

# UNIVERSITY OF BIRMINGHAM

## Research at Birmingham

### On density and regression estimation with incomplete data

Mojirsheibani, Majid; Manley, Kevin; Pouliot, William

DOI:

[10.1080/03610926.2016.1277751](https://doi.org/10.1080/03610926.2016.1277751)

License:

Other (please specify with Rights Statement)

*Document Version*

Peer reviewed version

*Citation for published version (Harvard):*

Mojirsheibani, M, Manley, K & Pouliot, W 2017, 'On density and regression estimation with incomplete data', *Communications in Statistics: Theory and Methods*. <https://doi.org/10.1080/03610926.2016.1277751>

[Link to publication on Research at Birmingham portal](#)

#### **Publisher Rights Statement:**

This is an Accepted Manuscript of an article published by Taylor & Francis in *Communications in Statistics - Theory and Methods* on 13/01/2017, available online: <http://www.tandfonline.com/10.1080/03610926.2016.1277751>

#### **General rights**

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

#### **Take down policy**

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact [UBIRA@lists.bham.ac.uk](mailto:UBIRA@lists.bham.ac.uk) providing details and we will remove access to the work immediately and investigate.

# On density and regression estimation with incomplete data

By

Majid Mojirsheibani<sup>1</sup> Kevin Manley<sup>2</sup>

Department of Mathematics, California State University, Northridge, CA 91330

and

William Pouliot<sup>3</sup>

Department of Economics, University of Birmingham, Birmingham, UK

## Abstract

We consider the problem of estimation of a density function in the presence of incomplete data and study the Hellinger distance between our proposed estimator and the true density function. Here the presence of incomplete data is handled by utilizing a Horvitz-Thompson-type inverse weighting approach, where the weights are estimates of the unknown *selection probabilities*. We also address the problem of estimating a regression function with incomplete data.

**Keywords:** Convergence, incomplete data, empirical process, kernel, density.

## 1 Introduction

There has been a recent growing interest in the problem of density estimation in the presence of incomplete data. See, for example, Dubnicka (2009), Müller (2012), Tang et al (2012), Wang (2008), and Hazelton (2000) for kernel density estimation when auxiliary variables (covariates) are available, and also Zou et al (2015) for wavelet density estimators with incomplete data. Most of the work in the above cited references can be viewed, primarily, as the counterparts of the classical problem of density estimation with complete data as discussed, for example, by Rosenblatt (1956), Parzen (1962), Prakasa Rao (1983), and Devroye and Györfi (1985). Our interest in this paper is in the problem of density estimation with incomplete data, but with an approach that is closer in spirit to the work of van de Geer's (1993, 2000), where the author uses empirical process theory, based on fully observable data, to establish convergence results for the Hellinger distance between the true and the estimated densities. We will also extend our approach to deal with the problem of regression function estimation.

### 1.1 Background tools from empirical process

The justification and presentation of our main results will be facilitated with the aid of some empirical process theory results that are appropriate for situations where the data may be incomplete. Therefore, out of necessity, our main contributions start with some new results on empirical processes with incomplete data. To fix notation, let  $\psi : \mathbb{R}^{d+p} \rightarrow \mathbb{R}$ ,  $d > 0, p > 0$ , and initially consider the estimation of the "mean"  $\nu(\psi) := \mathbb{E}(\psi(\mathbf{Z}))$  based on the independently and identically distributed (iid) data  $\mathbb{D}_n = \{\mathbf{Z}_1, \dots, \mathbf{Z}_n\}$ . However, here we are interested in the situation where the data  $\mathbf{Z}_i = (\mathbf{X}'_i, \mathbf{V}'_i)' \in \mathbb{R}^{d+p}$ , are not fully observed, more specifically  $\mathbf{X}_i \in \mathbb{R}^d$  is always available, but  $\mathbf{V}_i \in \mathbb{R}^p$  may be unobservable (unavailable) for various unknown reasons. To clarify our setup further, we also define the random variables  $\xi_i = 1$  if  $\mathbf{V}_i$  is observed, and  $\xi_i = 0$  otherwise, and represent the data as

$$\mathbb{D}_n = \{(\mathbf{Z}_1, \xi_1), \dots, (\mathbf{Z}_n, \xi_n)\} = \{(\mathbf{X}_1, \mathbf{V}_1, \xi_1), \dots, (\mathbf{X}_n, \mathbf{V}_n, \xi_n)\}.$$

Some important examples of the function  $\psi$  include  $\psi(\mathbf{Z}) = \psi(\mathbf{X}, V) = V \in \mathbb{R}^1$ , in which case the estimation of  $\mathbb{E}(\psi(\mathbf{Z}))$  reduces to the usual mean estimation for  $\mathbb{E}(V)$  when some of the  $V_i$ 's are not

---

<sup>1</sup>Corresponding author. Email: majid.mojirsheibani@csun.edu .

<sup>2</sup>Email: kevin.manley.66@my.csun.edu

<sup>3</sup>Email: w.pouliot@bham.ac.uk

available. This case has been addressed and studied extensively in the literature; see, for example, Cheng (1994), Wang and Rao (2002), Hirano and Ridder (2003), Wang et al (2004), Rueda et al (2006), Müller (2009), and Kim and Yu (2011). Another example involves estimation of higher moments as well as mixed moments of the components of the random vector  $\mathbf{Z}$ . On the other hand, when the function  $\psi$  is of the form  $\psi(\mathbf{Z}) \equiv \psi_{\mathbf{z}}(\mathbf{Z}) = \mathbb{I}\{\mathbf{Z} \leq \mathbf{z}\}$ ,  $\mathbf{z} \in \mathbb{R}^{d+p}$ , then  $\mathbb{E}(\psi(\mathbf{Z}))$  is the cumulative distribution function (cdf) of  $\mathbf{Z}$ . The case of  $\psi(\mathbf{Z}) \equiv \psi_{\mathbf{v}}(\mathbf{X}', \mathbf{V}') = \mathbb{I}\{\mathbf{V} \leq \mathbf{v}\}$  corresponds to the estimation of the marginal cdf of  $\mathbf{V}$ . These cases have been studied by, for example, Hu et al (2011), Liu et al (2011), Chenouri et al (2009), and Cheng and Chu (1996). When there are no missing  $\mathbf{V}_i$ 's in the data,  $\mathbb{E}(\psi(\mathbf{Z}))$  can be estimated by the classical nonparametric empirical version

$$\nu_n(\psi) = n^{-1} \sum_{i=1}^n \psi(\mathbf{Z}_i).$$

As for performance of this estimator, nonasymptotic exponential bounds are available on the uniform deviations on  $\nu_n(\psi)$  from  $\nu(\psi)$  (uniform in  $\psi$ ) under various conditions. In fact for our future reference, we state one such result (see, for example, Pollard (1984; pp. 26-27):

**Theorem 1** *Let  $\Psi$  be a class of functions  $\psi : \mathbb{R}^{d+p} \rightarrow [-B, B]$ , for some  $0 < B < \infty$ . Then, for every  $\epsilon > 0$  and every  $n \geq 1$ ,  $\mathbb{P}\{\sup_{\psi \in \Psi} |\nu_n(\psi) - \nu(\psi)| > \epsilon\} \leq 8 \mathbb{E}[\mathcal{N}_1(\epsilon/8, \Psi, \mathbb{D}_n)] e^{-n\epsilon^2/(128B^2)}$ .*

Here, the term  $\mathcal{N}_1(\epsilon, \Psi, \mathbb{D}_n)$ , called the  $\epsilon$ -covering number of  $\Psi$ , is the cardinality of the smallest subclass of functions  $\Psi_\epsilon = \{\psi_1, \dots, \psi_{\mathcal{N}_1(\epsilon, \Psi, \mathbb{D}_n)} : \mathbb{R}^{d+p} \rightarrow [-B, B]\}$ , with the property that for fixed points  $z_1, \dots, z_n$  and for each  $\psi \in \Psi$  there is a  $\psi^* \in \Psi_\epsilon$  satisfying  $\frac{1}{n} \sum_{i=1}^n |\psi(z_i) - \psi^*(z_i)| < \epsilon$ . Different extensions and variants of Theorem 1 are given by Alexander (1984), Massart (1990), Talagrand (1994), Giné (1996); also see the monograph by Vapnik (1998). Theorem 1 and its variants can provide tools to establish strongly consistency results for many important statistical estimation problems, including regression functions, density functions, and time series estimation; see, for example, the monograph by van de Geer (2000).

## 1.2 The difficulty with incomplete data

The situation can become quite different and challenging when not every  $\mathbf{Z}_i = (\mathbf{X}_i, \mathbf{V}_i)'$  is fully observable; in particular, some of the  $\mathbf{V}_i$ 's may be missing. Of course, one may decide to estimate  $\nu(\psi) = \mathbb{E}(\psi(\mathbf{Z}))$  based on the complete cases only, where a complete case refers to the fully observable  $\mathbf{Z}_i$  (i.e. when  $\xi_i = 1$ ). In this case the estimator can be expressed as  $\tilde{\nu}(\psi) = \frac{1}{n'} \sum_{i=1}^n \xi_i \psi(\mathbf{Z}_i)$ , where  $n' = \sum_{i=1}^n \xi_i$ . However, there are drawbacks with such estimators: (i) If a large proportion of the data (say 60 to 70 percent) have missing  $\mathbf{V}_i$ 's then, from a practical point of view, it makes sense to somehow revise  $\tilde{\nu}_n$  to take into account the information which is available from  $\mathbf{X}_i$ 's. (ii) There are also theoretical reasons for not using the estimator  $\tilde{\nu}_n$ . For example, this estimator is not in general unbiased for  $\nu(\psi)$  and therefore the corresponding empirical process,  $\{\tilde{\nu}(\psi) - \nu(\psi) | \psi \in \Psi\}$  is not centered in general (not even asymptotically), and this plays a crucial role in establishing the theoretical properties of our propose density and regression estimators in this paper.

## 1.3 Summary of main results

Our main contributions may be summarized as follows.

1. We present revised version of  $\tilde{\nu}(\psi)$  that take into account the missing covariates via an inverse weighting approach, where the weight functions are estimates of the selection probabilities. We propose a plug-in approach to replace these unknown selection probabilities with kernel regression and least-squares estimators
2. Under standard mild assumptions, we derive exponential bounds and inequalities similar to those of Theorem 1 that are suitable for our situation involving incomplete data.

3. We propose new density estimators in the presence of missing data. Furthermore, the new exponential bounds and inequalities in 2 above will be used to study the convergence properties of the proposed density estimator with respect to the Hellinger distance. Our methodology will also be applied to construct regression function estimators that are strongly optimal in the  $L_2$  sense.

## 2 Main results

In this section we introduce the notion of missingness at random and suggest estimators of *selection probabilities* which will be used to construct our density and regression function estimators. We first establish a counterpart of Theorem 1, corresponding to the case of incomplete data: this result will help us to study our density and regression estimators later in this section.

### 2.1 Revised estimation of $\nu(\psi) := \mathbb{E}(\psi(\mathbf{Z}))$

The function  $\mathbb{P}\{\xi = 1|\mathbf{Z}\}$ , called the selection probability, plays an important role in estimation theory with incomplete data. In practice, this function is usually unknown and must be estimated. Under the commonly used assumption of data *Missing At Random* (MAR), it is assumed that the selection probability does not depend on  $\mathbf{V}$  itself. In other words

$$\mathbb{P}\{\xi = 1|\mathbf{X}, \mathbf{V}\} = \mathbf{P}\{\xi = 1|\mathbf{X}\} (= \mathbf{E}(\xi|\mathbf{X})) \quad (1)$$

This assumption is essentially the baseline of analysis in the literature on incomplete data; see, for example, Cheng (1994), Cheng and Chu (1996), Wang and Rao (2002), Müller (2012), Tang et al (2012), and Bravo (2015). In what follows we shall focus on the case where the MAR assumption (1) holds. Define

$$\hat{\nu}(\psi) = n^{-1} \sum_{i=1}^n \xi_i \psi(\mathbf{Z}_i) / \hat{\pi}(\mathbf{X}_i) \quad (2)$$

The estimator of  $\nu(\psi)$ , where  $\hat{\pi}(\mathbf{X}_i)$  is an estimator of the selection probability

$$\pi(\mathbf{X}_i) := \mathbb{E}(\xi_i|\mathbf{X}_i) = \mathbb{P}\{\xi_i = 1|\mathbf{X}_i\} \quad (3)$$

The estimator  $\hat{\nu}(\psi)$  in (2) is in the spirit of the classical Horvitz-Thompson estimator (Horvitz and Thompson (1952)) in the sense that it works by weighting the complete cases by the inverse of the estimates of the selection probabilities,  $\pi(\mathbf{X}_i)$ . In fact, this approach has been used by many authors in the literature; see, for example, Robins et al. (1994) who propose a class of semiparametric estimators of regression coefficients based on inverse probability weighting estimating equations, with the imposed assumption that the missing probabilities are either known or can be modeled parametrically. Robins and Rotnitzky (1995) study efficient estimation in semiparametric multivariate regression models with missing response variables. Hirano et al. (2003) propose an estimator of the average treatment effect of a binary treatment using nonparametric methods for the missing probabilities. Wang et al. (2010) propose a class of kernel estimating equations to estimate a function  $\theta(\mathbf{x}) = g\{E(Y|\mathbf{X} = \mathbf{x})\}$  whenever one has access to some auxiliary covariate vector  $\mathbf{U}$ . Here  $g$  is a known link function. Their method works by weighting the units with complete data by either the inverse of the true selection probability  $\pi_i \equiv \pi(\mathbf{X}_i, \mathbf{U}_i) := P\{\zeta_i = 1|\mathbf{X}_i, \mathbf{U}_i\}$ , or an estimator of it, where  $\pi_i$  is assumed to have a known functional form. Unlike the above result, in this paper we do not assume the availability of any additional auxiliary covariates. Furthermore, we do not impose assumptions that the functional form of the missing probabilities must always be known or be such that the functional form of the missing probabilities must always be known or be such that it can be modeled parametrically. Of course, if such assumptions hold then our proposed least squares methods will produce more accurate results, but our approach can also tackle the nonparametric case (via kernel methods) where missing probabilities are completely unknown. In this paper, we also derive exponential bounds on the performance of the proposed estimators. Unlike the above authors, here we carry out a Horvitz-Thompson type inversely weighted least squares

criterion to estimate the underlying regression function, where the weights are the inverses of the estimated selection probabilities of the complete cases. In the case of density estimation with incomplete data, our results are quite new and may be viewed as the counterparts of the classical work of van de Geer (1993, 200) who establishes strong Hellinger consistency of maximum likelihood density estimation.

As for the estimator of  $\widehat{\pi}(\mathbf{X}_i)$  in (2), we consider two choices: (a) kernel regression and (b) the least squares.

(a) *The kernel regression estimation of  $\pi(\mathbf{X}_i) = \mathbb{E}(\xi_i|\mathbf{X}_i)$*

Our first estimator is given by the kernel regression estimator

$$\widehat{\pi}_{\text{ker}}(\mathbf{X}_i) = \frac{\sum_{k=1, \neq i}^n \xi_k \mathcal{K}\left(\frac{\mathbf{X}_i - \mathbf{X}_k}{h_n}\right)}{\sum_{k=1, \neq i}^n \mathcal{K}\left(\frac{\mathbf{X}_i - \mathbf{X}_k}{h_n}\right)}, \quad (4)$$

with the convention  $0/0 = 0$ , where  $\mathcal{K} : \mathbb{R}^d \rightarrow \mathbb{R}_+$  is the kernel used, and  $h_n$  is the smoothing parameter of the kernel ( $h_n \rightarrow 0$ , as  $n \rightarrow \infty$ ). Replacing  $\widehat{\pi}(\cdot)$  in (2) by  $\widehat{\pi}_{\text{ker}}(\cdot)$ , we find the following estimator of  $\nu(\psi) = \mathbb{E}(\psi(\mathbf{Z}))$

$$\widehat{\nu}_n^{(\text{ker})}(\psi) = n^{-1} \sum_{i=1}^n \xi_i \psi(\mathbf{Z}_i) / \widehat{\pi}_{\text{ker}}(\mathbf{X}_i). \quad (5)$$

To assess the performance of  $\widehat{\nu}_n^{(\text{ker})}(\psi)$ , we first state a number of assumptions.

(A1)  $\pi_0 := \inf_{\mathbf{x} \in \mathbb{R}^d} \mathbb{P}\{\xi = 1 | \mathbf{X} = \mathbf{x}\} > 0$ .

