for any $x$ by the fact that $0 \leq g(x) \leq 1$. Thus, by (12) and the mean value theorem, we

24

conclude that $g(x)$ is Lipschitz continuous in the sense that

$$|g(x_1) - g(x_2)| \leq |x_1 - x_2|,$$

for any $x_1$ and $x_2$ in $\mathbb{R}$. By noting the fact that $\widehat{\pi}(\boldsymbol{x}) = g\{\eta_{\widehat{\boldsymbol{\beta}}_{\lambda_2}}(\boldsymbol{x})\}$ and $\pi^*(\boldsymbol{x}) = g\{\eta_{\boldsymbol{\beta}_0}(\boldsymbol{x})\}$, by (9), we have

$$
\begin{aligned}
\mathbb{E}\{|\widehat{\pi}(\boldsymbol{X}) - \pi^*(\boldsymbol{X})|^2\} &\leq \mathbb{E}\{|\eta_{\widehat{\boldsymbol{\beta}}_{\lambda_2}}(\boldsymbol{X}) - \eta_{\boldsymbol{\beta}_0}(\boldsymbol{X})|^2\} \\
&= O_{\mathbb{P}}(n^{-2/3 - 2\zeta}),
\end{aligned}
\tag{13}
$$

where the expectation is taken conditional on $\widehat{\boldsymbol{\beta}}_{\lambda_2}$. By (13), we have shown Lemma 3.

$\square$

**Lemma 4.** *Suppose Assumptions 4–7 are satisfied. Then, given $\widehat{\boldsymbol{\beta}}_{\lambda_2}$, there holds*

$$\max\{|\widehat{\beta}_k - \beta_k^*| : k = 0, \ldots, p\} = O_{\mathbb{P}}(n^{-1/3 - \zeta}),$$

*where $\widehat{\beta}_k$ and $\beta_k^*$ are the $(k+1)$th component of $\widehat{\boldsymbol{\beta}}_{\lambda_2}$ and $\boldsymbol{\beta}^*$, respectively.*

*Proof of Lemma 4.* Given the estimated parameters for the response model, (9) can be re-expressed as

$$
\begin{aligned}
&\mathbb{E}\{|\eta_{\widehat{\boldsymbol{\beta}}_{\lambda_2}}(\boldsymbol{X}) - \eta_{\boldsymbol{\beta}^*}(\boldsymbol{X})|^2\} \\
&= \mathbb{E}\{(\widehat{\boldsymbol{\beta}}_{\lambda_2} - \boldsymbol{\beta}^*)^\top \boldsymbol{X}\boldsymbol{X}^\top (\widehat{\boldsymbol{\beta}}_{\lambda_2} - \boldsymbol{\beta}^*)\} \\
&= (\widehat{\boldsymbol{\beta}}_{\lambda_2} - \boldsymbol{\beta}^*)^\top \mathbb{E}(\boldsymbol{X}\boldsymbol{X}^\top)(\widehat{\boldsymbol{\beta}}_{\lambda_2} - \boldsymbol{\beta}^*)\} \\
&= O_{\mathbb{P}}(n^{-2/3 - 2\zeta}),
\end{aligned}
\tag{14}
$$

25

By Assumption A5 and (14), we conclude that

$$(\widehat{\boldsymbol{\beta}}_{\lambda_2} - \boldsymbol{\beta}^*)^\top (\widehat{\boldsymbol{\beta}}_{\lambda_2} - \boldsymbol{\beta}^*) \;=\; (\widehat{\beta}_0 - \beta_0^*)^2 + \sum_{k=1}^{p} (\widehat{\beta}_k - \beta_k^*)^2$$

$$= \; O_{\mathbb{P}}(n^{-2/3 - 2\zeta}). \tag{15}$$

Notice that

$$\max\{(\widehat{\beta}_k - \beta_k^*)^2 : k = 0, \ldots, p\} \leq (\widehat{\beta}_0 - \beta_0^*)^2 + \sum_{k=1}^{p} (\widehat{\beta}_k - \beta_k^*)^2. \tag{16}$$