(A2) The kernel  $\mathcal{K}$  in (4) satisfies  $\int_{\mathbb{R}^d} \mathcal{K}(\mathbf{x}) d\mathbf{x} = 1$  and  $\int_{\mathbb{R}^d} |x_i| \mathcal{K}(\mathbf{x}) d\mathbf{x} < \infty$ , for  $i = 1, \dots, d$ , where  $\mathbf{x} = (x_1, \dots, x_d)'$ . The smoothing parameter,  $h_n$ , of the kernel satisfies  $h_n \rightarrow 0$  and  $nh_n^d \rightarrow \infty$ , as  $n \rightarrow \infty$ .

(A3) The random vector  $\mathbf{X}$  has a compactly supported probability density function (pdf),  $f(\mathbf{x})$ , which is bounded away from zero on its compact support, i.e.,  $f_0 := \inf_{\mathbf{x}} f(\mathbf{x}) > 0$ . Furthermore,  $f$  and its first-order partial derivatives are uniformly bounded.

(A4) The partial derivatives  $\frac{\partial}{\partial x_i} \pi(\mathbf{x})$  exist for  $i = 1, \dots, d$  and are bounded uniformly, in  $\mathbf{x}$ , on the compact support of  $f$ , where  $\pi(\mathbf{x}) = \mathbb{E}(\xi | \mathbf{X} = \mathbf{x})$ .

The result below is a counterpart of Theorem 1 and is suitable for our incomplete data setup.

**Theorem 2** *Let  $\Psi$  be a class of functions  $\psi : \mathbb{R}^{d+p} \rightarrow [-B, B]$ ,  $0 < B < \infty$ , and let  $\widehat{\nu}_n^{(\text{ker})}(\psi)$  be as in (5). Then under assumptions (A1)–(A4), for every  $\epsilon > 0$  there is a  $n_0 > 0$  such that for all  $n > n_0$*

$$\mathbb{P} \left\{ \sup_{\psi \in \Psi} |\widehat{\nu}_n^{(\text{ker})}(\psi) - \nu(\psi)| > \epsilon \right\} \leq 6n e^{-nh^d C_1 \epsilon^2} + 8\mathbb{E} \left[ \mathcal{N}_1 \left( \frac{\pi_0 \epsilon}{16}, \Psi, \mathbb{D}_n \right) \right] e^{-nC_2 \epsilon^2},$$

where  $C_1$  and  $C_2$ , are positive constants not depending on  $n$  or  $\epsilon$ .

The proof of the above result appears in the Appendix. In passing we also note that Theorem 2 in conjunction with the Borel-Cantelli lemma imply that if  $(nh_n^d)^{-1} \log n \rightarrow 0$  and  $n^{-1} \log \mathbb{E} [\mathcal{N}_1(\pi_0 \epsilon / 16, \Psi, \mathbb{D}_n)] \rightarrow 0$ , as  $n \rightarrow \infty$  and  $h_n \rightarrow 0$ , then  $\sup_{\psi \in \Psi} |\widehat{\nu}_n^{(\text{ker})}(\psi) - \nu(\psi)| \xrightarrow{\text{a.s.}} 0$ .

Next, we consider the least squares estimator.

(b) *The least-squares estimator of  $\pi(\mathbf{X}_i) = \mathbb{E}(\xi_i|\mathbf{X}_i)$*

Let  $\mathcal{P}$  be a known class of functions of the form  $\tilde{\pi} : \mathbb{R}^d \rightarrow [\pi_0, 1]$ , where  $\pi_0$  is as in assumption (A1). The least-squares estimator of  $\pi$  is

$$\hat{\pi}_{\text{LS}} = \operatorname{argmin}_{\tilde{\pi} \in \mathcal{P}} \frac{1}{n} \sum_{i=1}^n (\xi_i - \tilde{\pi}(\mathbf{X}_i))^2. \quad (6)$$

Therefore, upon taking  $\hat{\pi}$  to be  $\hat{\pi}_{\text{LS}}$  in (2), we find the following estimator of  $\nu(\psi) = \mathbb{E}(\psi(\mathbf{Z}))$

$$\hat{\nu}^{(\text{LS})}(\psi) = \frac{1}{n} \sum_{i=1}^n \frac{\xi_i \psi(\mathbf{Z}_i)}{\hat{\pi}_{\text{LS}}(\mathbf{X}_i)}. \quad (7)$$

The following result is a version of Theorem 2, corresponding to the least square estimator.

**Theorem 3** *Let  $\Psi$  be a class of functions  $\psi : \mathbb{R}^{d+p} \rightarrow [-B, B]$ ,  $0 < B < \infty$  and let  $\hat{\nu}_n^{(\text{LS})}(\psi)$  be as in (7). Suppose that assumption (A1) holds and that  $\pi \in \mathcal{P}$ . Then for every  $\epsilon > 0$ , there is a  $n_0 > 0$  such that for all  $n > n_0$*

$$\begin{aligned} \mathbb{P} \left\{ \sup_{\psi \in \Psi} |\hat{\nu}_n^{(\text{LS})}(\psi) - \nu(\psi)| > \epsilon \right\} &\leq 8 \mathbb{E} [\mathcal{N}_1(C_{31}\epsilon, \Psi, \mathbb{D}_n)] e^{-nC_{32}\epsilon^2} \\ &\quad + 8 \mathbb{E} [\mathcal{N}_1(C_{33}\epsilon, \mathcal{P}, (\mathbf{X}_i)_{i=1}^n)] e^{-nC_{34}\epsilon^2} \\ &\quad + 8 \mathbb{E} [\mathcal{N}_1(C_{35}\epsilon^2, \mathcal{P}, (\mathbf{X}_i)_{i=1}^n)] e^{-nC_{36}\epsilon^4}, \end{aligned}$$

where  $C_{31}$  to  $C_{36}$  are positive constants not depending on  $n$  or  $\epsilon$ .

The proof of this theorem appears in the Appendix.

In the next two sections we introduce our density and regression function estimators based on the approach and results of this section.

## 2.2 Density estimation

Once again let  $\mathbf{Z} = (\mathbf{X}', \mathbf{V}')'$ , where  $\mathbf{X} \in \mathbb{R}^d$  is always observable but  $\mathbf{V} \in \mathbb{R}^p$  may be missing at random in the sense that  $\mathbb{P}\{\xi = 1 | \mathbf{Z}\} \stackrel{\text{MAR}}{=} \mathbb{P}\{\xi = 1 | \mathbf{X}\} =: \pi(\mathbf{X})$ . Here, as before,  $\xi = 0$  if  $\mathbf{V}$  is missing (and  $\xi = 1$ , otherwise). We are interested in estimating the probability density function of  $\mathbf{Z}$ , based on the data  $\mathbb{D}_n = \{(\mathbf{Z}_1, \xi_1), \dots, (\mathbf{Z}_n, \xi_n)\} = \{(\mathbf{X}_1, \mathbf{V}_1, \xi_1), \dots, (\mathbf{X}_n, \mathbf{V}_n, \xi_n)\}$ . When  $(X, V) \in \mathbb{R}^2$ , Hazelton (2000) constructs a kernel density estimator of the marginal distribution of  $V$ , using the  $X$ 's as the auxiliary variables. This estimator, which is shown to be strongly uniformly consistent, is based on the distribution function estimator of Cheng and Chu (1996). A more recent result along these lines is the Horvitz-Thomson inverse weighting type density estimator of Dubnicka (2009) for the random variable  $V$ , based on the availability of an auxiliary random variable  $X$ . Our approach here, which does not assume the availability of any auxiliary variables, works as follows.

Suppose that the true pdf,  $g_0$ , of  $\mathbf{Z}$  belongs to a class of densities  $\mathcal{G}$ . Clearly, when there are no missing data, the classical maximum likelihood estimator (MLE) of  $g_0$  is  $\hat{g}_n = \operatorname{argmax}_{g \in \mathcal{G}} \sum_{i=1}^n \log g(\mathbf{Z}_i)$ . Now, to tackle the presence of missing data define

$$\hat{L}_n(g) = n^{-1} \sum_{i=1}^n \frac{\xi_i \log g(\mathbf{Z}_i)}{\hat{\pi}(\mathbf{X}_i)}, \quad g \in \mathcal{G}, \quad (8)$$

where  $\hat{\pi}$  can be either  $\hat{\pi}_{\text{ker}}$  as defined in (4) or  $\hat{\pi}_{\text{LS}}$  as defined in (6). We consider the following estimator

$$\hat{g}_n = \operatorname{argmax}_{g \in \mathcal{G}} \hat{L}_n(g).$$

How good is this MLE-type density estimator? To answer this question we first recall that the Hellinger distance  $\rho_H$  between two densities  $g_1, g_2 \in \mathcal{G}$  is given by

$$\rho_H(g_1, g_2) = \sqrt{\frac{1}{2} \int \left[ \sqrt{g_1(\mathbf{z})} - \sqrt{g_2(\mathbf{z})} \right]^2 d\mathbf{z}}.$$

Hellinger-consistency of certain nonparametric density estimators have been studied by van de Geer (1993). For more on  $\rho_H$  and some of its properties see, for example, van de Geer (2000, Ch.4). The following result gives exponential performance bounds on the distance (Hellinger) between  $\hat{g}_n$  and the target density  $g_0$ .

**Theorem 4** *Let  $\rho_H(\hat{g}_n, g_0)$  be the Hellinger distance between  $\hat{g}_n$  and  $g_0$ , and suppose that assumptions (A1)-(A4) hold. Then for every  $\epsilon > 0$  there is an  $n_0 > 0$  such that for all  $n > n_0$*

$$\mathbb{P} \left\{ \rho_H^2(\hat{g}_n, g_0) > \epsilon \right\} \leq 8\mathbb{E} \left[ \mathcal{N}_1 \left( C_{16}\epsilon, \frac{\mathcal{G}}{g_0}, \mathbb{D}_n \right) \right] e^{-nC_{17}\epsilon^2} + \tau_n(\epsilon),$$

in which

$$\tau_n(\epsilon) = \begin{cases} 6n e^{-nh^d C_{18}\epsilon^2} & \text{if using } \hat{\pi}_{\ker} \text{ in (8),} \\ 8\mathbb{E} [\mathcal{N}_1(C_{19}\epsilon, \mathcal{Q}, (\mathbf{X}_i)_{i=1}^n)] e^{-nC_{20}\epsilon^2} \\ \quad + 8\mathbb{E} [\mathcal{N}_1(C_{21}\epsilon^2, \mathcal{Q}, (\mathbf{X}_i)_{i=1}^n)] e^{-nC_{22}\epsilon^4} & \text{if using } \hat{\pi}_{LS} \text{ in (8),} \end{cases}$$

where  $C_{16} - C_{22}$  are positive constants not depending on  $n$  or  $\epsilon$  and

$$\frac{\mathcal{G}}{g_0} = \left\{ \frac{g}{g_0} \mathbb{I}\{g_0 > 0\} \mid g \in \mathcal{G} \right\}.$$

The above result can be used to establish almost-sure convergence results for  $\hat{g}_n$  with respect to the Hellinger distance. For example, with  $\hat{\pi} = \hat{\pi}_{\ker}$ , if  $(nh_n^d)^{-1} \log n \rightarrow 0$  and  $n^{-1} \log \mathbb{E} [\mathcal{N}_1(C_{16}\epsilon, \mathcal{G}/g_0, \mathbb{D}_n)] \rightarrow 0$ , as  $n \rightarrow \infty$  and  $h_n \rightarrow 0$ , then by the Borel-Cantelli lemma one has  $\rho_H(\hat{g}_n, g_0) \xrightarrow{a.s.} 0$ .

#### PROOF OF THEOREM 4

We prove the theorem for the case where  $\hat{\pi} = \hat{\pi}_{\ker}$ . The proof for the case where  $\hat{\pi} = \hat{\pi}_{LS}$  is almost identical and will not be given. Define the quantities

$$\begin{aligned} \bar{g}_n &= \frac{\hat{g}_n + g_0}{2}, \quad \bar{g} = \frac{g + g_0}{2}, \quad g \in \mathcal{G} \\ L_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) &= \mathbb{E} \left[ \log \frac{\bar{g}_n(\mathbf{Z})}{g_0(\mathbf{Z})} \mathbb{I}\{g_0(\mathbf{Z}) > 0\} \mid \mathbb{D}_n \right] \\ L \left( \frac{\bar{g}}{g_0} \mathbb{I}\{g_0 > 0\} \right) &= \mathbb{E} \left[ \log \frac{\bar{g}(\mathbf{Z})}{g_0(\mathbf{Z})} \mathbb{I}\{g_0(\mathbf{Z}) > 0\} \right]. \end{aligned}$$

We first show that

$$\rho_H^2(\bar{g}_n, g_0) \leq \frac{1}{2} \left\{ \hat{L}_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) - L_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) \right\}, \quad (9)$$

where  $\hat{L}_n(\cdot)$  is as in (8). Here, (9) may be viewed as a version of Lemma 4.1 of van de Geer (2000) tailored to fit our current situation where  $\mathbf{Z}_i$ 's are allowed to have missing components; in fact, to prove (9), we borrow the arguments used in the proof of the cited result. First note that by the definition of  $\hat{g}_n$  we have

$$\forall 0 < g \in \mathcal{G} : \hat{L}_n \left( \frac{\hat{g}_n}{g} \right) = \hat{L}_n(\hat{g}_n) - \hat{L}_n(g) \geq 0$$

Also, by the concavity of the logarithmic function

$$\frac{1}{2} \log \frac{\hat{g}_n(\mathbf{z})}{g_0(\mathbf{z})} \mathbb{I}\{g_0(\mathbf{z}) > 0\} \leq \log \frac{\bar{g}_n(\mathbf{z})}{g_0(\mathbf{z})} \mathbb{I}\{g_0(\mathbf{z}) > 0\}, \quad \text{where } \bar{g}_n = \frac{\hat{g}_n + g_0}{2}.$$

Thus

$$\begin{aligned} 0 &\leq \widehat{L}_n \left( \frac{\widehat{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) \leq 2 \widehat{L}_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) \\ &= 2 \left[ \widehat{L}_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) - L_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) \right] + 2 L_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right). \end{aligned} \quad (10)$$

But  $\bar{g}_n = (\widehat{g}_n + g_0)/2$  is also a density and therefore by Lemma 1.3 of van de Geer (2000)

$$L_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) \leq -2 \rho_{\mathbb{H}}^2(\bar{g}_n, g_0).$$

This last inequality together with (10) imply (9). Therefore

$$\begin{aligned} \rho_{\mathbb{H}}^2(\widehat{g}_n, g_0) &\leq 16 \rho_{\mathbb{H}}^2(\bar{g}_n, g_0) \\ &\leq 8 \left[ \widehat{L}_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) - L_n \left( \frac{\bar{g}_n}{g_0} \mathbb{I}\{g_0 > 0\} \right) \right], \quad (\text{by (9)}) \\ &\leq 8 \sup_{g \in \mathcal{G}} \left| \widehat{L}_n \left( \frac{\bar{g}}{g_0} \mathbb{I}\{g_0 > 0\} \right) - L \left( \frac{\bar{g}}{g_0} \mathbb{I}\{g_0 > 0\} \right) \right|, \quad (\text{where } \bar{g} = \frac{g+g_0}{2}). \end{aligned}$$