Thus, we have proved Lemma 4 by (15) and (16). $\qquad\square$

*Proof of Theorem 2.* For simplicity, denote $\pi_i^* = \pi^*(\boldsymbol{x}_i)$ and $\widehat{\pi}_i = \widehat{\pi}(\boldsymbol{x}_i) = (1 + \exp[-\{(1, \boldsymbol{x}_i^\top)\widehat{\boldsymbol{\beta}}_{\lambda_2}\}])^{-1}$. What if we consider

$$\widehat{\theta}_{AIPW} = \frac{1}{n} \sum_{i=1}^{n} \left[ \widehat{f}_0(\boldsymbol{x}_i) + \frac{\delta_i}{\widehat{\pi}_i} \left\{ y_i - \widehat{f}_0(\boldsymbol{x}_i) \right\} \right],$$

where $\widehat{\pi}_i$ is an estimator of $\Pr(\delta_i = 1 \mid \boldsymbol{x}_i)$ by the group lasso for logistic regression.

Then, we have

$$
\begin{aligned}
\widehat{\theta}_{AIPW} &= \frac{1}{n}\sum_{i=1}^{n}\left[ f^*(\boldsymbol{x}_i) + \left\{\widehat{f}_0(\boldsymbol{x}_i) - f^*(\boldsymbol{x}_i)\right\} + \frac{\delta_i}{\pi_i^*}\{y_i - f^*(\boldsymbol{x}_i)\} \right. \\
&\quad + \frac{\delta_i}{\widehat{\pi}_i}\{y_i - f^*(\boldsymbol{x}_i)\} - \frac{\delta_i}{\pi_i^*}\{y_i - f^*(\boldsymbol{x}_i)\} \\
&\quad \left. + \frac{\delta_i}{\widehat{\pi}_i}\{f^*(\boldsymbol{x}_i) - \widehat{f}_0(\boldsymbol{x}_i)\} \right] \\
&= \frac{1}{n}\sum_{i=1}^{n}\left[ f^*(\boldsymbol{x}_i) + \frac{\delta_i}{\pi_i^*}\{y_i - f^*(\boldsymbol{x}_i)\} \right] \\
&\quad + \frac{1}{n}\sum_{i=1}^{n}\left\{1 - \frac{\delta_i}{\widehat{\pi}_i}\right\}\left\{\widehat{f}_0(\boldsymbol{x}_i) - f^*(\boldsymbol{x}_i)\right\} \\
&\quad + \frac{1}{n}\sum_{i=1}^{n}\left[\frac{\delta_i}{\widehat{\pi}_i} - \frac{\delta_i}{\pi_i^*}\right]\{y_i - f^*(\boldsymbol{x}_i)\} \\
&= \frac{1}{n}\sum_{i=1}^{n}\left[ f^*(\boldsymbol{x}_i) + \frac{\delta_i}{\pi_i^*}\{y_i - f^*(\boldsymbol{x}_i)\} \right] \\
&\quad + \frac{1}{n}\sum_{i=1}^{n}\left\{1 - \frac{\delta_i}{\widehat{\pi}_i}\right\}\left\{\widehat{f}_0(\boldsymbol{x}_i) - f^*(\boldsymbol{x}_i)\right\} \\
&\quad + \frac{1}{n}\sum_{i=1}^{n}\left[\frac{\delta_i}{\widehat{\pi}_i} - \frac{\delta_i}{\pi_i^*}\right]\epsilon_i,
\end{aligned}
\tag{17}
$$

where $\epsilon_i = y_i - f^*(\boldsymbol{x}_i)$.

First, we consider the first term of (17), and we have

$$
\begin{aligned}
&\mathbb{E}\left[ f^*(\boldsymbol{X}) + \frac{\delta}{\pi^*(\boldsymbol{X})}\{Y - f^*(\boldsymbol{X})\} \right] \\
&= \mathbb{E}\left( \mathbb{E}\left[ f^*(\boldsymbol{X}) + \frac{\delta}{\pi^*(\boldsymbol{X})}\{Y - f^*(\boldsymbol{X})\} \right] \mid \boldsymbol{X}, Y \right) \\
&= \mathbb{E}(Y),
\end{aligned}
$$

where $\delta$ is a binary random variable with success probability $\pi^*(\boldsymbol{X})$ conditional on $\boldsymbol{X}$. Since $\mathbb{E}\{|f^*(\boldsymbol{X}) + \delta\pi^*(\boldsymbol{X})^{-1}\{Y - f^*(\boldsymbol{X})\}|^2\} < \infty$, by the classical central limit theorem

(Van der Vaart; 2000, Example 2.1), we have

$$\sqrt{n}\left(\frac{1}{n}\sum_{i=1}^{n}\left[f^*(\boldsymbol{x}_i)+\frac{\delta_i}{\pi_i^*}\{y_i-f^*(\boldsymbol{x}_i)\}\right]-\theta^*\right)\to \mathrm{N}(0,\sigma^2), \tag{18}$$

in distribution under regularity conditions, where $\theta = \mathbb{E}(y)$ and $\sigma^2$ is to be estimated.