Next, define the class of functions

$$\mathcal{M} = \left\{ m : \mathbb{R}^{d+p} \rightarrow \mathbb{R} \mid m(\mathbf{z}) = \log \frac{g(\mathbf{z}) + g_0(\mathbf{z})}{2g_0(\mathbf{z})} \mathbb{I}\{g_0(\mathbf{z}) > 0\}, g \in \mathcal{G} \right\}$$

and note that, by Theorem 2, for every  $\epsilon > 0$

$$\begin{aligned} \mathbb{P} \left\{ \rho_{\mathbb{H}}^2(\widehat{g}_n, g_0) > \epsilon \right\} &\leq \mathbb{P} \left\{ \sup_{g \in \mathcal{G}} \left| \widehat{L}_n \left( \frac{\bar{g}}{g_0} \mathbb{I}\{g_0 > 0\} \right) - L \left( \frac{\bar{g}}{g_0} \mathbb{I}\{g_0 > 0\} \right) \right| > \frac{\epsilon}{8} \right\} \\ &= \mathbb{P} \left\{ \sup_{m \in \mathcal{M}} \left| \widehat{L}_n(m) - L(m) \right| > \frac{\epsilon}{8} \right\} \\ &\leq 6n e^{-nh_n^d C_{14} \epsilon^2} + 8 \mathbb{E} \left[ \mathcal{N}_1 \left( \frac{\pi_0 \epsilon}{128}, \mathcal{M}, \mathbb{D}_n \right) \right] e^{-nC_{15} \epsilon^2}, \end{aligned}$$

for  $n$  large enough, where  $C_{14}$  and  $C_{15}$  are positive constants not depending on  $n$  or  $\epsilon$ . Now, let  $m_1, m_2 \in \mathcal{M}$ , where  $m_k = \log \frac{g_k + g_0}{2g_0} \mathbb{I}\{g_0 > 0\}$ ,  $g_k \in \mathcal{G}$ ,  $k = 1, 2$ , and observe that

$$\begin{aligned} |m_1(\mathbf{z}) - m_2(\mathbf{z})| &= \left| \log \frac{g_1(\mathbf{z}) + g_0(\mathbf{z})}{g_0(\mathbf{z})} - \log \frac{g_2(\mathbf{z}) + g_0(\mathbf{z})}{g_0(\mathbf{z})} \right| \mathbb{I}\{g_0(\mathbf{z}) > 0\} \\ &\leq \left| \frac{g_1(\mathbf{z})}{g_0(\mathbf{z})} - \frac{g_2(\mathbf{z})}{g_0(\mathbf{z})} \right| \mathbb{I}\{g_0(\mathbf{z}) > 0\}, \quad (\text{see van de Geer (2000, p.50)}). \end{aligned}$$

Therefore,  $n^{-1} \sum_{i=1}^n \frac{\xi_i}{\pi(\mathbf{X}_i)} |m_1(\mathbf{Z}_i) - m_2(\mathbf{Z}_i)| \leq n^{-1} \sum_{i=1}^n \frac{\xi_i}{\pi(\mathbf{X}_i)} \left| \frac{g_1(\mathbf{Z}_i)}{g_0(\mathbf{Z}_i)} - \frac{g_2(\mathbf{Z}_i)}{g_0(\mathbf{Z}_i)} \right| \mathbb{I}\{g_0(\mathbf{Z}_i) > 0\}$ , which implies that if  $\left\{ \frac{g_1}{g_0} \mathbb{I}\{g_0 > 0\}, \dots, \frac{g_N}{g_0} \mathbb{I}\{g_0 > 0\} \right\}$  is a minimal  $\epsilon$ -cover of  $\frac{\mathcal{G}}{g_0}$ , then  $\{m_1, \dots, m_N\}$  is an  $\epsilon$ -cover of  $\mathcal{M}$ , where  $m_k = \log \frac{g_k + g_0}{2g_0} \mathbb{I}\{g_0 > 0\}$ . Thus  $\mathcal{N}_1(\epsilon, \mathcal{M}, \mathbb{D}_n) \leq \mathcal{N}_1(\epsilon, \frac{\mathcal{G}}{g_0}, \mathbb{D}_n)$ , and this completes the proof of Theorem 4.  $\square$

## 2.3 Regression function estimation

Let  $(\mathbf{Z}, Y)$  be an  $\mathbb{R}^{d+p} \times \mathbb{R}$ -valued random vector with an unknown distribution. Here  $\mathbf{Z} = (\mathbf{X}', \mathbf{V}')'$ , where  $\mathbf{X} \in \mathbb{R}^d$  is always observable but  $\mathbf{V} \in \mathbb{R}^p$  may be missing at random in the sense that

$$\mathbb{P}\{\xi = 1 | \mathbf{Z}, Y\} \stackrel{\text{MAR}}{=} \mathbb{P}\{\xi = 1 | \mathbf{X}, Y\} =: \pi(\mathbf{X}, Y).$$



We note here that  $\pi$  is a map of the form  $\mathbb{R}^{d+1} \rightarrow [0, 1]$ . Next, let  $\mathbb{D}_n = \{(\mathbf{X}_1, \mathbf{V}_1, Y_1, \xi_1), \dots, (\mathbf{X}_n, \mathbf{V}_n, Y_n, \xi_n)\}$  represent the data (iid) and consider the following kernel regression estimator of  $\pi(\mathbf{X}_i, Y_i) = \mathbb{P}\{\xi_i = 1 | \mathbf{X}_i, Y_i\}$

$$\widehat{\pi}_{\text{ker}}(\mathbf{X}_i, Y_i) = \sum_{j=1, \neq i}^n \xi_j \mathcal{H}((\mathbf{U}_i - \mathbf{U}_j)/\lambda_n) / \sum_{j=1, \neq i}^n \mathcal{H}((\mathbf{U}_i - \mathbf{U}_j)/\lambda_n), \text{ where } \mathbf{U}_j = (\mathbf{X}_j', Y_j)', \quad (11)$$

with the convention that  $0/0 = 0$ , where  $\mathcal{H} : \mathbb{R}^{d+1} \rightarrow \mathbb{R}_+$  is the kernel used with the smoothing parameter  $\lambda_n$  satisfying  $\lambda_n \rightarrow 0$  as  $n \rightarrow \infty$ . Alternatively, as a second choice, one may decide to consider the least-squares estimator of  $\pi(\mathbf{X}_i, Y_i)$ . More specifically, let  $\mathcal{Q}$  be a known class of functions of the form  $\tilde{\pi} : \mathbb{R}^{d+1} \rightarrow [\pi_{\min}, 1]$ , where  $\pi_{\min}$  is as in assumption (A1'), given below. Then the least-squares estimator is given by

$$\widehat{\pi}_{\text{LS}}(\mathbf{X}_i, Y_i) = \operatorname{argmin}_{\tilde{\pi} \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n (\xi_i - \tilde{\pi}(\mathbf{X}_i, Y_i))^2. \quad (12)$$

Let  $\Phi$  be a class of candidate regression functions of the form  $\phi : \mathbb{R}^{d+p} \rightarrow \mathbb{R}$ . Also, let  $\phi^*(\mathbf{z}) = \mathbb{E}(Y | \mathbf{Z} = \mathbf{z})$  be the true regression function. Put

$$L(\phi) = \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 \quad \text{and} \quad \widehat{L}_n(\phi) = \frac{1}{n} \sum_{i=1}^n \frac{\xi_i |\phi(\mathbf{Z}_i) - Y_i|^2}{\widehat{\pi}(\mathbf{X}_i, Y_i)}, \quad \forall \phi \in \Phi, \quad (13)$$

where  $\widehat{\pi}$  is either  $\widehat{\pi}_{\text{ker}}$  as defined in (11) or  $\widehat{\pi}_{\text{LS}}$  as defined in (12). Our proposed Horvitz-Thompson-based least squares estimator of the unknown regression function  $\phi^*$  is given by

$$\phi_n = \operatorname{argmin}_{\phi \in \Phi} \widehat{L}_n(\phi). \quad (14)$$

To study the properties of the  $L_2$  error of our estimator  $\phi_n$ , we derive exponential performance bounds on the deviations of the  $L_2$  error of  $\phi_n$  from that of the best member of the class of candidate functions  $\Phi$ , i.e., the quantity

$$\mathbb{E}([\phi_n(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2 | \mathbb{D}_n) - \inf_{\phi \in \Phi} \mathbb{E}[\phi(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2.$$

We first state the following counterparts of Assumptions (A1) - (A4). Let  $\mathbf{U} = (\mathbf{X}', Y)'$ , and  $\mathbf{u} = (\mathbf{x}', y)'$ . Then the following assumption are exactly the same as those in (A1) - (A4), but with  $d$  replaced by  $d + 1$ ,  $\mathcal{K}$  by  $\mathcal{H}$ ,  $h$  with  $\lambda_n$ , and  $\pi_0$  with  $\pi_{\min}$ :

(A1')  $\pi_{\min} := \inf_{\mathbf{u} \in \mathbb{R}^{d+1}} \mathbb{P}\{\xi = 1 | \mathbf{U} = \mathbf{u}\} > 0$ .

(A2') The kernel  $\mathcal{H}$  in (11) satisfies  $\int_{\mathbb{R}^{d+1}} \mathcal{H}(\mathbf{u}) d\mathbf{u} = 1$  and  $\int_{\mathbb{R}^{d+1}} |u_i| \mathcal{H}(\mathbf{u}) d\mathbf{u} < \infty$ , for  $i = 1, \dots, (d + 1)$ . Also, the smoothing parameter  $\lambda_n$  satisfies  $\lambda_n \rightarrow 0$  and  $n\lambda_n^{d+1} \rightarrow \infty$ , as  $n \rightarrow \infty$ .

(A3') The random vector  $\mathbf{U}$  has a compactly supported pdf, which is bounded away from zero on its compact support. Furthermore, the pdf of  $\mathbf{U}$  and its first-order partial derivatives are uniformly bounded.

(A4') The partial derivatives  $\frac{\partial}{\partial u_i} \pi(\mathbf{u})$  exist for  $i = 1, \dots, (d + 1)$  and are bounded uniformly, in  $\mathbf{u}$ , on the compact support of the pdf of  $\mathbf{U}$ .

**Theorem 5** *Let  $\Phi$  be a class of functions  $\phi : \mathbb{R}^{d+p} \rightarrow [-B, B]$ , for some  $B < \infty$ , and let  $\phi_n$  be as in (14). Suppose that  $|Y| \leq A < \infty$ . Then, under assumptions (A1')-(A4'), for every  $\epsilon > 0$  and  $n$  large enough*

$$\begin{aligned} & \mathbb{P} \left\{ \left| \mathbb{E}([\phi_n(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2 | \mathbb{D}_n) - \inf_{\phi \in \Phi} \mathbb{E}[\phi(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2 \right| > \epsilon \right\} \\ & \leq 8 \mathbb{E} \left[ \mathcal{N}_1 \left( \frac{\pi_{\min} \epsilon}{64(A+B)}, \Phi, \mathbb{D}_n \right) \right] e^{-nC_{10}\epsilon^2} + \delta_n(\epsilon), \end{aligned}$$

in which

$$\delta_n(\epsilon) = \begin{cases} 6ne^{-n\lambda_n^{d+1}C_{11}\epsilon^2} & \text{if using } \widehat{\pi}_{\ker} \text{ in (13),} \\ 8\mathbb{E}[\mathcal{N}_1(C_{12}\epsilon, \mathcal{Q}, (\mathbf{X}_i, Y_i)_{i=1}^n)] e^{-nC_{13}\epsilon^2} \\ \quad + 8\mathbb{E}[\mathcal{N}_1(C_{14}\epsilon^2, \mathcal{Q}, (\mathbf{X}_i, Y_i)_{i=1}^n)] e^{-nC_{15}\epsilon^4} & \text{if using } \widehat{\pi}_{LS} \text{ in (13),} \end{cases}$$

where  $C_{10} - C_{15}$  are positive constants not depending on  $n$  or  $\epsilon$ .

We also note that the above result can be used to establish various almost-sure convergence results. For example, suppose that  $\widehat{\pi}_{\ker}$  is used in (13). Now, if  $n^{-1} \log\{\mathbb{E}[\mathcal{N}_1(\pi_{\min}\epsilon/(64(A+B)), \Phi, \mathbb{D}_n)]\} \rightarrow 0$  and  $(n\lambda_n^{d+1})^{-1} \log n \rightarrow 0$ , as  $n \rightarrow \infty$ , then by an application of the Borel-Cantelli lemma  $\mathbb{E}([\phi_n(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2 | \mathbb{D}_n) \xrightarrow{\text{a.s.}} \inf_{\phi \in \Phi} \mathbb{E}[\phi(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2$ . Clearly, if  $\phi^* \in \Phi$  then  $\inf_{\phi \in \Phi} \mathbb{E}[\phi(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2 = 0$ , which yields  $\mathbb{E}([\phi_n(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2 | \mathbb{D}_n) \xrightarrow{\text{a.s.}} 0$ , under the above conditions.

#### PROOF OF THEOREM 5.

We give a proof for the case where  $\widehat{\pi}$  in (13) is the kernel estimator  $\widehat{\pi}_{\ker}$ . The proof for the case where  $\widehat{\pi}_{LS}$  is used is virtually the same and will not be given. First note that using the decomposition  $\mathbb{E}([\phi_n(\mathbf{Z}) - Y]^2 | \mathbb{D}_n) = \mathbb{E}([\phi_n(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2 | \mathbb{D}_n) + \mathbb{E}[\phi^*(\mathbf{Z}) - Y]^2$  we can write

$$\begin{aligned} \mathbb{E}([\phi_n(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2 | \mathbb{D}_n) &= \left\{ \mathbb{E}([\phi_n(\mathbf{Z}) - Y]^2 | \mathbb{D}_n) - \inf_{\phi \in \Phi} \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 \right\} \\ &\quad + \left\{ \inf_{\phi \in \Phi} \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 - \mathbb{E}[\phi^*(\mathbf{Z}) - Y]^2 \right\}. \end{aligned} \quad (15)$$

But

$$\inf_{\phi \in \Phi} \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 - \mathbb{E}[\phi^*(\mathbf{Z}) - Y]^2 = \inf_{\phi \in \Phi} \mathbb{E}[\phi(\mathbf{Z}) - \phi^*(\mathbf{Z})]^2. \quad (16)$$

Furthermore,

$$\begin{aligned} &\mathbb{E}([\phi_n(\mathbf{Z}) - Y]^2 | \mathbb{D}_n) - \inf_{\phi \in \Phi} \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 \\ &= \sup_{\phi \in \Phi} \left\{ \mathbb{E}([\phi_n(\mathbf{Z}) - Y]^2 | \mathbb{D}_n) - \widehat{L}_n(\phi_n) + \widehat{L}_n(\phi_n) - \widehat{L}_n(\phi) + \widehat{L}_n(\phi) - \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 \right\} \\ &\leq \sup_{\phi \in \Phi} \left\{ \mathbb{E}([\phi_n(\mathbf{Z}) - Y]^2 | \mathbb{D}_n) - \widehat{L}_n(\phi_n) + \widehat{L}_n(\phi) - \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 \right\} \\ &\quad (\text{because } \widehat{L}_n(\phi_n) - \widehat{L}_n(\phi) \leq 0, \quad \forall \phi \in \Phi) \\ &\leq \left| \mathbb{E}([\phi_n(\mathbf{Z}) - Y]^2 | \mathbb{D}_n) - \widehat{L}_n(\phi_n) \right| + \sup_{\phi \in \Phi} \left| \widehat{L}_n(\phi) - \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 \right| \\ &\leq 2 \sup_{\phi \in \Phi} \left| \widehat{L}_n(\phi) - \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 \right|, \end{aligned} \quad (17)$$

where  $\widehat{L}_n(\phi)$  is as in (13). Therefore in view of (15), (16), and (17), it is sufficient to show that the bound in Theorem 5 is also a bound on

$$\mathbb{P} \left\{ \sup_{\phi \in \Phi} \left| \widehat{L}_n(\phi) - \mathbb{E}[\phi(\mathbf{Z}) - Y]^2 \right| > \frac{\epsilon}{2} \right\}. \quad (18)$$

Now consider the class of functions of the form  $\varphi(\mathbf{z}, y) = [y - \phi(\mathbf{z})]^2$  indexed by members of  $\Phi$ , i.e., the class of functions

$$\mathcal{F} \equiv \mathcal{F}_{\Phi} = \left\{ \varphi : \mathbb{R}^{d+p} \times \mathbb{R} \longrightarrow [-(A+B), A+B] \mid \varphi(\mathbf{z}, y) = [y - \phi(\mathbf{z})]^2, \phi \in \Phi \right\}.$$

Then, by the definition of  $\widehat{L}_n(\phi)$  in (13), we have

$$(18) = \mathbb{P} \left\{ \sup_{\varphi \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \frac{\xi_i \varphi(\mathbf{Z}_i, Y_i)}{\widehat{\pi}_{\ker}(\mathbf{X}_i, Y_i)} - \mathbb{E}[\varphi(\mathbf{Z}, Y)] \right| > \frac{\epsilon}{2} \right\}. \quad (19)$$

Furthermore, by an application of Theorem 2, for  $n$  large enough

$$(19) \leq 6n e^{-n\lambda_n^{d+1}C_{10}\epsilon^2} + 8\mathbb{E} \left[ \mathcal{N}_1 \left( \frac{\pi_{\min}\epsilon}{32}, \mathcal{F}, \mathbb{D}_n \right) \right] e^{-nC_{11}\epsilon^2},$$

where  $\pi_{\min} = \inf_{\mathbf{x}, y} \pi(\mathbf{x}, y) = \inf_{\mathbf{x}, y} \mathbb{P}(\xi = 1 | \mathbf{X} = \mathbf{x}, Y = y)$ , and where  $C_{10}$  and  $C_{11}$  are positive constants not depending on  $n$  or  $\epsilon$ . To complete the proof, observe that for any  $\varphi_1, \varphi_2 \in \mathcal{F}$

$$\begin{aligned} |\varphi_1(\mathbf{z}, y) - \varphi_2(\mathbf{z}, y)| &= |[y - \phi_1(\mathbf{z})]^2 - [y - \phi_2(\mathbf{z})]^2| \\ &\leq |[y - \phi_1(\mathbf{z})] + [y - \phi_2(\mathbf{z})]| \times |[y - \phi_1(\mathbf{z})] - [y - \phi_2(\mathbf{z})]| \\ &\leq 2(A + B) |\phi_1(\mathbf{z}) - \phi_2(\mathbf{z})|. \end{aligned}$$

Therefore  $n^{-1} \sum_{i=1}^n |\varphi_1(\mathbf{Z}_i, Y_i) - \varphi_2(\mathbf{Z}_i, Y_i)| \leq 2(A + B)n^{-1} \sum_{i=1}^n |\phi_1(\mathbf{Z}_i) - \phi_2(\mathbf{Z}_i)|$ , which implies that if  $\{\phi_1, \dots, \phi_N\}$  is a minimal  $\epsilon/(2(A + B))$ -cover of  $\Phi$  with respect to the empirical  $L_1$  norm, then  $\{\varphi_1, \dots, \varphi_N\}$  is an  $\epsilon$ -cover of  $\mathcal{F}$ . Consequently, for every  $t > 0$  we have  $\mathcal{N}_1(t, \mathcal{F}, \mathbb{D}_n) \leq \mathcal{N}_1(t/(2(A + B)), \Phi, \mathbb{D}_n)$ , and this completes the proof of Theorem 5.  $\square$

### Remarks.

Our proposed kernel estimators of the selection probabilities, as given by (4) or (14) play a crucial role in the development of our proposed density and regression function estimators. This can be noticed from the presence of the function  $\hat{\pi}$  in the denominator of the expressions in (8) and (13). The main issue with kernel type estimators is usually the choice of the bandwidth. In the case of kernel regression estimators, a popular choice of the bandwidth is the one that minimizes the Integrated Squared Error (ISE) of the corresponding kernel regression estimator. However, since ISE depends on the underlying unknown regression and density functions, Härdle and Marron (1985) replace them with "leave-one-out" estimators which are then used to de

ne their cross-validation bandwidth selection rule. A more recent approach is based on the cross-validation method of Racine and Li (2004), which is implemented in the 'R' package called "np" (see Racine and Hay

eld (2008)); in fact, we have used this method in our numerical studies of the next section. In the case of density estimation with missing data, one may also consider choosing the bandwidth as the minimizer of the mean Hellinger distance (MHD) or the mean weighted Hellinger distance (MWHD) proposed and studied by Ibrahim A. Ahmad and A. R. Mugdadi (2006). Unfortunately, the fact that our setup involves missing variables makes it very difficult to study MHD or MWHD analytically here. Alternatively, one may choose the bandwidth (from a grid of values) as the minimizer of an empirical version of MWHD, but we have not pursued that path in this paper.

## 3 Numerical examples

In what follows, we provide a number of numerical examples in order to assess the performance of our proposed estimators.

**Example A.** [*Density estimation.*] Here we consider the performance of the density estimator  $\hat{g}_n$  as the maximizer of  $\hat{L}(g)$  defined by (8), where  $\hat{\pi}$  in (8) can be either  $\hat{\pi}_{ker}$ , defined in (4) or by  $\hat{\pi}_{LS}$ , defined in (6). We denote the corresponding density estimators by  $\hat{g}_{\hat{\pi}_{ker}}$  and  $\hat{g}_{\hat{\pi}_{LS}}$ , respectively. We have also considered the complete case (cc) density estimator, which uses the fully observed data only; this estimator will be denoted by  $\hat{g}_{cc}$ . Next, we carry our the numerical work, we generated  $n$  iid observations,  $\mathbf{Z}_n, 1 = 1, \dots, n$ , from  $d$ -dimensional normal distributions; here we have considered two different samples sizes,  $n = 100$  and  $n = 200$ . Also, we use two different values of  $d$ .  $d = 2$  and  $d = 5$ . For  $d = 2$ , we generated the actual data  $\mathbf{Z}_n = (X_i, V_i, i = 1, \dots, n$ , from a bivariate Gaussian distribution with mean vector  $(0, 1)'$  and covariance matrix  $\Sigma = (\sigma_{if})$ , where  $\sigma_{1,1} = 1, \sigma_{2,2} = 2, \sigma_{1,2} = \sigma_{2,1} = 0$ . Here  $X_i$  is always observable but  $V_i$  may be missing at random based on one of the following two missing

probability models for the function  $\pi$  defined in Here  $X_i$  is always observable but  $V_i$  may be missing at random based on one of the following two missing probability models for the function  $\pi$  defined in (3):

*Model A.*  $\pi(x) := \mathbb{P}\{\zeta = 1|X = x\} = \exp(1 + 0.2x)/1 + \exp(1 + 0.2x)\}$ ,

*Model B.*  $\pi(x) := \mathbb{P}\{\zeta = 1|X = x\} = 0.4[\exp(1 + 0.2x)/1 + \cos\{\exp(x) + \cos(x \sin(x))\}]$ .

For  $d = 5$ , the data were generated from a 5-dimensional Gaussian distribution with vector  $(1, 1, 1, 1, 1)'$  and the covariance matrix  $\Sigma = (\sigma_{ij})$ , where  $\sigma_{j,k} = 2^{-|j-k|}$ ,  $1 \leq j, k \leq 5$ . Here, the last four components of  $\mathbf{V}_i$  are allowed to be missing: more specifically, writing  $\mathbf{Z}_i = (X_i, V_{i1}, V_{i2}, V_{i3}, V_{i4})'$ , the component  $X_i$  is always observable, but  $(V_{i1}, V_{i2}, V_{i3}, V_{i4})$  may be missing according to one of the following two missing probability models:

*Model C.*  $\pi(x) := \mathbb{P}\{\zeta = 1|X = x\} = \exp(1 - 0.2x)/1 + \exp(1 - 0.2x)\}$ ,

*Model D.*  $\pi(x) := \mathbb{P}\{\zeta = 1|X = x\} = 0.5[1 + \cos(2x \sin(x))]$ .

The above choices of the missing probability mechanism result, roughly, in 30% missing data for Mode A and about 50% missing data for models B, C, and D. Let  $g$  be the true probability density function of  $\mathbf{Z}$ . In order to and the density estimators  $\hat{g}_{\hat{\pi}_{ker}}$  and  $\hat{g}_{\hat{\pi}_{LS}}$ , we first used the cross-validation method of Racine and Li (2004) in the R package called "np" (see Racine and Hayfield 2008) to find the kernel estimator  $\hat{\pi}_{ker}$ , defined in (4). As for the parameters of the logistic missing probability mechanism  $\hat{\pi}_{LS}$ , defined in (6), we used nonlinear least squares regression (based on the R package "nls2"). To assess the performance of the estimators  $\hat{g}_{\hat{\pi}_{ker}}$ ,  $\hat{g}_{\hat{\pi}_{LS}}$ , and the complete case estimator  $\hat{g}_{cc}$ , we computed the Hellinger distance between each estimator and the true density  $g$  based on two different sample sizes,  $n = 100$  and  $n = 200$ . As a point of reference, we have also included the usual maximum likelihood estimator of  $g$  based on the full data of size  $n$  (i.e., when there are no missing data); this estimator is denoted by  $\tilde{g}$ . The entire above process was repeated a total of 500 times, each time using a sample of size  $n$ , and the average Hellinger distance were computed. The results for the case  $d = 2$ , which correspond to models A and B, appear in the first two rows of Table 1. The numbers appearing in brackets are the standard errors over 500 Monte Carlo runs. Table 1 shows both  $\hat{g}_{\pi_{ker}}$  and  $\hat{g}_{\pi_{LS}}$  tend to outperform the complete case estimator  $\hat{g}_{cc}$  for models A and B. Under Model A, the estimator  $\hat{g}_{\pi_{LS}}$  is slightly superior to the kernel based estimator  $\hat{g}_{\pi_{ker}}$  which is not surprising because we are assuming that we know that the true underlying missing probability mechanism follows a logistic model. Similarly, under Model B, the kernel estimator does a better job in estimating the highly nonlinear trigonometric function  $\pi(x)$  than the least squares method which is still assumes a logistic model. The estimator  $\tilde{g}$ , which is based on no missing data, is included only as a point of reference.

Table 1: Average Hellinger error, over 500 Monte Carlo runs, corresponding to the four models A, B, C, and D, for the proposed density estimators,  $\hat{g}_{\hat{\pi}_{ker}}$  and  $\hat{g}_{\hat{\pi}_{LS}}$ . Here  $\hat{g}_{cc}$  is the complete case estimator (that discards all incomplete data points), and  $\tilde{g}$  represents the estimator based on no missing data (it is included as a point of reference). The numbers in brackets are the standard errors.

	n=100				n=200			
	$\hat{g}_{\hat{\pi}_{ker}}$	$\hat{g}_{\hat{\pi}_{LS}}$	$\hat{g}_{cc}$	$\tilde{g}$	$\hat{g}_{\hat{\pi}_{ker}}$	$\hat{g}_{\hat{\pi}_{LS}}$	$\hat{g}_{cc}$	$\tilde{g}$
Model A	.788 (.00127)	.0762 (.00128)	.0810 (.00133)	.0622 (.00108)	.544 (.00085)	.0538 (.00083)	.0585 (.00093)	.0465 (.00074)
Model B	.1124 (.00183)	.1159 (.00204)	.1227 (.00220)	.0662 (.00108)	.0832 (.00131)	.0904 (.00150)	.0981 (.00164)	.0465 (.00074)
Model C	.1061 (.00161)	.1031 (.00159)	.1124 (.00171)	.0755 (.00111)	.0736 (.00106)	.0715 (.00105)	.0812 (.00119)	.0527 (.00079)
Model D	.1508 (.00191)	.2388 (.00719)	.1854 (.00148)	.0755 (.00111)	.1208 (.00145)	.2082 (.00664)	.1729 (.00106)	.0527 (.00079)

The results for the case  $d = 5$ , which corresponds to models C and D, appear in the last two rows of

Table 2: Average  $L_2$  error, over 500 Monte Carlo runs, corresponding to the four models A, B, C, and D, for the proposed density estimators,  $\hat{g}_{\hat{\pi}_{ker}}$  and  $\hat{g}_{\hat{\pi}_{LS}}$ . Here  $\hat{g}_{cc}$  is the complete case estimator (that discards all incomplete data points), and  $\tilde{g}$  represents the estimator based on no missing data (it is included as a point of reference). The numbers in brackets are the standard errors.

	n=100				n=200			
	$\hat{g}_{\hat{\pi}_{ker}}$	$\hat{g}_{\hat{\pi}_{LS}}$	$\hat{g}_{cc}$	$\tilde{g}$	$\hat{g}_{\hat{\pi}_{ker}}$	$\hat{g}_{\hat{\pi}_{LS}}$	$\hat{g}_{cc}$	$\tilde{g}$
Model A	.0349 (.00061)	.0338 (.00060)	.0364 (.00063)	.0295 (.00051)	.0241 (.00039)	.0237 (.00037)	.0261 (.00043)	.0206 (.00034)
Model B	.0519 (.00094)	.0533 (.00104)	.0565 (.00111)	.0295 (.00051)	.0377 (.00063)	.0408 (.00072)	.0445 (.00080)	.0206 (.00034)
Model C	.1076 (.00083)	.1060 (.00083)	.1107 (.00085)	.0909 (.00069)	.0898 (.00066)	.0885 (.00066)	.0942 (.00072)	.0759 (.00058)
Model D	.1286 (.00084)	.1540 (.00205)	.1434 (.00058)	.0909 (.00069)	.1153 (.00074)	.1433 (.00204)	.1387 (.00042)	.0759 (.00058)

Table 3: Average  $L_1$  error, over 500 Monte Carlo runs, corresponding to the four models A, B, C, and D, for the proposed density estimators,  $\hat{g}_{\hat{\pi}_{ker}}$  and  $\hat{g}_{\hat{\pi}_{LS}}$ . Here  $\hat{g}_{cc}$  is the complete case estimator (that discards all incomplete data points), and  $\tilde{g}$  represents the estimator based on no missing data (it is included as a point of reference). The numbers in brackets are the standard errors.

	n=100				n=200			
	$\hat{g}_{\hat{\pi}_{ker}}$	$\hat{g}_{\hat{\pi}_{LS}}$	$\hat{g}_{cc}$	$\tilde{g}$	$\hat{g}_{\hat{\pi}_{ker}}$	$\hat{g}_{\hat{\pi}_{LS}}$	$\hat{g}_{cc}$	$\tilde{g}$
Model A	.1603 (.00265)	.1561 (.00264)	.1674 (.00280)	.1361 (.00226)	.1122 (.00180)	.1105 (.00171)	.1214 (.00201)	.0958 (.00153)
Model B	.2290 (.00374)	.2336 (.00410)	.2472 (.00438)	.1361 (.00226)	.1686 (.00265)	.1809 (.00298)	.1960 (.00325)	.0958 (.00153)
Model C	.2392 (.00362)	.2323 (.00358)	.2533 (.00384)	.1705 (.00249)	.1660 (.00238)	.1612 (.00236)	.1832 (.00269)	.1186 (.00178)
Model D	.3396 (.00428)	.5313 (.01555)	.4170 (.00332)	.1705 (.00249)	.2723 (.00333)	.4657 (.01482)	.3892 (.00238)	.1186 (.00178)

Table 1. As described earlier, the data were generated from a 5-dimensional Gaussian distribution with a unit mean vector and covariance  $\Sigma = (\sigma_{jk})$  are known to be of the form  $\sigma_{j,k} = 2^{-|j-k|}$ ,  $1 \leq j, k \leq 5$ . Once again Table 1 shows that with the logistic missing probability mechanism of Model C, the estimator  $\hat{g}_{\hat{\pi}_{LS}}$  is superior to  $\hat{g}_{\hat{\pi}_{ker}}$  and that both of these estimators outperform the complete case estimator  $\hat{g}_{cc}$ . The situation is completely reversed for Model D (the fourth row of Table 1), where  $\hat{g}_{\hat{\pi}_{ker}}$  is by fare the best among the proposed estimators. In fact, we also note that here  $\hat{g}_{\hat{\pi}_{LS}}$  is a very poor estimator with a rather large standard deviation. Some of these facts are reflected in the boxplots of the 500 Hellinger errors (ie., distances) that appear in the first row of Figure 3 as well as the first row of Figure 3. Although we have used the Hellinger distance to assess the performance of proposed density

estimators, we have also studied the  $L_2$  and  $L_1$  errors of each estimator. The results appear in Tables 2 and 3. As these tables show, one again  $\hat{g}_{\hat{\pi}_{LS}}$  can be the best estimator when the assumptions of a logistic missing probability mechanism is indeed true. Otherwise,  $\hat{g}_{\hat{\pi}_{ker}}$  is the best estimator.

We can draw the following conclusions form the results in tables 1, 2 and 3: if the missing probability mechanism,  $\pi(x)$  has a known form (such as the logistic model), then  $\hat{g}_{\hat{\pi}_{LS}}$  can be the best estimator. But in the more realistic case where one has no information about functional form  $\pi(x)$ , the density estimator  $\hat{g}_{\hat{\pi}_{LS}}$  is, in general, the most appropriate one.