Next, we consider the third term of (17) .

$$\frac{1}{n}\sum_{i=1}^{n}\left[\frac{\delta_i}{\widehat{\pi}_i}-\frac{\delta_i}{\pi_i^*}\right]\epsilon_i \;=\; \frac{1}{n}\sum_{i=1}^{n}\frac{\delta_i\epsilon_i}{\pi_i^*}\frac{\pi_i^*-\widehat{\pi}_i}{\widehat{\pi}_i}.$$

By Assumption 4 and Lemma 3, we conclude that $(\pi_i^* - \widehat{\pi}_i)\widehat{\pi}_i^{-1} = o_{\mathbb{P}}(1)$ uniformly for $i = 1, \ldots, n$. Since

$$\frac{1}{n}\sum_{i=1}^{n}\frac{\delta_i\epsilon_i}{\pi_i^*} = O_{\mathbb{P}}(n^{-1/2}),$$

we conclude that

$$\frac{1}{n}\sum_{i=1}^{n}\left[\frac{\delta_i}{\widehat{\pi}_i}-\frac{\delta_i}{\pi_i^*}\right]\epsilon_i = o_{\mathbb{P}}(n^{-1/2}). \tag{19}$$

Now, we focus on the second term of (17), and consider

$$\frac{1}{n}\sum_{i=1}^{n}\left\{1-\frac{\delta_i}{\widehat{\pi}_i}\right\}\left\{\widehat{f}_0(\boldsymbol{x}_i)-f^*(\boldsymbol{x}_i)\right\}$$

$$=\quad \frac{1}{n}\sum_{i=1}^{n}\left\{1-\frac{\delta_i}{\pi_i^*}\right\}\left\{\widehat{f}_0(\boldsymbol{x}_i)-f^*(\boldsymbol{x}_i)\right\}$$

$$+\frac{1}{n}\sum_{i=1}^{n}\frac{\delta_i}{\pi_i^*\widehat{\pi}_i}(\widehat{\pi}_i-\pi_i^*)\left\{\widehat{f}_0(\boldsymbol{x}_i)-f^*(\boldsymbol{x}_i)\right\}.$$

$$\tag{20}$$

By Theorem 1, $\widehat{f}_0(\boldsymbol{x}_i)-f^*(\boldsymbol{x}_i) = O_{\mathbb{P}}(\log(n)n^{-1/6})$ uniformly for $i = 1, \ldots, n$. Since $n^{-1}\sum_{i=1}^{n}\{1-\delta_i(\pi_i^*)^{-1}\} = O_{\mathbb{P}}(n^{-1/2})$, the first term of is of (20) is of the order $o_{\mathbb{P}}(n^{-1/2})$. Besides, to show

28

the second term of (20) is also of the order $o_{\mathbb{P}}(n^{-1/2})$, it is enough to show

$$\frac{1}{n}\sum_{i=1}^{n}\frac{\delta_i}{\pi_i^*\widehat{\pi}_i}(\widehat{\pi}_i - \pi_i^*) = O_{\mathbb{P}}(n^{-1/3-\zeta/2}), \tag{21}$$

where $\zeta$ is in Assumption 7.

By Lemma 4, we conclude that $\max\{|\widehat{\beta}_k - \beta_k^*| : k = 0, \ldots, p\} = o_p(1)$. Denote $A_n$ to be the event that $\{\max\{|\widehat{\beta}_k - \beta_k^*| : k = 0, \ldots, p\} \geq C_\kappa\}$, where $C_\kappa$ is a positive constant such that $\min\{\widehat{\boldsymbol{\beta}}_{\lambda_2}(\boldsymbol{x}) \geq \kappa/2 : \boldsymbol{x} \in \mathcal{X}\}$. The existence of $C_\kappa$ is guaranteed by the compactness of $\mathcal{X}$ and Assumption A4. Then, we have $P(A_n) \to 0$ as $n \to \infty$. On $\boldsymbol{x} \in A_n^C$, we conclude $\widehat{\pi}(\boldsymbol{x}) \geq \kappa/2$.