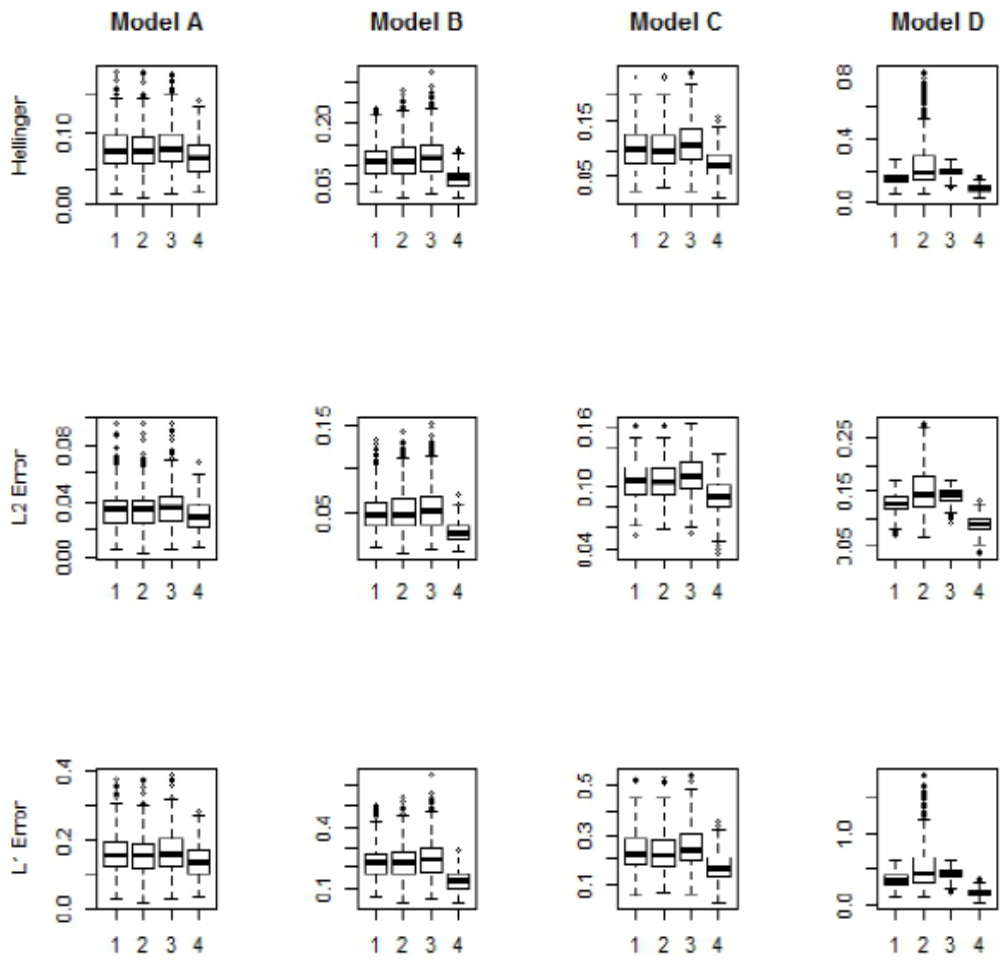


Figure 1: Boxplots of the errors of various estimator under different models, when  $n = 100$ . Within each of these 12 plots, boxplot 1 corresponds to  $\hat{g}_{\hat{\pi}_{ker}}$ , 2 corresponds to  $\hat{g}_{\hat{\pi}_{LS}}$ , 3 corresponds to  $\hat{g}_{cc}$  and 4 corresponds to the case with no missing data (i.e.  $\tilde{g}$ ).

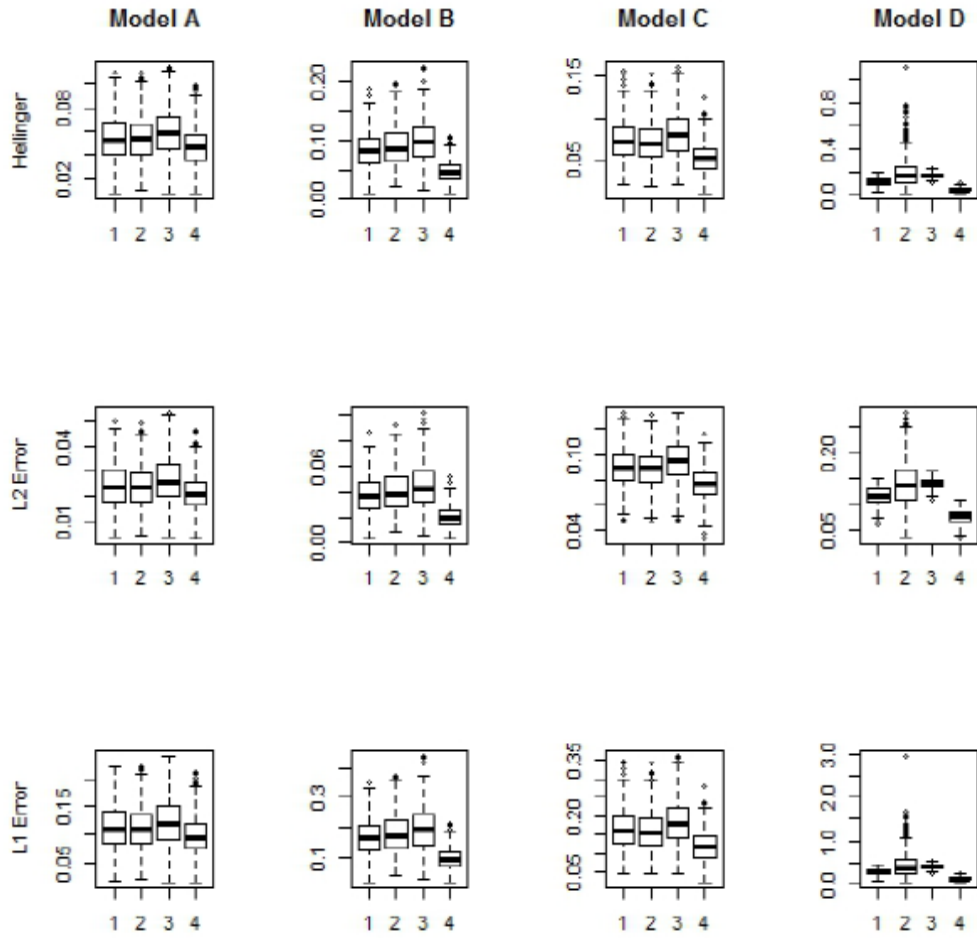


Figure 2: Boxplots of the errors of various estimator under different models, when  $n = 200$ . Within each of these 12 plots, boxplot 1 corresponds to  $\hat{g}_{\hat{\pi}_{ker}}$ , 2 corresponds to  $\hat{g}_{\hat{\pi}_{LS}}$ , 3 corresponds to  $\hat{g}_{cc}$  and 4 corresponds to the case with no missing data (i.e.  $\hat{g}$ ).

**Example B.** [*Regression function estimation.*] Here we consider the performance of the following two versions of the regression function estimator  $\phi_n$  defined via (13) and (14). The first estimator, denoted by  $\phi_{n,ker}$  is the minimizer of (13) when  $\hat{\pi}$  is taken to be the kernel estimator in (14). The second estimator, denoted by  $\phi_{n,LS}$ , is the minimizer of (14) when  $\hat{\pi}$  is taken to be the least squares estimator in (14). We also consider the complete case estimator, denoted by  $\phi_{n,cc}$ , and the estimator based on the full data of size  $n$  (i.e., when there are no missing values). To perform our numerical studies, we generated samples  $(Y_i, \mathbf{Z}_i)$ ,  $i = 1, \dots, n$  of sizes  $n = 100$  and  $n = 200$  form

$$Y = \sin(2Z_1) + Z_2^2 + Z_3 - \exp(-Z_4) + \varepsilon \quad \text{with } \varepsilon \approx N(0,0.5). \quad (20)$$

where  $\varepsilon$  is independent of  $\mathbf{Z} = (Z_1, Z_2, Z_3, Z_4)$  and  $\mathbf{Z}$  has a Gaussian distribution with mean  $\mathbf{0}$  and the covariance matrix  $\Sigma = (\sigma_{ij})_{i,j>1}$ , where  $\sigma_{ij} = 2^{-|i-j|}$ . In passing we also note that (20) is similar to the model used in Meier et al. (2009). Here,  $Z_1$  and  $Z_2$  are always observable, but  $Z_3$  and  $Z_4$  are allowed to be missing at random according to one of the following two missing probability mechanism:

*Model (I).*  $\pi(z_1, z_2) := P\{\zeta = 1 | Z_1 = z_1, Z_2 = z_2\} = 0.4[1 + \cos(z_1 + z_2)]$ .

*Model (II).*  $\pi(z_1, z_2) := P\{\zeta = 1 | Z_1 = z_1, Z_2 = z_2\} = \exp(2 - 3z_1 - z_2) / [1 + \exp(2 + z_1 + z_2)]$ .

Model (I) results in approximately 45% missing data where as model (II), which is logistic results in about 75% missing data. Since the true underlying data generating model (20) is never known in practice, we decided to fit the partial second order model  $D(Y|\mathbf{Z} = z) = z_1 + z_2 + z_3 + z_4 + z_4 + z_1^2 + z_2^2$ . As in Example A, we used the cross-validation method of RAcine and Li (2004) in the R package "np" (Racine and Hayfield 2008) to find the kernel regression estimator  $\hat{\pi}$  of  $\pi$  which is then used to find  $\phi_{n,ker}$  via (13) and (14). Similarly, the parameters of the logistic missing probability mechanism were estimated using nonlinear least squares regression (based on the R package "nls2"), which are then used to find  $\phi_{n,LS}$ . To assess the performance of these two estimators we computed the empirical  $L_2$  error of each estimator. This simulation process was repeated a total of 500 times (each time using samples of size  $n = 100$  and  $n = 200$  observations to find the least squares estimators  $\phi_{n,ker}$ ,  $\phi_{n,LS}$  and  $\phi_{cc}$ ) and the average  $L_2$  errors were computed. The results appear in Table 4; the numbers appearing in brackets are the standard errors computed over 500 Monte Carlo runs. As a point of reference, we have also included the estimator corresponding to the case with no missing data this appears as  $\tilde{\phi}_n$  in Table 4.

Table 4: Average  $L_2$  error, over 500 Monte Carlo runs corresponding to models (I) and (II) for regression function estimators  $\phi_{n,LS}$ ,  $\phi_{n,LS}$ ,  $\phi_{n,cc}$  and  $\tilde{\phi}_n$ . Here  $\tilde{\phi}_n$  is based on the data with no missing values.

	n=100				n=200			
	$\hat{\phi}_{n,ker}$	$\hat{\phi}_{n,LS}$	$\hat{\phi}_{cc}$	$\tilde{\phi}_n$	$\hat{\phi}_{n,ker}$	$\hat{\phi}_{n,LS}$	$\hat{\phi}_{cc}$	$\tilde{\phi}_n$
Model (I)	.9733 (.02463)	1.0741 (.02284)	1.0537 (.02849)	.7324 (.00598)	.8008 (.00758)	.8814 (.00730)	.8690 (.01039)	.6704 (.00337)
Model (II)	4.6476 (.27083)	4.4798 (.27177)	6.0594 (.30785)	.7324 (.00598)	2.8740 (.14868)	2.5824 (.14347)	5.2582 (.20042)	.6704 (.00337)

Table 4 shows,  $\phi_{n,ker}$  outperforms both  $\phi_{n,LS}$  and  $\phi_{n,cc}$  under Model (I). On the other hand when Model (II) is correct,  $\phi_{n,LS}$  is better. It is also important to notice that, under Model (II), the error of  $\phi_{n,cc}$  is substantially larger than every other estimator. In general, since the popular logistic model does not necessarily hold true, it would be safer to use  $\phi_{n,ker}$  instead of  $\phi_{n,LS}$  in practice. Figure 3 gives the boxplots of the 500  $L_2$  errors of various estimators. These boxplots show that the estimator are much more variable under Model (II).



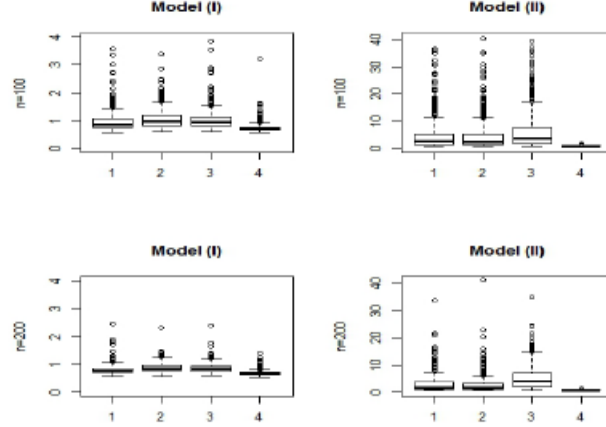


Figure 3: Boxplots of the  $L_2$  errors of various regression estimators under Models (I) and (II). Within each of these 4 plots, boxplot 1 corresponds to  $\hat{\phi}_{n,ker}$ , 2 corresponds to  $\hat{\phi}_{n,LS}$ , 3 corresponds to  $\hat{\phi}_{n,cc}$  and 4 corresponds to the case with no missing data (i.e,  $p\tilde{h}i_n$ ).

## 4 Appendix

To prove theorems 2 and 3 we first state two technical lemmas which may be of some independent theoretical interests as well.

**Lemma 1** Define  $\check{\nu}_n(\psi) = n^{-1} \sum_{i=1}^n \xi_i \psi(\mathbf{Z}_i) / \pi(\mathbf{X}_i)$ , where the function  $\pi(\mathbf{X}_i)$  is as in (3) and  $\xi_i$ 's are the Bernoulli random variable that appear in (1). Suppose that Assumption (A1) holds. Then for every  $\epsilon > 0$  and  $n \geq 1$

$$\mathbb{P} \left\{ \sup_{\psi \in \Psi} |\check{\nu}_n(\psi) - \nu(\psi)| > \epsilon \right\} \leq 8 \mathbb{E} \left[ \mathcal{N}_1 \left( \frac{\pi_0 \epsilon}{8}, \Psi, \mathbb{D}_n \right) \right] e^{-n\pi_0^2 \epsilon^2 / (128B^2)},$$

where  $B$  is as in Theorem 2 and  $\pi_0 = \inf_{\mathbf{x}} \pi(\mathbf{x}) > 0$ .

**Lemma 2** Let  $f$  be the pdf of the random vector  $\mathbf{X}$  and put  $T(\mathbf{x}) = \pi(\mathbf{x})f(\mathbf{x})$ , where  $\pi(\mathbf{x}) = \mathbb{E}[\xi | \mathbf{X} = \mathbf{x}]$ . Let  $\hat{T}(\mathbf{x}) = (nh_n^d)^{-1} \sum_{j=1}^n \xi_j \mathcal{K}((\mathbf{x} - \mathbf{X}_j)/h_n)$ . Then

(i) Under conditions (A2), (A3), and (A4),

$$\left| \mathbb{E} \left[ \hat{T}(\mathbf{X}) | \mathbf{X} \right] - T(\mathbf{X}) \right| \stackrel{a.s.}{\leq} ch_n,$$

where  $c > 0$  is a constant not depending on  $n$ .

(ii) For every constant  $\beta > 0$ ,

$$\mathbb{P} \left\{ \left| \hat{T}(\mathbf{X}) - \mathbb{E} \left[ \hat{T}(\mathbf{X}) | \mathbf{X} \right] \right| > \beta \mid \mathbf{X} = \mathbf{x} \right\} \leq 2 \exp \left\{ \frac{-nh_n^d \beta^2}{2\|\mathcal{K}\|_\infty [2\|f\|_\infty + h_n^{2d} \beta / 3]} \right\}.$$

PROOF OF THEOREM 2.