Since $\{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n\}$ is a random sample, given $\widehat{\boldsymbol{\beta}}_{\lambda_2}$, for any positive constant $C$, we consider

$$
\begin{aligned}
&\Pr\left(\left|\frac{1}{n}\sum_{i=1}^{n}\frac{\delta_i}{\pi^*(\boldsymbol{X}_i)\widehat{\pi}(\boldsymbol{X}_i)}\{\widehat{\pi}(\boldsymbol{X}_i) - \pi^*(\boldsymbol{X}_i)\}\right| \geq Cn^{-1/3-\zeta/2}\right) \\
\leq\ & Pr\left(\left|\frac{1}{n}\sum_{i=1}^{n}\frac{2\delta_i}{\kappa^2}\{\widehat{\pi}(\boldsymbol{X}_i) - \pi^*(\boldsymbol{X}_i)\}\right| \geq Cn^{-1/3-\zeta/2}\right) + P(A_n) \\
\leq\ & \frac{\mathbb{E}[2\kappa^{-2}n^{-1}\sum_{i=1}^{n}[\delta_i\{\pi^*(\boldsymbol{X}_i)\widehat{\pi}(\boldsymbol{X}_i)\}^{-1}\{\widehat{\pi}(\boldsymbol{X}_i) - \pi^*(\boldsymbol{X}_i)\}]^2}{C^2 n^{-2/3-\zeta}} \\
& +P(A_n) \\
\leq\ & \frac{2\sum_{i=1}^{n}\mathbb{E}\{\widehat{\pi}(\boldsymbol{X}_i) - \pi^*(\boldsymbol{X}_i)\}^2}{n\kappa^2 C^2 n^{-2/3-\zeta}} + P(A_n) \\
\leq\ & \frac{\mathbb{E}\{\widehat{\pi}(\boldsymbol{X}_1) - \pi^*(\boldsymbol{X}_1)\}^2}{\kappa^2 C^2 n^{-2/3-\zeta}} + P(A_n) \\
=\ & o_{\mathbb{P}}(1), \tag{22}
\end{aligned}
$$

where $\boldsymbol{X}_i$ is the random variable associated with $\boldsymbol{x}_i$, the first inequality holds by the Markov inequality, the second inequality holds by Assumption 4 and the fact that $\delta_i \leq 1$ for $i = 1, \ldots, n$, the third inequality holds since $\boldsymbol{X}_1, \ldots, \boldsymbol{X}_n$ are identically distributed, and the

fourth inequality holds by Lemma 3. By (22), we have validated (21), so we have

$$\frac{1}{n}\sum_{i=1}^{n}\left\{1 - \frac{\delta_i}{\widehat{\pi}_i}\right\}\left\{\widehat{f}_0(\boldsymbol{x}_i) - f^*(\boldsymbol{x}_i)\right\} = o_{\mathbb{P}}(n^{-1/2}). \tag{23}$$

By (18), (19) and (23), we have proved Theorem 2 by the Slutsky's theorem (Athreya and Lahiri; 2006, Theorem 9.1.6). □

# References

Abadie, A. and Imbens, G. W. (2016). Matching on the estimated propensity score, *Econometrica* **84**(2): 781–807.

Athreya, K. B. and Lahiri, S. N. (2006). *Measure Theory and Probability Theory*, Springer, New York.

Banerjee, O., El Ghaoui, L. and d'Aspremont, A. (2008). Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data, *Journal of Machine Learning Research* **9**(15): 485–516.

Barber, R. and Candes, E. (2019). A knockoff filter for high-dimensional selective inference, *Annals of Statistics* **47**(5): 2504–2537.

Belloni, A., Chernozhukov, V. and Hansen, C. (2013). Inference on treatment effects after selection among high-dimensional controls, *The Review of Economic Studies* **81**(2): 608–650.

Bertsimas, D. and King, A. (2017). Logistic regression: From art to science, *Statistical Science* **32**(3): 367–384.

Chen, M., Ibrahim, J. and Yiannoutsos, C. (1999). Prior elicitation, variable selection and bayesian computation for logistic regression models, *Journal of the Royal Statistical Society: Series B* **61**(1): 223–242.