Let  $\check{\nu}_n(\psi) = n^{-1} \sum_{i=1}^n \xi_i \psi(\mathbf{Z}_i) / \pi(\mathbf{X}_i)$  and observe that

$$\begin{aligned} \mathbb{P} \left\{ \sup_{\psi \in \Psi} \left| \hat{\nu}_n^{(ker)}(\psi) - \nu(\psi) \right| > \epsilon \right\} &\leq \mathbb{P} \left\{ \sup_{\psi \in \Psi} \left| \hat{\nu}_n^{(ker)}(\psi) - \check{\nu}_n(\psi) \right| > \frac{\epsilon}{2} \right\} + \mathbb{P} \left\{ \sup_{\psi \in \Psi} \left| \check{\nu}_n(\psi) - \nu(\psi) \right| > \frac{\epsilon}{2} \right\} \\ &:= \Delta_{n,1}(\epsilon) + \Delta_{n,2}(\epsilon). \end{aligned} \quad (21)$$

Now by Lemma 1, for every  $\epsilon > 0$  and every  $n \geq 1$ ,

$$\Delta_{n,2}(\epsilon) \leq 8 \mathbb{E}[\mathcal{N}_1(\pi_0\epsilon/16, \Psi, \mathbb{D}_n)] e^{-n\pi_0^2\epsilon^2/(512B^2)}. \quad (22)$$

To deal with the term  $\Delta_{n,1}(\epsilon)$  in (21) first note that

$$|\widehat{\nu}_n^{(\ker)}(\psi) - \check{\nu}_n(\psi)| \leq \frac{1}{n} \sum_{i=1}^n \xi_i |\psi(\mathbf{Z}_i)| \cdot \left| \frac{1}{\widehat{\pi}_{\ker}(\mathbf{X}_i)} - \frac{1}{\pi(\mathbf{X}_i)} \right| \leq \frac{B}{n} \sum_{i=1}^n \left| \frac{1}{\widehat{\pi}_{\ker}(\mathbf{X}_i)} - \frac{1}{\pi(\mathbf{X}_i)} \right|,$$

where  $B = \|\psi\|_\infty$ . Thus

$$\begin{aligned} \Delta_{n,1}(\epsilon) &\leq \mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \left| \frac{1}{\widehat{\pi}_{\ker}(\mathbf{X}_i)} - \frac{1}{\pi(\mathbf{X}_i)} \right| > \frac{\epsilon}{2B} \right\} \\ &\leq \mathbb{P} \left\{ \left[ \frac{1}{n} \sum_{i=1}^n \left| \frac{\pi(\mathbf{X}_i) - \widehat{\pi}_{\ker}(\mathbf{X}_i)}{\widehat{\pi}_{\ker}(\mathbf{X}_i) \pi(\mathbf{X}_i)} \right| > \frac{\epsilon}{2B} \right] \cap \bigcap_{i=1}^n [\widehat{\pi}_{\ker}(\mathbf{X}_i) \geq 2^{-1}\pi_0] \right\} \\ &\quad + \mathbb{P} \left\{ \bigcup_{i=1}^n [\widehat{\pi}_{\ker}(\mathbf{X}_i) < 2^{-1}\pi_0] \right\} \\ &\leq \mathbb{P} \left\{ n^{-1} \sum_{i=1}^n |\widehat{\pi}_{\ker}(\mathbf{X}_i) - \pi(\mathbf{X}_i)| > (4B)^{-1}\pi_0^2\epsilon \right\} + \mathbb{P} \left\{ \bigcup_{i=1}^n [\widehat{\pi}_{\ker}(\mathbf{X}_i) < 2^{-1}\pi_0] \right\} \\ &:= \Delta_{n,1}^{(i)}(\epsilon) + \Delta_{n,1}^{(ii)}(\epsilon). \end{aligned} \quad (23)$$

Let  $T(\mathbf{X}) = \pi(\mathbf{X})f(\mathbf{X})$  and define  $\widehat{T}(\mathbf{X}_i) = ((n-1)h_n^d)^{-1} \sum_{j=1, \neq i}^n \xi_j \mathcal{K}((\mathbf{X}_i - \mathbf{X}_j)/h_n)$  and  $\widehat{f}(\mathbf{X}_i) = ((n-1)h_n^d)^{-1} \sum_{j=1, \neq i}^n \mathcal{K}((\mathbf{X}_i - \mathbf{X}_j)/h_n)$ , and note that since  $|\widehat{T}(\mathbf{X}_i)/\widehat{f}(\mathbf{X}_i)| \leq 1$ , one finds  $|\widehat{\pi}_{\ker}(\mathbf{X}_i) - \pi(\mathbf{X}_i)| \leq [|\widehat{T}(\mathbf{X}_i) - T(\mathbf{X}_i)| + |\widehat{f}(\mathbf{X}_i) - f(\mathbf{X}_i)|]/f(\mathbf{X}_i)$ . Therefore

$$\begin{aligned} \Delta_{n,1}^{(i)}(\epsilon) &\leq \sum_{i=1}^n \mathbb{P} \left\{ \frac{1}{f(\mathbf{X}_i)} \left| \widehat{T}(\mathbf{X}_i) - T(\mathbf{X}_i) \right| > (8B)^{-1}\pi_0^2\epsilon \right\} \\ &\quad + \sum_{i=1}^n \mathbb{P} \left\{ \frac{1}{f(\mathbf{X}_i)} \left| \widehat{f}(\mathbf{X}_i) - f(\mathbf{X}_i) \right| > (8B)^{-1}\pi_0^2\epsilon \right\} \\ &:= \sum_{i=1}^n p_{n,i}(\epsilon) + \sum_{i=1}^n q_{n,i}(\epsilon). \end{aligned} \quad (24)$$

But

$$\begin{aligned} p_{n,i}(\epsilon) &\leq \mathbb{P} \left\{ \left| \widehat{T}(\mathbf{X}_i) - \mathbb{E}(\widehat{T}(\mathbf{X}_i)|\mathbf{X}_i) + \mathbb{E}(\widehat{T}(\mathbf{X}_i)|\mathbf{X}_i) - T(\mathbf{X}_i) \right| > (8B)^{-1}f_0\pi_0^2\epsilon \right\} \\ &\quad (\text{because } f_0 := \inf f(\mathbf{x}) > 0 \text{ by Assumption (A3)}), \\ &\leq \mathbb{P} \left\{ \left| \widehat{T}(\mathbf{X}_i) - \mathbb{E}(\widehat{T}(\mathbf{X}_i)|\mathbf{X}_i) \right| + (16B)^{-1}f_0\pi_0^2\epsilon > (8B)^{-1}f_0\pi_0^2\epsilon \right\} \\ &\quad (\text{by Part (i) of Lemma 2, for } n \text{ large enough}) \\ &= \mathbb{E} \left[ \mathbb{P} \left\{ \left| \widehat{T}(\mathbf{X}_i) - \mathbb{E}(\widehat{T}(\mathbf{X}_i)|\mathbf{X}_i) \right| > (16B)^{-1}f_0\pi_0^2\epsilon \mid \mathbf{X}_i \right\} \right] \\ &\leq 2 \exp \left\{ \frac{-(n-1)h_n^d f_0^2 \pi_0^4 \epsilon^2}{(16B)^2 (2\|\mathcal{K}\|_\infty) [2\|f\|_\infty + h_n^{2d} f_0 \pi_0^2 \epsilon / (48B)]} \right\} \\ &\quad (\text{by Part (ii) of Lemma 2}) \\ &\leq 2 \exp \left\{ \frac{-(n-1)h_n^d f_0^2 \pi_0^4 \epsilon^2}{2(16B)^2 \|\mathcal{K}\|_\infty [2\|f\|_\infty + f_0/12]} \right\}, \quad (\text{for large } n), \end{aligned} \quad (25)$$

where we have used the fact that in bounding  $\mathbb{P}\{|\widehat{\pi}_{\ker}(\mathbf{X}_i) - \pi(\mathbf{X}_i)| > (4B)^{-1}\pi_0^2\epsilon\}$ , one only needs to consider  $0 < \epsilon \leq 4B/\pi_0^2$  because  $|\widehat{\pi}_{\ker}(\mathbf{X}_i) - \pi(\mathbf{X}_i)| \leq 1$ . Similarly, since  $\widehat{f}(\mathbf{X}_i)$  is the special case of

$\widehat{T}(\mathbf{X}_i)$  (take  $\xi_j = 1$  in the definition of  $\widehat{T}(\mathbf{X}_i)$ , for all  $j$ ), the above arguments leading to (26) give

$$q_{n,i}(\epsilon) \leq 2 \exp \left\{ \frac{-(n-1)h_n^d f_0^2 \pi_0^4 \epsilon^2}{2(16B)^2 \|\mathcal{K}\|_\infty [2\|f\|_\infty + f_0/12]} \right\},$$

for  $n$  large enough. Thus, in view of (25), for  $n$  large enough

$$\Delta_{n,1}^{(i)}(\epsilon) \leq 4n e^{-(n-1)h_n^d C_{15} \epsilon^2}, \quad (27)$$

where  $C_{15} = \{2(16B)^2 \|\mathcal{K}\|_\infty [2\|f\|_\infty + f_0/12]\}^{-1} f_0^2 \pi_0^4$ . To complete the proof of Theorem 2, we also need to bound the term  $\Delta_{n,1}^{(ii)}(\epsilon)$  in (24). Since  $\mathbb{P}\{\widehat{\pi}_{\ker}(\mathbf{X}_i) < \pi_0/2\} \leq \mathbb{P}\{|\widehat{\pi}_{\ker}(\mathbf{X}_i) - \pi(\mathbf{X}_i)| > \pi_0/2\}$ , the arguments that lead to the bound on  $\Delta_{n,1}^{(i)}(\epsilon)$  (see (25) and (27)) yield

$$\begin{aligned} \Delta_{n,1}^{(ii)}(\epsilon) &\leq \sum_{i=1}^n \mathbb{P}\{\widehat{\pi}_{\ker}(\mathbf{X}_i) < \pi_0/2\} \leq \sum_{i=1}^n \mathbb{P}\{|\widehat{\pi}_{\ker}(\mathbf{X}_i) - \pi(\mathbf{X}_i)| > \pi_0/2\} \\ &\leq 2n e^{-(n-1)h_n^d C_{16}}, \end{aligned} \quad (28)$$

where  $C_{16} = \{128\|\mathcal{K}\|_\infty (2\|f\|_\infty + f_0\pi_0/24)\}^{-1} f_0^2 \pi_0^2$ . Therefore,  $\Delta_{n,1}(\epsilon) \leq \Delta_{n,1}^{(i)}(\epsilon) + \Delta_{n,1}^{(ii)}(\epsilon) \leq 6n e^{-nh^d C_{17}}$ , with  $C_{17} = \frac{1}{2} \min(C_{15}, C_{16})$ . The theorem now follows from the bounds in (21), (22), (24), (27), and (28).  $\square$

### PROOF OF THEOREM 3.

We start by writing,

$$\mathbb{P} \left\{ \sup_{\psi \in \Psi} |\widehat{\nu}_n^{(\text{LS})}(\psi) - \nu(\psi)| > \epsilon \right\} \leq \mathbb{P} \left\{ \sup_{\psi \in \Psi} |\widehat{\nu}_n^{(\text{LS})}(\psi) - \check{\nu}_n(\psi)| > \frac{\epsilon}{2} \right\} + \mathbb{P} \left\{ \sup_{\psi \in \Psi} |\check{\nu}_n(\psi) - \nu(\psi)| > \frac{\epsilon}{2} \right\},$$

where  $\check{\nu}_n(\psi) = n^{-1} \sum_{i=1}^n \xi_i \psi(\mathbf{Z}_i) / \pi(\mathbf{X}_i)$ . But, by Lemma 1, for every  $n \geq 1$ ,

$$\mathbb{P} \left\{ \sup_{\psi \in \Psi} |\check{\nu}_n(\psi) - \nu(\psi)| > \frac{\epsilon}{2} \right\} \leq 8 \mathbb{E} [\mathcal{N}_1(\pi_0 \epsilon / 16, \Psi, \mathbb{D}_n)] e^{-n\pi_0^2 \epsilon^2 / (512B^2)}.$$

Furthermore, since  $|\widehat{\nu}_n^{(\text{LS})}(\psi) - \check{\nu}_n(\psi)| \leq Bn^{-1} \sum_{i=1}^n |(\widehat{\pi}_{\text{LS}}(\mathbf{X}_i))^{-1} - (\pi(\mathbf{X}_i))^{-1}|$ , one finds

$$\begin{aligned} &\mathbb{P} \left\{ \sup_{\psi \in \Psi} |\widehat{\nu}_n^{(\text{LS})}(\psi) - \check{\nu}_n(\psi)| > \frac{\epsilon}{2} \right\} \\ &\leq \mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \left| \frac{1}{\widehat{\pi}_{\text{LS}}(\mathbf{X}_i)} - \frac{1}{\pi(\mathbf{X}_i)} \right| > \frac{\epsilon}{2B} \right\} \leq \mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n |\pi(\mathbf{X}_i) - \widehat{\pi}_{\text{LS}}(\mathbf{X}_i)| > \frac{\pi_0^2 \epsilon}{2B} \right\} \\ &\quad (\text{the second inequality follows because } \widehat{\pi}_{\text{LS}}(\mathbf{X}_i) \geq \pi_0 \text{ and } \pi(\mathbf{X}_i) \geq \pi_0) \\ &\leq \mathbb{P} \left\{ \sup_{\tilde{\pi} \in \mathcal{P}} \left| \frac{1}{n} \sum_{i=1}^n |\tilde{\pi}(\mathbf{X}_i) - \pi(\mathbf{X}_i)| - \mathbb{E} |\tilde{\pi}(\mathbf{X}_i) - \pi(\mathbf{X}_i)| \right| > \frac{\pi_0^2 \epsilon}{4B} \right\} \\ &\quad + \mathbb{P} \left\{ \mathbb{E} \left[ |\widehat{\pi}_{\text{LS}}(\mathbf{X}) - \pi(\mathbf{X})| \mid \mathbb{D}_n \right] > \frac{\pi_0^2 \epsilon}{4B} \right\} \\ &:= I_n(\epsilon) + II_n(\epsilon). \end{aligned} \quad (29)$$

Using Theorem 1, with  $\Psi$  replaced by the class  $\mathcal{P}$ , it is straightforward to see that

$$I_n(\epsilon) \leq 8 \mathbb{E} \left[ \mathcal{N}_1 \left( \frac{\pi_0^2 \epsilon}{32B}, \mathcal{P}, (\mathbf{X}_i)_{i=1}^n \right) \right] e^{-nC_{19} \epsilon^2},$$

where  $C_{19} = \pi_0^4/(2048B^2)$ . As for the term  $II_n(\epsilon)$  in (29), let

$$\widehat{L}_n(\tilde{\pi}) = n^{-1} \sum_{i=1}^n (\xi_i - \tilde{\pi}(\mathbf{X}_i))^2, \forall \tilde{\pi} \in \mathcal{P},$$

and note that by the Cauchy-Schwartz inequality

$$\begin{aligned} II_n(\epsilon) &\leq \mathbb{P} \left\{ \mathbb{E} \left[ (\widehat{\pi}_{\text{LS}}(\mathbf{X}) - \pi(\mathbf{X}))^2 \mid \mathbb{D}_n \right] > \frac{\pi_0^4 \epsilon^2}{16B^2} \right\} \\ &= \mathbb{P} \left\{ \mathbb{E} \left[ (\widehat{\pi}_{\text{LS}}(\mathbf{X}) - \xi)^2 \mid \mathbb{D}_n \right] - \mathbb{E}(\pi(\mathbf{X}) - \xi)^2 > \frac{\pi_0^4 \epsilon^2}{16B^2} \right\}, \quad (\text{since } \pi(\mathbf{X}) = \mathbb{E}(\xi \mid \mathbf{X})) \\ &\leq \mathbb{P} \left\{ 2 \sup_{\tilde{\pi} \in \mathcal{P}} \left| \widehat{L}(\tilde{\pi}) - \mathbb{E}(\tilde{\pi}(\mathbf{X}) - \xi)^2 \right| > \frac{\pi_0^4 \epsilon^2}{16B^2} \right\}, \end{aligned} \quad (30)$$

where (30) follows from the fact that, since  $\mathbb{E}(\pi(\mathbf{X}) - \xi)^2 = \inf_{\tilde{\pi} \in \mathcal{P}} \mathbb{E}(\tilde{\pi}(\mathbf{X}) - \xi)^2$ , one has

$$\begin{aligned} &\mathbb{E} \left[ (\widehat{\pi}_{\text{LS}}(\mathbf{X}) - \xi)^2 \mid \mathbb{D}_n \right] - \mathbb{E}(\pi(\mathbf{X}) - \xi)^2 \\ &\leq \sup_{\tilde{\pi} \in \mathcal{P}} \left\{ \mathbb{E} \left[ (\widehat{\pi}_{\text{LS}}(\mathbf{X}) - \xi)^2 \mid \mathbb{D}_n \right] - \widehat{L}(\widehat{\pi}_{\text{LS}}) + \widehat{L}(\widehat{\pi}_{\text{LS}}) - \widehat{L}(\tilde{\pi}) + \widehat{L}(\tilde{\pi}) - \mathbb{E}(\tilde{\pi}(\mathbf{X}) - \xi)^2 \right\} \\ &\leq 2 \sup_{\tilde{\pi} \in \mathcal{P}} \left| \widehat{L}(\tilde{\pi}) - \mathbb{E}(\tilde{\pi}(\mathbf{X}) - \xi)^2 \right|, \quad (\text{because } \widehat{L}(\widehat{\pi}_{\text{LS}}) - \widehat{L}(\tilde{\pi}) \leq 0, \forall \tilde{\pi} \in \mathcal{P}). \end{aligned}$$