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, *The Econometrics Journal* **21**(1): C1–C68.

Fan, J. and Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties, *Journal of the American Statistical Association* **96**(456): 1348–1360.

Fan, J. and Lv, J. (2008). Sure independence screening for ultrahigh dimensional feature space (with discussion), *Journal of the Royal Statistical Society, Series B* **70**(5): 849–911.

Fan, J., Xue, L. and Zou, H. (2014). Strong oracle optimality of folded concave penalized estimation, *Annals of statistics* **42**(3): 819.

Fan, J., Yang, F. and Song, R. (2011). Nonparametric independence screening in sparse ultrahigh dimensional additive models, *Journal of the American Statistical Association* **106**(494): 544–557.

Fischer, S. and Steinwart, I. (2020). Sobolev norm learning rates for regularized least-squares algorithms, *Journal of Machine Learning Research* **21**(205): 1–38.

Fukumizu, K. and Leng, C. (2014). Gradient-based kernel dimension reduction for regression, *Journal of the American Statistical Association* **109**(505): 359–370.

Genkin, A., Lewis, D. D. and Madigan, D. (2007). Large-scale Bayesian logistic regression for text categorization, *Technometrics* **49**(3): 291–304.

Graham, B. S., de Xavier Pinto, C. C. and Egel, D. (2012). Inverse probability tilting for moment condition models with missing data, *The Review of Economic Studies* **79**(3): 1053–1079.

He, X., Wang, J. and Lv, S. (2021). Efficient kernel-based variable selection with sparsistency, *Statistica Sinica* **31**: 2123–2151.

Huang, J., Horowitz, J. and Wei, F. (2010). Variable selection in nonparametric additive models, *Annals of Statistics* **38**(4): 2282–2313.

Jaakkola, T., Diekhans, M. and Haussler, D. (1999). Using the Fisher kernel method to detect remote protein homologies, *In Proceedings of Seventh International Conference on Intelligent Systems for Molecular Biology* **99**: 149–158.

Kim, J. and Shao, J. (2013). *Statistical Methods for Handling Incomplete Data*, CRC press, New York.

Kim, Y., Kim, J. and Kim, Y. (2006). Blockwise sparse regression, *Statistica Sinica* **16**(2): 375–390.

Lin, H., Zhou, F., Wang, Q., Zhou, L. and Qin, J. (2018). Robust and efficient estimation for the treatment effect in causal inference and missing data problems, *Journal of Econometrics* **205**: 363–380.

Lin, Y. and Zhang, H. (2006). Component selection and smoothing in multivariate nonparametric regression, *Annals of Statistics* **34**(5): 2272–2297.

Lv, S., Lin, H., Lian, H. and Huang, J. (2018). Oracle inequalities for sparse additive quantile regression in reproducing kernel Hilbert space, *Annals of Statistics* **46**(2): 781–813.

Meier, L., Van de Geer, S. and Bühlmann, P. (2008). The group lasso for logistic regression, *Journal of the Royal Statistical Society, Series B* **70**(1): 53–71.

Mendelson, S. and Neeman, J. (2010). Regularization in kernel learning, *Annals of Statistics* **38**(1): 526–565.

Mercer, J. (1909). Functions of positive and negative type and their connection with the theory of integral equations, *Philosophical Transactions of the Royal Society A* **209**: 415–446.

Ning, Y., Sida, P. and Imai, K. (2020). Robust estimation of causal effects via a high-dimensional covariate balancing propensity score, *Biometrika* **107**(3): 533–554.

Qin, J., Zhang, B. and Leung, D. H. (2017). Efficient augmented inverse probability weighted estimation in missing data problems, *Journal of Business & Economic Statistics* **35**(1): 86–97.

Qu, A., Lindsay, B. G. and Lu, L. (2010). Highly efficient aggregate unbiased estimating functions approach for correlated data with missing at random, *Journal of the American Statistical Association* **105**(489): 194–204.

Ravikumar, P., Wainwright, M. and Lafferty, J. (2010). High-dimensional Ising model selection using $\ell^1$-regularized logistic regression, *Annals of Statistics* **38**(3): 1287 – 1319.

Robins, J. M., Rotnitzky, A. and Zhao, L. P. (1994). Estimation of regression coefficients when some regressors are not always observed, *Journal of the American Statistical Association* **89**(427): 846–866.

Rosasco, L., Villa, S., Mosci, S., Santoro, M. and Verri, A. (2013). Nonparametric sparsity and regularization, *Journal of Machine Learning Research* **14**(16): 1665–1714.

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

Rubin, D. (1976). Inference and missing data, *Biometrika* **63**(3): 581–592.

Shen, X., Pan, W. and Zhu, Y. (2012). Likelihood-based selection and sharp parameter estimation, *Journal of the American Statistical Association* **107**(497): 223–232.

Shen, X., Pan, W., Zhu, Y. and Zhou, H. (2013). On constrained and regularized high-dimensional regression, *Annals of the Institute of Statistical Mathematics* **65**(5): 807–832.

Shevade, S. and Keerthi, S. (2003). A simple and efficient algorithm for gene selection using sparse logistic regression, *Bioinformatics* **19**(17): 2246–2253.

Smale, S. and Zhou, D. (2007). Learning theory estimates via integral operators and their approximations, *Constructive Approximation* **26**(2): 153–172.

Steinwart, I. (2005). Consistency of support vector machines and other regularized kernel classifiers, *IEEE Transactions on Information Theory* **51**(1): 128–142.

Steinwart, I. and Christmann, A. (2008). *Support Vector Machine*, Springer, New York.

Sun, W., Wang, J. and Fang, Y. (2013). Consistent selection of tuning parameters via variable selection stability, *Journal of Machine Learning Research* **14**(71): 3419–3440.

Tan, Z. (2010). Bounded, efficient and doubly robust estimation with inverse weighting, *Biometrika* **97**(3): 661–682.

Tan, Z. (2020). Model-assisted inference for treatment effects using regularized calibrated estimation with high-dimensional data, *Annals of Statistics* **48**(2): 811 – 837.

Tang, J., Alelyani, S. and Liu, H. (2014). Feature selection for classification: A review, *Data classification: Algorithms and applications* pp. 37–64.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso, *Journal of the Royal Statistical Society, Series B* **58**(1): 267–288.

Van der Vaart, A. (2000). *Asymptotic Statistics*, Cambridge University Press, New York.

Wainwright, M., Ravikumar, P. and Lafferty, J. (2007). High-dimensional graphical model selection using $\ell^1$-regularized logistic regression, *Advances in neural information processing systems* **19**: 1465–1472.

Wang, H. (2009). Forward regression for ultra-high dimensional variable screening, *Journal of the American Statistical Association* **104**(488): 1512–1524.

Wang, H. and Kim, J. (2021). Statistical inference after kernel ridge regression imputation under item nonresponse, *Technical Report (Available at https://arxiv.org/pdf/2102.00058.pdf)* pp. 1–26.

Wang, X. and Leng, C. (2016). High dimensional ordinary least squares projection for screening variables, *Journal of the Royal Statistical Society, Series B* **78**(3): 589–611.

Wooldridge, J. (2007). Inverse probability weighted estimation for general missing data problems, *Journal of Econometrics* **141**(2): 1281–1301.

Yang, L., Lv, S. and Wang, J. (2016). Model-free variable selection in reproducing kernel Hilbert space, *Journal of Machine Learning Research* **17**(82): 1–24.

Yang, S., Kim, J. and Song, R. (2020). Doubly robust inference when combining probability and non-probability samples with high-dimensional data, *Journal of the Royal Statistical Society: Series B* **82**(2): 445–465.

Yuan, M. and Lin, Y. (2006). Model selection and estimation in regression with grouped variables, *Journal of the Royal Statistical Society: Series B* **68**(1): 49–67.

Zhang, C., Liu, Y. and Wu, Y. (2016). On quantile regression in reproducing kernel Hilbert spaces with data sparsity constraint, *Journal of Machine Learning Research* **17**(40): 1–45.

Zhao, P., Tang, N., Qu, A. and Jiang, D. (2017). Semiparametric estimating equations inference with nonignorable missing data, *Statistica Sinica* **27**(1): 89–113.

Zhou, D. (2007). Derivative reproducing properties for kernel methods in learning theory, *Journal of Computational and Applied Mathematics* **220**(1–2): 456–463.

Zou, H. (2006). The adaptive lasso and its oracle properties, *Journal of the American Statistical Association* **101**(476): 1418–1429.