Finally, using Theorem 1, we find

$$(30) \leq 8 \mathbb{E} \left[ \mathcal{N}_1 \left( \frac{\pi_0^4 \epsilon^2}{(16)^2 B}, \mathcal{P}, (\mathbf{X}_i)_{i=1}^n \right) \right] e^{-nC_{20}\epsilon^4},$$

where  $C_{20} = \pi_0^8/((128)(32)^2 B^2)$ . This completes the proof of Theorem 3.  $\square$

#### PROOF OF LEMMA 1.

The proof is based on the *symmetrization* arguments of Dudley and Pollard (Dudley (1978, P.925) and Pollard (1984, Sec. II.3)); also see van der Vaart and Wellner (1996, Sec. 2.3). Let  $\mathbb{D}'_n = \{(\mathbf{X}'_1, \mathbf{V}'_1, \xi'_1), \dots, (\mathbf{X}'_n, \mathbf{V}'_n, \xi'_n)\}$  be a hypothetical sample (a *ghost* sample) independent of the data  $\mathbb{D}_n$ , where  $(\mathbf{X}'_i, \mathbf{V}'_i, \xi'_i) \stackrel{\text{iid}}{=} (\mathbf{X}_1, \mathbf{V}_1, \xi_1)$ ,  $i = 1, \dots, n$ , and put

$$\check{\nu}'_n(\psi) = n^{-1} \sum_{i=1}^n \xi'_i \psi(\mathbf{Z}'_i) / \pi(\mathbf{X}'_i).$$

Next, fix the data  $\mathbb{D}_n$  and observe that if  $\sup_{\psi \in \Psi} |\check{\nu}'_n(\psi) - \nu(\psi)| > \epsilon$ , then there is at least one  $\psi_\epsilon \in \Psi$ , which depends on  $\mathbb{D}_n$  (but not  $\mathbb{D}'_n$ ), such that  $|\check{\nu}'_n(\psi_\epsilon) - \nu(\psi_\epsilon | \mathbb{D}_n)| > \epsilon$ , where  $\nu(\psi_\epsilon | \mathbb{D}_n) = \mathbb{E}[\psi_\epsilon(\mathbf{Z}) | \mathbb{D}_n]$ . Now observe that for  $n\epsilon^2 \geq 8B^2/\pi_0^2$

$$\begin{aligned} \mathbb{P} \left\{ |\check{\nu}'_n(\psi_\epsilon) - \nu(\psi_\epsilon | \mathbb{D}_n)| < \frac{\epsilon}{2} \mid \mathbb{D}_n \right\} &= 1 - \mathbb{P} \left\{ |\check{\nu}'_n(\psi_\epsilon) - \nu(\psi_\epsilon | \mathbb{D}_n)| \geq \frac{\epsilon}{2} \mid \mathbb{D}_n \right\} \\ &\geq 1 - \sup_{\psi \in \Psi} \mathbb{P} \left\{ |\check{\nu}'_n(\psi) - \nu(\psi)| \geq \frac{\epsilon}{2} \right\} \\ &\geq 1 - \sup_{\psi \in \Psi} \mathbb{E}(\check{\nu}'_n(\psi) - \nu(\psi))^2 / (\epsilon/2)^2, \quad (\text{Markov's inequality}) \\ &\geq 1 - \frac{4}{n\epsilon^2} \sup_{\psi \in \Psi} \text{Var}(\xi_1 \psi(\mathbf{Z}_1) / \pi(\mathbf{X}_1)) \\ &\geq 1 - \frac{4B^2}{n\pi_0^2 \epsilon^2} \geq \frac{1}{2}, \quad (\text{because } n\epsilon^2 \geq 8B^2/\pi_0^2), \end{aligned} \quad (31)$$

where (31) follows from the MAR assumption (??) and the fact that

$$\mathbb{E}[\xi \psi(\mathbf{Z}) / \pi(\mathbf{X})] = \mathbb{E}[\mathbb{E}(\xi \psi(\mathbf{Z}) / \pi(\mathbf{X}) | \mathbf{Z})] \stackrel{\text{by (??)}}{=} \mathbb{E}(\psi(\mathbf{Z})) =: \nu(\psi).$$

Therefore, for  $n\epsilon^2 \geq 8B^2/\pi_0^2$ , we have

$$\begin{aligned}
\frac{1}{2} &\leq \mathbb{P} \left\{ |\check{\nu}'_n(\psi_\epsilon) - \nu(\psi_\epsilon|\mathbb{D}_n)| < \frac{\epsilon}{2} \middle| \mathbb{D}_n \right\} \\
&\leq \mathbb{P} \left\{ |\check{\nu}_n(\psi_\epsilon) - \nu(\psi_\epsilon|\mathbb{D}_n)| - |\check{\nu}'_n(\psi_\epsilon) - \check{\nu}_n(\psi_\epsilon)| < \frac{\epsilon}{2} \middle| \mathbb{D}_n \right\} \\
&\leq \mathbb{P} \left\{ |\check{\nu}'_n(\psi_\epsilon) - \check{\nu}_n(\psi_\epsilon)| > \frac{\epsilon}{2} \middle| \mathbb{D}_n \right\} \\
&\quad (\text{because by the definition of } \psi_\epsilon, \text{ conditional on } \mathbb{D}_n, |\check{\nu}_n(\psi_\epsilon) - \nu(\psi_\epsilon|\mathbb{D}_n)| > \epsilon) \\
&\leq \mathbb{P} \left\{ \sup_{\psi \in \Psi} |\check{\nu}'_n(\psi) - \check{\nu}_n(\psi)| > \frac{\epsilon}{2} \middle| \mathbb{D}_n \right\}. \tag{32}
\end{aligned}$$

Now observe that the far left and the far right sides of (32) do not depend on  $\psi_\epsilon$  and that the chain of inequalities between them remain valid on the set  $\{\sup_{\psi \in \Psi} |\check{\nu}_n(\psi) - \nu(\psi)| > \epsilon\}$ . Therefore, integrating the two far sides of (32) with respect to the distribution of  $\mathbb{D}_n$ , over this set, we find

$$\begin{aligned}
\mathbb{P}\left\{\sup_{\psi \in \Psi} |\check{\nu}_n(\psi) - \nu(\psi)| > \epsilon\right\} &\leq 2 \mathbb{P} \left\{ \sup_{\psi \in \Psi} |\check{\nu}'_n(\psi) - \check{\nu}_n(\psi)| > \frac{\epsilon}{2} \right\} \\
&\leq 2 \mathbb{P} \left\{ \sup_{\psi \in \Psi} \frac{1}{n} \left| \sum_{i=1}^n \left[ \frac{\xi'_i \psi(\mathbf{Z}'_i)}{\pi(\mathbf{X}'_i)} - \frac{\xi_i \psi(\mathbf{Z}_i)}{\pi(\mathbf{X}_i)} \right] \right| > \frac{\epsilon}{2} \right\}. \tag{33}
\end{aligned}$$

Next, let  $\sigma_1, \dots, \sigma_n$  be iid random variables, independent of  $\mathbb{D}_n$  and  $\mathbb{D}'_n$ , where  $\mathbb{P}\{\sigma_i = +1\} = \mathbb{P}\{\sigma_i = -1\} = 1/2$ . Observing that  $\sigma_i$ 's are random signs, we have

$$\begin{aligned}
(33) &= 2 \mathbb{P} \left\{ \sup_{\psi \in \Psi} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i \left[ \frac{\xi'_i \psi(\mathbf{Z}'_i)}{\pi(\mathbf{X}'_i)} - \frac{\xi_i \psi(\mathbf{Z}_i)}{\pi(\mathbf{X}_i)} \right] \right| > \frac{\epsilon}{2} \right\} \\
&\leq 4 \mathbb{P} \left\{ \sup_{\psi \in \Psi} \frac{1}{n} \left| \sum_{i=1}^n \frac{\sigma_i \xi_i \psi(\mathbf{Z}_i)}{\pi(\mathbf{X}_i)} \right| > \frac{\epsilon}{4} \right\} \\
&\leq 4 \mathbb{E} \left[ \mathbb{P} \left\{ \sup_{\psi \in \Psi} \frac{1}{n} \left| \sum_{i=1}^n \frac{\sigma_i \xi_i \psi(\mathbf{Z}_i)}{\pi(\mathbf{X}_i)} \right| > \frac{\epsilon}{4} \middle| \mathbb{D}_n \right\} \right]. \tag{34}
\end{aligned}$$

Now, put  $\epsilon' = \epsilon/8$  and, for fixed  $\mathbb{D}_n$ , let  $\Psi_{\epsilon'}$  be a weighted empirical  $L_1$   $\epsilon'$ -cover of  $\Psi$  based on the weights  $W_i = \xi_i/\pi(\mathbf{X}_i)$ . That is, for each  $\psi \in \Psi$  there is a  $\psi^* \in \Psi_{\epsilon'}$  such that  $n^{-1} \sum_{i=1}^n W_i |\psi(\mathbf{Z}_i) - \psi^*(\mathbf{Z}_i)| < \epsilon' = \epsilon/8$ . Let  $\Gamma_1(\epsilon', \Psi, \mathbb{D}_n)$  be the  $\epsilon'$ -covering number of  $\Psi$  with respect to the weighted empirical  $L_1$  norm. Then for some  $\psi^* \in \Psi_{\epsilon'}$  we have

$$\begin{aligned}
\frac{1}{n} \left| \sum_{i=1}^n \frac{\sigma_i \xi_i \psi(\mathbf{Z}_i)}{\pi(\mathbf{X}_i)} \right| &\leq \frac{1}{n} \left| \sum_{i=1}^n \frac{\sigma_i \xi_i}{\pi(\mathbf{X}_i)} [\psi(\mathbf{Z}_i) - \psi^*(\mathbf{Z}_i)] \right| + \frac{1}{n} \left| \sum_{i=1}^n \frac{\sigma_i \xi_i \psi^*(\mathbf{Z}_i)}{\pi(\mathbf{X}_i)} \right| \\
&\leq \frac{\epsilon}{8} + \frac{1}{n} \left| \sum_{i=1}^n \frac{\sigma_i \xi_i \psi^*(\mathbf{Z}_i)}{\pi(\mathbf{X}_i)} \right|.
\end{aligned}$$

Consequently

$$\begin{aligned}
(34) &\leq 4 \mathbb{E} \left[ \mathbb{P} \left\{ \sup_{\psi \in \Psi_{\epsilon'}} \frac{1}{n} \left| \sum_{i=1}^n \frac{\sigma_i \xi_i \psi(\mathbf{Z}_i)}{\pi(\mathbf{X}_i)} \right| > \frac{\epsilon}{8} \middle| \mathbb{D}_n \right\} \right] \\
&\leq 4 \mathbb{E} \left[ \Gamma_1(\epsilon', \Psi, \mathbb{D}_n) \cdot \max_{\psi \in \Psi_{\epsilon'}} \mathbb{P} \left\{ \frac{1}{n} \left| \sum_{i=1}^n \frac{\sigma_i \xi_i \psi(\mathbf{Z}_i)}{\pi(\mathbf{X}_i)} \right| > \frac{\epsilon}{8} \middle| \mathbb{D}_n \right\} \right] \\
&\leq 8 \mathbb{E} [\Gamma_1(\epsilon/8, \Psi, \mathbb{D}_n)] e^{-n\pi_0^2 \epsilon^2 / (128B^2)}, \quad (\text{by Hoeffding's inequality}).
\end{aligned}$$

However, for all functions  $\psi_1, \psi_2 : \mathbb{R}^{d+p} \rightarrow \mathbb{R}$  one has  $\sum_{i=1}^n (\xi_i / \pi(\mathbf{X}_i)) |\psi_1(\mathbf{Z}_i) - \psi_2(\mathbf{Z}_i)| \leq (1/\pi_0) \sum_{i=1}^n |\psi_1(\mathbf{Z}_i) - \psi_2(\mathbf{Z}_i)|$ . Therefore, if  $\{\phi_1, \dots, \phi_N\}$  is a minimal  $(\pi_0 \epsilon)$ -cover of  $\Psi$  with respect to the empirical  $L_1$  norm, then it is an  $\epsilon$ -cover of  $\Psi$  with respect to the weighted empirical  $L_1$  norm. Thus, for every  $\epsilon > 0$ , we find  $\Gamma_1(\epsilon, \Psi, \mathbb{D}_n) \leq \mathcal{N}_1(\epsilon \pi_0, \Psi, \mathbb{D}_n)$ . Putting all the above together, we have, for  $n\epsilon^2 \geq 8B^2/\pi_0^2$ ,

$$\mathbb{P}\left\{\sup_{\psi \in \Psi} |\check{\nu}_n(\psi) - \nu(\psi)| > \epsilon\right\} \leq 8\mathbb{E}\left[\mathcal{N}_1\left(\frac{\pi_0 \epsilon}{8}, \Psi, \mathbb{D}_n\right)\right] e^{-n\pi_0^2 \epsilon^2 / (128B^2)}.$$

When  $n\epsilon^2 < 8B^2/\pi_0^2$  the lemma is trivially true (because the bound in the lemma will exceed 1).  $\square$

## PROOF OF LEMMA 2.

Part (i). The proof is similar to (and in fact easier than) that of Lemma 2.2 of Mojirsheibani and Montazeri (2007) and goes as follows. First note that

$$\mathbb{E}\left[\widehat{T}(\mathbf{X}) \mid \mathbf{X}\right] = h^{-d} \mathbb{E}\left[\xi_1 \mathcal{K}((\mathbf{X} - \mathbf{X}_1)/h) \mid \mathbf{X}\right] \stackrel{a.s.}{=} h_n^{-d} \mathbb{E}\left[\mathcal{K}((\mathbf{X} - \mathbf{X}_1)/h_n) \mathbb{E}[\xi_1 \mid \mathbf{X}, \mathbf{X}_1] \mid \mathbf{X}\right].$$

Since  $\mathbb{E}[\xi_1 \mid \mathbf{X}, \mathbf{X}_1] = \mathbb{E}[\xi_1 \mid \mathbf{X}_1] = \pi(\mathbf{X}_1)$  (because  $\mathbf{X}$  is independent of  $\xi_1$  and  $\mathbf{X}_1$ ), we find

$$\begin{aligned} \mathbb{E}\left[\widehat{T}(\mathbf{X}) \mid \mathbf{X}\right] - T(\mathbf{X}) &= h_n^{-d} \mathbb{E}\left[(\pi(\mathbf{X}_1) - \pi(\mathbf{X})) \mathcal{K}((\mathbf{X} - \mathbf{X}_1)/h_n) \mid \mathbf{X}\right] \\ &\quad + \mathbb{E}\left[\pi(\mathbf{X}) \left\{h_n^{-d} \mathcal{K}((\mathbf{X} - \mathbf{X}_1)/h_n) - f(\mathbf{X})\right\} \mid \mathbf{X}\right] \\ &:= R_{n,1}(\mathbf{X}) + R_{n,2}(\mathbf{X}). \end{aligned}$$

Now a one-term Taylor expansion gives

$$R_{n,1}(\mathbf{X}) = h_n^{-d} \mathbb{E}\left[\left(\sum_{i=1}^d (X_i - X_{1,i}) \partial \pi(\mathbf{X}^*) / \partial X_i\right) \times \mathcal{K}((\mathbf{X} - \mathbf{X}_1)/h_n) \mid \mathbf{X}\right],$$

where  $X_i$  and  $X_{1,i}$  are the  $i^{\text{th}}$  components of  $\mathbf{X}$  and  $\mathbf{X}_1$ , respectively, and  $\mathbf{X}^*$  is a point on the interior of the line segment joining the points  $\mathbf{X}$  and  $\mathbf{X}_1$ . Therefore,

$$\begin{aligned} |R_{n,1}(\mathbf{X})| &\leq C_{10} \sum_{i=1}^d \mathbb{E}\left[|X_{1,i} - X_i| h_n^{-d} \mathcal{K}((\mathbf{X} - \mathbf{X}_1)/h_n) \mid \mathbf{X}\right] \\ &\quad (\text{where } C_{10} = \max_{1 \leq i \leq d} \sup_{\mathbf{x}} |\partial \pi(\mathbf{x}) / \partial x_i| < \infty, \text{ by Assumption A4}) \\ &= C_{10} \sum_{i=1}^d \int |x_i - X_i| h_n^{-d} \mathcal{K}((\mathbf{X} - \mathbf{x})/h_n) f(\mathbf{x}) d\mathbf{x} \\ &\leq C_{10} \|f\|_{\infty} \sum_{i=1}^d h_n \int |u_i| \mathcal{K}(\mathbf{u}) d\mathbf{u} = C_{11} h_n, \quad (\text{by Assumptions A2 and A4}), \end{aligned}$$

where  $0 < C_{11} < \infty$ . As for the term  $R_{n,2}(\mathbf{X})$ , we have

$$\begin{aligned} |R_{n,2}(\mathbf{X})| &= \left| \pi(\mathbf{X}) \int h^{-d} \mathcal{K}((\mathbf{X} - \mathbf{x})/h_n) [f(\mathbf{x}) - f(\mathbf{X})] d\mathbf{x} \right| \\ &= \left| \pi(\mathbf{X}) \int [f(\mathbf{X} - h_n \mathbf{v}) - f(\mathbf{X})] \mathcal{K}(\mathbf{v}) d\mathbf{v} \right| \\ &\leq \left( \sum_{i=1}^d \sup_{\mathbf{x}} |\partial f(\mathbf{x}) / \partial x_i| \cdot \int |v_i| \mathcal{K}(\mathbf{v}) d\mathbf{v} \right) h_n = C_{12} h_n, \end{aligned}$$

where  $0 < C_{12} < \infty$ , by Assumptions A2. This completes the proof of part (i).

Part (ii). For  $j = 1, \dots, n$ , let

$$S_j(\mathbf{X}) = h_n^{-d} \{ \xi_j \mathcal{K}((\mathbf{X} - \mathbf{X}_j)/h_n) - \mathbb{E}[\xi_j \mathcal{K}((\mathbf{X} - \mathbf{X}_j)/h_n) | \mathbf{X}] \},$$

and observe that, conditional on  $\mathbf{X}$ , the terms  $S_j(\mathbf{X})$  are independent, zero-mean random variables, bounded by  $-h_n^{-d} \|\mathcal{K}\|_\infty$  and  $+h_n^{-d} \|\mathcal{K}\|_\infty$ . Furthermore,  $\text{Var}(S_j(\mathbf{X}) | \mathbf{X}) = \mathbb{E}[S_j^2(\mathbf{X}) | \mathbf{X}] \leq 2h_n^{-d} \|\mathcal{K}\|_\infty \|f\|_\infty$ . Therefore, by Bernstein's (1946) inequality,

$$\begin{aligned} \mathbb{P} \left\{ \left| \hat{T}(\mathbf{X}) - \mathbb{E} \left[ \hat{T}(\mathbf{X}) | \mathbf{X} \right] \right| > \beta \middle| \mathbf{X} = \mathbf{x} \right\} &= \mathbb{P} \left\{ \frac{1}{n} \left| \sum_{j=1}^n S_j(\mathbf{X}) \right| > \beta \middle| \mathbf{X} = \mathbf{x} \right\} \\ &\leq 2 \exp \left\{ \frac{-n\beta^2}{2 \left[ 2h_n^{-d} \|\mathcal{K}\|_\infty \|f\|_\infty + h_n^d \|\mathcal{K}\|_\infty \beta / 3 \right]} \right\}, \end{aligned}$$

which does not depend on  $\mathbf{x}$ . The lemma now follows upon integrating both sides with respect to the distribution of  $\mathbf{X}$ . □

## References

- Ahmad, I.A. and Mugdadi, A.R. 2006. Weighted Hellinger distance as an error criterion for bandwidth selection in kernel estimation. *J. Nonpar. Statist.* 18, 215-236.
- Alexander, K. 1984. Probability inequalities for empirical processes and a law of the iterated logarithm. *Ann. Probab.* 12, 1041-1067.
- Bernstein, S. 1946. *The theory of probabilities.* Gastehizdat Publishing House, Moscow.
- Bravo, F. 2015. Semiparametric estimation with missing covariates. *Journal of Multivariate Analysis*, 139, 329-346.
- Cheng P.E. 1994. Nonparametric estimation of mean functionals with data missing at random. *Journal of the American Statistical Association*, 87, 81-87.
- Cheng, P.E. and Chu, C.K. 1996. Kernel estimation of distribution functions and quantiles with missing data. *Statistica Sinica*, 6, 63-78.
- Chenouri, S., Mojirsheibani, M., and Montazeri, Z. 2009. Empirical measures for incomplete data with applications. *Electron. J. Stat.* 3, 1021-1038.
- Devroye, L., Györfi, L., Lugosi, G. 1996. *A probabilistic theory of pattern recognition.* Springer-Verlag, New York.
- Dubnicka, S. 2009. Kernel density estimation with missing data and auxiliary variables. *Aust. N. Z. J. Stat.* 51, 247-270.
- Dudley, R.M. 1978. Central limit theorems for empirical measures. *Ann. Probab.* 6, 899-929.
- Giné, E. 1996. Empirical processes and applications: an overview. *Bernoulli* 2, 1-8.
- Györfi, L., Kohler, M., Krzyżak, A., and Walk, H. 2002. *A distribution-free theory of nonparametric regression.* Springer-Verlag, New York.
- Härdle, W., Marron, J. 1995 Optimal bandwidth selection in nonparametric regression function estimation. *Ann Stat* 13:1465-1481.
- Hazelton, M.L. 2000. Marginal density estimation from incomplete bivariate data. *Statistics & Probability Letters*, 47, 75-84.
- Hirano, K.I. and Ridder, G. 2003. Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica*, 71, 1161-1189.
- Horvitz, D.G. and Thompson, D.J. 1952. A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47, 663-685.
- Hu, Z., Follmann, D., and Qin, J. 2011. Dimension reduced kernel estimation for distribution function with incomplete data. *J. Statist. Plann. Inference* 141, 3084-3093.
- Kim, J.K. and Yu, C.L. 2011. A semiparametric estimation of mean functionals with nonignorable missing data. *J. Amer. Statist. Assoc.* 106, 157-165.
- Little, R.J.A. and Rubin, D.B. 2002. *Statistical Analysis With Missing Data*, 2nd ed. Wiley, New York.
- Liu, X., Liu, P., and Zhou, Y. 2011. Distribution estimation with auxiliary information for missing data. *J. Statist. Plann. Inference* 141, 711-724.
- Massart, P. 1990. The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *Ann. Probab.* 18, 1269-1283.
- Mojirsheibani, M. and Montazeri, Z. 2007. On nonparametric classification with missing covariates. *J. Multivariate Anal.* 98, 1051-1071.
- Müller, U. 2012. Estimating the density of a possibly missing response variable in nonlinear regression. *J. Statist. Plann. Inference*, 142, 1198-1214.
- Müller, U. 2009. Estimating linear functionals in nonlinear regression with responses missing at random. *Ann. Statist.* 37, 2245-2277.



- Parzen, E. 1962. On estimation of a probability density function and mode. *Ann. Statist.* 33, 1065-1076.
- Pollard, D. 1984. *Convergence of Stochastic Processes*. Springer, New York.
- Prakasa Rao, B.L.S. 1983. *Nonparametric functional estimation*. Academic Press, Orlando.
- Racine, J. and Hay  
eld, T. 2008. Nonparametric Econometrics: The np Package. *J. Statistical Software*, 27, 1-32.
- Robins, J., Rotnitzky, A. 1995. Semiparametric efficiency in multivariate regression models with missing data. *J. Amer. Statist. Assoc.*, 90, 122-129.
- Robins, J., Rotnitzky, A., and Zhao, Lue, P. 1994. Estimation of regression coefficients when some regressors are not always observed. *J. Amer. Statist. Assoc.* 89, 846-866.
- Rosenblatt, M. 1956. Remarks on some nonparametric estimates of a density function. *Ann. Math. Statist.* 27, 832-837.
- Rueda, M., Martínez, S., Martínez, H., and Arcos, A. (2006) Mean estimation with calibration techniques in presence of missing data. *Comput. Statist. Data Anal.* 50, 3263-3277.
- Talagrand, M. 1994. Sharper bounds for Gaussian and empirical processes. *Ann. Probab.* 22, 280-76.
- van de Geer, S. 2000. *Applications of empirical process theory*. Cambridge University Press, Cambridge.
- van der Vaart, A. and Wellner, J. 1996. *Weak convergence and empirical processes*. Springer-Verlag, New York.
- Vapnik, V. 1998. *Statistical learning theory*. John Wiley & Sons, Inc., New York.
- Wang, Q. 2008. Probability density estimation with data missing at random when covariables are present. *J. Statist. Plann. Inference*, 138, 568-587.
- Wang, Q. and Rao, J. 2002. Empirical likelihood-based inference under imputation for missing response data. *Ann. Statist.* 30, 890-924.
- Wang, Q., Linton, O., Härdle, W. 2004. Semiparametric regression analysis with missing response at random. *J. Amer. Statist. Assoc.*, 99, 334-345.
- Wang, Q. and Qin, Y. 2010. Empirical likelihood confidence bands for distribution functions with missing responses. *Journal of Statistical Planning and Inference*, 140, 2778-2789.
- Wang, L., Rotnitzky, A., Lin, X. 2010. Nonparametric regression with missing outcomes using weighted kernel estimating equations. *J. Amer. Statist. Assoc.*, 105, 1135-1146.
- Zou, Y., Liang, H., Zhang, J. 2015. Nonlinear wavelet density estimation with data missing at random when covariates are present. *Metrika*, 78, 967-995

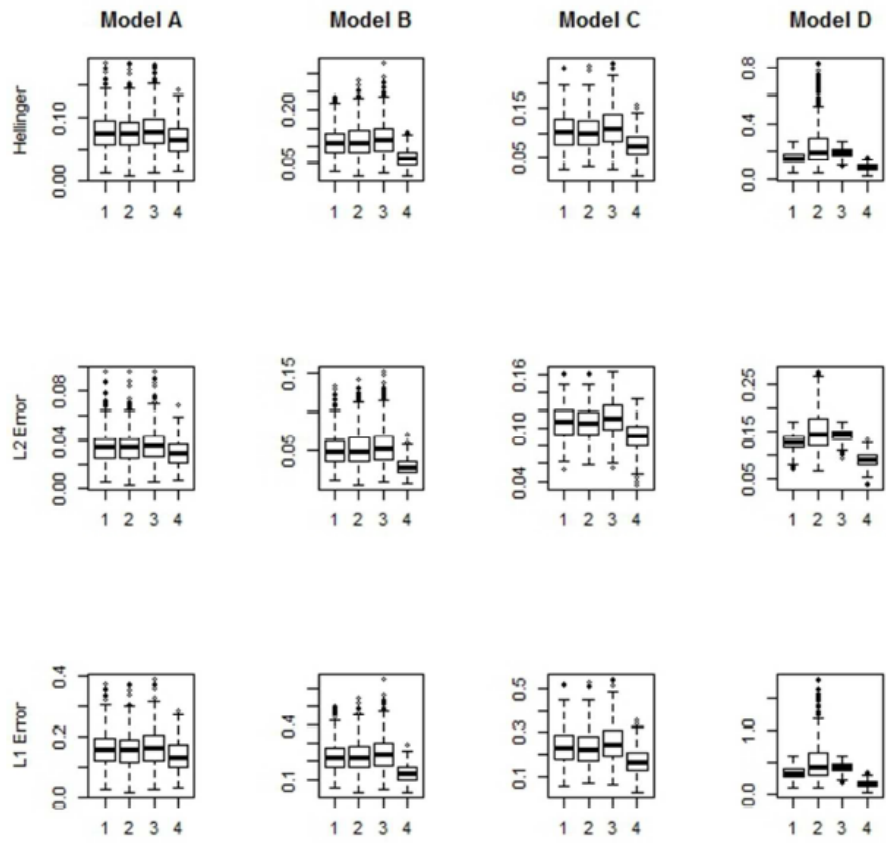
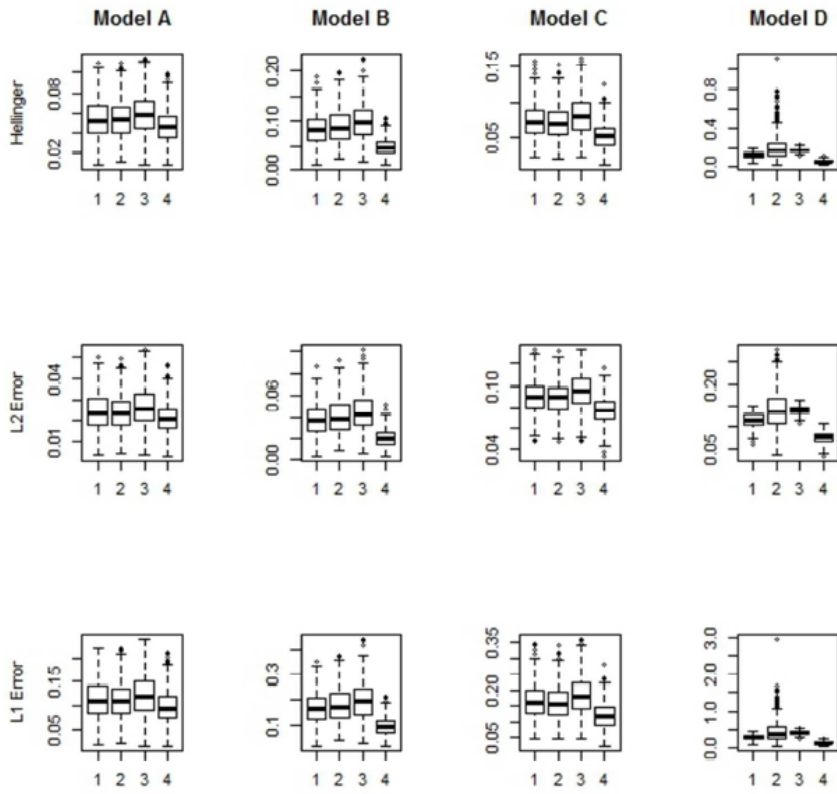
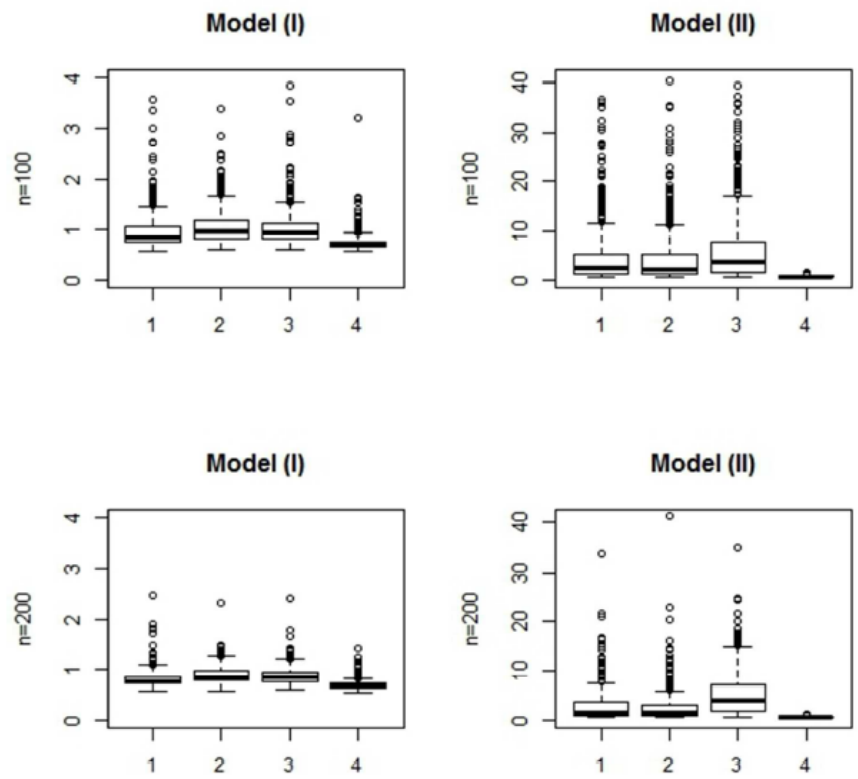


Figure 4:



Boxplots

Figure 5:



Boxplots for Example B

Figure 6: