

RESEARCH

Open Access



# Estimation of the mean of the partially linear single-index errors-in-variables model with missing response variables

Xin Qi<sup>1\*</sup> and ZhuoXi Yu<sup>2</sup>

\*Correspondence: [xinerqi@sina.com](mailto:xinerqi@sina.com)  
<sup>1</sup>Guangdong Polytechnic of Science and Technology, Zhuhai, P.R. China  
Full list of author information is available at the end of the article

## Abstract

In this paper, we estimate the mean of the partially linear single-index errors-in-variables model with missing response variables. The linear covariate is measured with additive error, therefore missing is not random. Two special estimators are defined that include a semiparametric regression imputation estimator and a marginal average estimator. These estimators are shown to be asymptotically normal and have the same asymptotic variance. A simulation experiment is used to illustrate our proposed method.

**Keywords:** Partially linear single-index model; Asymptotic normality; Missing response

## 1 Introduction

Semiparametric errors-in-variables (EV) models have attracted broad attention and have been deeply studied during the last two decades. Relevant studies include partially linear EV models (Liang et al. [5], He and Liang [4]), varying coefficient EV models (You et al. [16], Zhao and Xue [17]), partially linear varying coefficient EV models (You and Chen [15], Wei and Mei [13]), partially linear additive EV models (Wei et al. [11, 12]). Here, we consider the following partially linear single-index EV model:

$$\begin{cases} Y = g(Z^T \alpha) + X^T \beta + \varepsilon, \\ V = X + e, \end{cases} \quad (1.1)$$

where  $Y$  is a response variable, the single covariate  $Z \in \mathbb{R}^p$  is observed completely, the linear covariate  $X \in \mathbb{R}^q$  is observed with additive error, and only its substitute  $V$  can be observed;  $g(\cdot)$  is an unknown smooth link function,  $\varepsilon$  is the random error with  $E(\varepsilon|Z, X) = 0$ ,  $\text{Var}(\varepsilon|Z, X) < \infty$ ;  $(\alpha, \beta)$  is an unknown vector in  $\mathbb{R}^p \times \mathbb{R}^q$  with  $\|\alpha\| = 1$  which ensures identifiability, and the first nonzero component of  $\alpha$  is positive, where  $\|\cdot\|$  denotes the Euclidean norm. The measurement error  $e$  is independent of  $(Y, Z, X)$  with  $E(e) = 0$  and  $\text{Cov}(e) = \Sigma_e$ . Here, we assume that  $\Sigma_e$  is known. If it is unknown, the estimation method is analogous to the partial replication method of Liang et al. [5] in a partially linear EV model.

© The Author(s) 2020. This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

For the complete data set, the partially linear single-index EV model has been discussed by Liang and Wang [6] and Chen and Cui [1].

It is well known that the studies on the mean  $E(Y) = \theta$  are very important in regression models. If all the responses in the sample are available, the response variable mean can be usually obtained. However, in fact, some responses may be missing. This missing response problem may be caused by various reasons. For example, it may be too expensive to acquire the response  $Y$ 's and only part of  $Y$ 's are available. In practice, missing-data problems frequently occur in epidemiology studies, survey sampling, social science, and many other fields. Therefore, it is necessary to study the mean  $E(Y) = \theta$  based on the missing data set.

However, there's little research about the response variable mean in the partially linear single-index model. In this paper, we focus on the mean  $E(Y) = \theta$ , when there are missing responses in the partially linear single-index EV model (1.1). An indicator variable  $\delta$  is introduced in order to indicate whether an observation of  $Y$  is missing or observed, i.e.,  $\delta = 0$  indicates that  $Y$  is missing and  $\delta = 1$  indicates that  $Y$  is observed. Throughout this paper, if  $X$  is observable, we assume the data missing mechanism is as follows:

$$p(\delta = 1|Y, Z, X) = p(\delta = 1|Z, X) = \pi(Z, X)$$

for some unknown  $\pi(Z, X)$ . In addition, we also assume that the measurement error  $e$  is independent of  $\delta$ ,  $p(\delta = 1|Y, Z, X, V) = \pi(Z, X)$ . Since  $X$  is observed with measurement errors,  $Y$  is not missing at random if there are no further assumptions. The details can be seen in the paper of Liang et al. [7].

The imputation method is a common method of dealing with missing data, which fills in a plausible value for each missing data and then analyzes the result as if they were complete data. When some responses are missing, Cheng [2] applied kernel regression imputation to estimate  $\theta$  in a Nonparametric Model. Similar to the method of Cheng [2], Wei [10] estimated  $\theta$  in a partially linear varying-coefficient EV model with missing responses. In addition, the marginal average method also can be used in a missing data set in place of the imputation method. When some responses are missing in a partially linear model, Wang et al. [9] and Liang et al. [7] used the above two methods to estimate the mean of the responses with the covariates  $X$  being observed and not observed, respectively. In this paper, we extend the method in Liang et al. [7] to the partially linear single index EV models, propose two estimators of  $\theta$  in model (1.1) with missing response. The estimators are shown to be asymptotically normal and have the same asymptotic variance.

The rest of this paper is organized as follows. In Sect. 2, two estimators of  $\theta$  are proposed and a relative asymptotic result is presented. In Sect. 3, some simulation results are reported. All proofs are shown in Sect. 4.

## 2 Methodology and result

### 2.1 Estimation of the mean $E(Y) = \theta$

In order to derive the estimators of  $\theta$ , first we use the complete method of Qi and Wang [8] to estimate the regression coefficients, the single-index coefficients and the nonparametric function. By the least-squares method and the correction for attenuation technique, an

estimator of can be defined as

$$\hat{\beta}_n = \left\{ \frac{1}{n} \sum_{i=1}^n \delta_i [V_i - \hat{m}_V(Z_i^T \alpha)]^{\otimes 2} - \Sigma_e \right\}^{-1} \cdot \left\{ \frac{1}{n} \sum_{i=1}^n \delta_i [V_i - \hat{m}_V(Z_i^T \alpha)] [Y_i - \hat{m}_Y(Z_i^T \alpha)] \right\}, \tag{2.1}$$

where  $\hat{m}_Y(t) = \sum_{i=1}^n \frac{\delta_i K_{h_1}(Z_i^T \alpha - t)}{\sum_{i=1}^n \delta_i K_{h_1}(Z_i^T \alpha - t)} Y_i$  and  $\hat{m}_V(t) = \sum_{i=1}^n \frac{\delta_i K_{h_1}(Z_i^T \alpha - t)}{\sum_{i=1}^n \delta_i K_{h_1}(Z_i^T \alpha - t)} V_i$  are the estimators of  $m_Y(t) = \frac{E(\delta_i Y_i | Z_i^T \alpha = t)}{E(\delta_i | Z_i^T \alpha = t)}$  and  $m_V(t) = \frac{E(\delta_i V_i | Z_i^T \alpha = t)}{E(\delta_i | Z_i^T \alpha = t)}$ ,  $K_{h_1}(t) = \frac{K_1(\frac{t}{h_1})}{h_1}$ , with  $K_1(\cdot)$  being a kernel function and  $h_1$  being a suitable bandwidth.

After obtaining the estimator of  $\beta$ , we can obtain the estimators  $\hat{g}_n(\cdot)$  and  $\hat{g}'_n(\cdot)$  of  $g(\cdot)$  and  $g'(\cdot)$  for any fixed  $\alpha$ . By the locally linear method of Fan and Gijbels [3], we approximate  $g(t)$  within the neighborhood of  $t_0$ ,  $g(t) \approx g(t_0) + g'(t_0)(t - t_0)$  and minimize

$$\min_{g(t_0), g'(t_0)} \sum_{i=1}^n [Y_i - V_i^T \hat{\beta}_n - g(t_0) - g'(t_0)(t_i - t_0)]^2 K_{h_2}(t_i - t_0) \delta_i, \tag{2.2}$$

where  $K_{h_2}(t) = \frac{K_2(\frac{t}{h_2})}{h_2}$ , with  $K_2(\cdot)$  being a kernel function and  $h_2$  being a suitable bandwidth.

However, (2.1) and (2.2) cannot be applied directly in practice, since  $\alpha$  is unknown. So we need to estimate by minimizing

$$\min_{\alpha} \sum_{i=1}^n \delta_i [Y_i - V_i^T \hat{\beta}_n - \hat{g}_n(Z_i^T \alpha)]^2, \tag{2.3}$$

which yields, say  $\hat{\alpha}_n$ . Note that  $\hat{\beta}_n$  and  $\hat{g}_n(\cdot)$  can also be used to obtain  $\hat{\alpha}_n$  in (2.3). The complete estimation procedure is decomposed in an iterative process with the following steps:

Step 1. Acquire an initial value  $\hat{\alpha}_0$ , for example, by the method of Xia and Härdle [14], and let  $\hat{\alpha}_n = \frac{\hat{\alpha}_0}{\|\hat{\alpha}_0\|}$ .

Step 2. When  $\alpha = \hat{\alpha}_n$ , we can obtain  $\hat{\beta}_{nk}, \hat{g}_{nk}(\cdot)$  based on (2.1) and (2.2).

Step 3. The solution of (2.3) is denoted as  $\hat{\alpha}_{n(k+1)}$ . Let  $\hat{\alpha}_n = \frac{\hat{\alpha}_{n(k+1)}}{\|\hat{\alpha}_{n(k+1)}\|}$ .

Step 4. Iterate Steps 2 and 3 until convergence is achieved.

Next, we turn to estimate the mean  $E(Y) = \theta$ . Similar to Wang et al. [9] and Liang et al. [7], we construct two estimators of  $\theta$ . First, each missing  $Y_i$  is imputed by the estimated regression function  $V_i^T \hat{\beta}_n + \hat{g}_n(Z_i^T \hat{\alpha}_n)$ . Consequently, we obtain the semiparametric regression imputation estimator of  $\theta$ , which is designed as

$$\hat{\theta}_1 = n^{-1} \sum_{i=1}^n \delta_i Y_i + n^{-1} \sum_{i=1}^n (1 - \delta_i) [V_i^T \hat{\beta}_n + \hat{g}_n(Z_i^T \hat{\alpha}_n)].$$

Second, we only consider the sample average of the estimated regression function, that is, every  $Y_i$  is ignored. Accordingly, we get the marginal average estimator of  $\theta$ , which is defined as

$$\hat{\theta}_2 = n^{-1} \sum_{i=1}^n [V_i^T \hat{\beta}_n + \hat{g}_n(Z_i^T \hat{\alpha}_n)].$$

### 2.2 Asymptotic result

In this section, the asymptotic normality of  $\theta$ s will be summarized. And it will be shown that they are asymptotically equivalent.

For a concise representation, let  $P(t_0, \delta) = \delta/E(\delta|Z^T\alpha = t_0)$  and  $\tilde{S} = S - \frac{E(\delta S|Z^T\alpha)}{E(\delta|Z^T\alpha)}$ , for example,  $\tilde{X}_i = X_i - \frac{E(\delta X_i|Z^T\alpha)}{E(\delta|Z^T\alpha)}$ . Moreover, in order to state the asymptotic results, the following assumptions will be used:

- (C1) Let  $\Gamma_{\tilde{X}} = E\{\delta\tilde{X}^{\otimes 2}\}$ ,  $\Gamma_{\tilde{Z}} = E\{\delta[\tilde{Z}g'(Z^T\alpha)]^{\otimes 2}\}$  and  $\Gamma_{\tilde{Z}\tilde{X}} = E\{\delta\tilde{Z}\tilde{X}^Tg'(Z^T\alpha)\}$ .
- (C2) The bandwidth satisfies  $h_1 = h_0n^{-\frac{1}{p+4}}$  for some positive constant  $h_0$ ,  $\frac{nh_2^p}{\log n} \rightarrow \infty$ , where  $p$  is the dimension of  $Z$ .
- (C3) The kernels  $K_i(\cdot)$  ( $i = 1, 2$ ) are bounded symmetric density functions with compact support  $(-1, 1)$ , and they satisfy  $\int uK_i(u) du = 0, \int u^2K_i(u) du \neq 0$ .
- (C4) The density function  $f(t)$  of  $Z^T\alpha$  is bounded away from 0 and has two bounded derivatives on its support.
- (C5)  $g(\cdot), m_Y(\cdot), m_V(\cdot)$  have two bounded, continuous derivatives on their supports.
- (C6) The probability function  $\pi(Z, X)$  has bounded continuous second partial derivatives, and is bounded away from zero on the support of  $(Z, X)$ .
- (C7)  $E(|\varepsilon|^4 < \infty), E(|e|^3 < \infty)$ .

Now we give the following asymptotical result.

**Theorem 2.1** *Assume that conditions (C1)–(C7) are satisfied. Then we obtain*

$$\sqrt{n}(\hat{\theta}_i - \theta) \rightarrow N(0, \Theta_1 + \Theta_2), \quad i = 1, 2,$$

where  $\Theta_1 = E[P(Z^T\alpha, \delta)\varepsilon + [1 - P(Z^T\alpha, \delta)]e^T\beta + E[g'(Z^T\alpha)\tilde{Z}^T] \cdot \Gamma_{\tilde{Z}}^{-1}\delta g'(Z^T\alpha)\tilde{Z}(\varepsilon - e^T\beta) + E[\tilde{V}^T - g'(Z^T\alpha)\tilde{Z}^T\Gamma_{\tilde{Z}}^{-1}\Gamma_{\tilde{Z}\tilde{X}}] \cdot \Gamma_{\tilde{X}}^{-1}\{\delta[\tilde{V}(\varepsilon - e^T\beta) + \Sigma_{uu}\beta]\}^2]$  and  $\Theta_2 = E[X^T\beta + g(Z^T\alpha) - \theta]^2$ .

### 3 Simulation

In this section, we present a simulation study to analyze the finite sample performance of the regression imputation estimator  $\theta_1$  and the marginal average estimator  $\theta_2$ .

The simulation uses the partial linear single-index EV model (1.1) with a specific link function:

$$\begin{cases} Y = \sin(2\pi \cdot Z^T\alpha) + X^T\beta + \varepsilon, \\ V = X + e, \end{cases} \tag{3.1}$$

where  $X$  is generated from the standard normal distribution, trivariate  $Z$  is simulated from the uniform distribution  $U[0, 1]$ ,  $e$  is generated from the normal distribution  $N(0, 0.25^2)$ ,  $\varepsilon$  is simulated from the normal distribution with mean 0 and variance 0.01, and  $\alpha = (\frac{\sqrt{3}}{3}, \frac{\sqrt{3}}{3}, \frac{\sqrt{3}}{3})^T, \beta = 1$ . The kernel functions were taken to be  $K_i(t) = \frac{3}{4}(1 - t^2)^2$  if  $\|t\| \leq 1$ , and 0 otherwise,  $i = 1, 2$ .

The choices of bandwidths are quite crucial. In this paper, we use the least-squares delete-one cross-validation (CV) method to select bandwidths:  $\hat{h}_1$  and  $\hat{h}_2$  are chosen as

$$(\hat{h}_1, \hat{h}_2) = \arg \min_{h_1, h_2} \frac{1}{n} \sum_{i=1}^n \delta_i \{Y_i - V_i^T \hat{\beta}_n^{(-i)} - \hat{g}_n^{(-i)}(Z_i^T \hat{\alpha}_n^{(-i)})\}^2, \tag{3.2}$$

**Table 1** Biases and SE of  $\hat{\theta}_1, \hat{\theta}_2$  under different missing functions and different sample sizes

Missing rate	$n$	$\hat{\theta}_1$	$\hat{\theta}_2$
0.40	100	-0.0305 (0.1290)	0.0088 (0.1299)
	150	0.0101 (0.1088)	0.0072 (0.1079)
	200	0.0162 (0.0904)	0.0055 (0.0904)
0.30	100	-0.0234 (0.1270)	0.0109 (0.1282)
	150	0.0228 (0.1078)	0.0063 (0.1065)
	200	0.0166 (0.0892)	0.0064 (0.0893)
0.20	100	-0.0394 (0.1247)	0.0122 (0.1272)
	150	0.0141 (0.1055)	0.0068 (0.1056)
	200	0.0120 (0.0875)	0.0069 (0.0886)
0.10	100	0.0123 (0.1241)	0.0136 (0.1261)
	150	0.0107 (0.1038)	0.0073 (0.1047)
	200	0.0089 (0.0871)	0.0074 (0.0879)

where  $\hat{\beta}_n^{(-i)}, \hat{g}_n^{(-i)}$  and  $\hat{\alpha}_n^{(-i)}$  are the “leave-one-out” versions of  $\hat{\beta}_n, \hat{g}_n$  and  $\hat{\alpha}_n$ , respectively. However, the  $h_i, i = 1, 2$  from (3.2) may not be the optimal bandwidths because they may not satisfy the conditions imposed in the theorems. According to their conditions, the optimal bandwidth according to (3.2) is to choose a constant  $h_0$ .

Based on model (3.1), we considered the following four response probabilities of missing, namely:

Case 1:  $P(\delta = 1|Z = z, X = x) = \frac{\exp(0.6+z^T\phi+\varphi x)}{1+\exp(0.6+z^T\phi+\varphi x)}$ , where  $\phi = (-0.12, -0.012, -0.12)^T, \varphi = 0.35$ ;

Case 2:  $P(\delta = 1|Z = z, X = x) = \frac{\exp(0.6+z^T\phi+\varphi x)}{1+\exp(0.6+z^T\phi+\varphi x)}$ , where  $\phi = (0.2, 0.2, 0.2)^T, \varphi = 0.45$ ;

Case 3:  $P(\delta = 1|Z = z, X = x) = \frac{\exp(0.6+z^T\phi+\varphi x)}{1+\exp(0.6+z^T\phi+\varphi x)}$ , where  $\phi = (0.65, 0.65, 0.65)^T, \varphi = 0.8$ ;

Case 4:  $P(\delta = 1|Z = z, X = x) = 0.9$  for all  $z$  and  $x$ . The average missing rates were 0.4, 0.3, 0.2, and 0.1, respectively. From the 1000 simulated values of  $\hat{\theta}_1, \hat{\theta}_2$ , we calculated the biases and standard errors (SE) of the two estimators. The simulated results are reported in Table 1.

From Table 1, we observe that

- (a) Biases and SE decrease as  $n$  increases for every fixed missing rate. Also, SE increase as the missing rate increases for every fixed sample size  $n$ .
- (b) The SE of  $\hat{\theta}_1, \hat{\theta}_2$  are nearly the same for every fixed missing rate and sample size.

#### 4 Proof of the main result

In order to prove the main result, we first give some lemmas.

**Lemma 4.1** *Under conditions (C1)–(C7), we have*

$$\begin{aligned}
 & \hat{g}_n(t_0, \hat{\alpha}_n, \hat{\beta}_n) - g(t_0) \\
 &= \frac{1}{n} \cdot \frac{1}{f(t_0)E(\delta|Z^T\alpha = t_0)} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^T\alpha - t_0) (\varepsilon_i - e_i^T\beta) \\
 & \quad - (\hat{\beta}_n - \beta)^T \frac{E(\delta V|Z^T\alpha = t_0)}{E(\delta|Z^T\alpha = t_0)} - (\hat{\alpha}_n - \alpha)^T \frac{E(\delta Z g'(Z^T\alpha)|Z^T\alpha = t_0)}{E(\delta|Z^T\alpha = t_0)} \\
 & \quad + o_p\left(\frac{1}{\sqrt{n}}\right) + O_p(h_2^2). \tag{4.1}
 \end{aligned}$$

*Proof of Lemma 4.1* When  $\alpha = \hat{\alpha}_n$ , the estimators of  $g(\cdot)$  and  $g'(\cdot)$  can be obtained from (2.2). By a straightforward calculation,

$$0 = \frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \hat{\alpha}_n - t_0) \begin{pmatrix} 1 \\ Z_i^\top \hat{\alpha}_n - t_0 \end{pmatrix} \cdot [Y_i - V_i^\top \hat{\beta}_n - \hat{g}_n(t_0) - \hat{g}'_n(t_0)(Z_i^\top \hat{\alpha}_n - t_0)].$$

Then focusing on the top equation only and using Taylor expansion, we have

$$\begin{aligned} 0 &= \frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) [Y_i - V_i^\top \beta - g(Z_i^\top \alpha)] \\ &\quad - \frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) [\hat{g}_n(t_0) - g(t_0)] \\ &\quad - (\hat{\beta}_n - \beta)^\top \frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) V_i \\ &\quad - (\hat{\alpha}_n - \alpha)^\top \frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) Z_i g'(t_0) + o_p\left(\frac{1}{\sqrt{n}}\right) + O_p(h_2^2), \end{aligned}$$

that is,

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) [\hat{g}_n(t_0) - g(t_0)] \\ &= \frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) (\varepsilon_i - e_i^\top \beta) - (\hat{\beta}_n - \beta)^\top \frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) V_i \\ &\quad - (\hat{\alpha}_n - \alpha)^\top \frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) Z_i g'(t_0) + o_p\left(\frac{1}{\sqrt{n}}\right) + O_p(h_2^2). \tag{4.2} \end{aligned}$$

Note that  $\frac{1}{n} \sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0) = f(t_0) + o_p(1)$ . Dividing all terms in (4.2) by  $\frac{1}{n} \times \sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0)$ , we obtain

$$\begin{aligned} &[\hat{g}_n(t_0) - g(t_0)] \frac{\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0)}{\frac{1}{n} \sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0)} \\ &= \frac{1}{n} \frac{\sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) (\varepsilon_i - e_i^\top \beta)}{\sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0)} \\ &\quad - (\hat{\beta}_n - \beta)^\top \frac{\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) V_i}{\sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0)} \\ &\quad - (\hat{\alpha}_n - \alpha)^\top \frac{\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) Z_i g'(t_0)}{\sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0)} + o_p\left(\frac{1}{\sqrt{n}}\right) + O_p(h_2^2). \end{aligned}$$

Noting that  $\frac{\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0)}{\frac{1}{n} \sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0)} = E(\delta | Z_i^\top \alpha = t_0) + o_p(1)$ , we get

$$\begin{aligned} & \frac{\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) V_i}{\frac{1}{n} \sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0)} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) V_i}{\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0)} \times \frac{\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0)}{\frac{1}{n} \sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0)} \\ &= E(\delta V | Z^\top \alpha = t_0) (1 + o_p(1)) \\ &= E(\delta X | Z^\top \alpha = t_0) (1 + o_p(1)). \end{aligned}$$

Similarly, we also have

$$\frac{\frac{1}{n} \sum_{i=1}^n \delta_i K_{h_2}(Z_i^\top \alpha - t_0) Z_i g'(t_0)}{\frac{1}{n} \sum_{i=1}^n K_{h_2}(Z_i^\top \alpha - t_0)} = E(\delta Z g'(t_0) | Z^\top \alpha = t_0) (1 + o_p(1)).$$

Thus we get equation (4.1). □

**Lemma 4.2** *Under conditions (C1)–(C7), we have*

$$\hat{\alpha}_n - \alpha = \frac{1}{n} \Gamma_{\tilde{Z}}^{-1} \sum_{i=1}^n \delta_i g'(Z_i^\top \alpha) \tilde{Z}_i (\varepsilon_i - e_i^\top \beta) - \Gamma_{\tilde{Z}}^{-1} \Gamma_{\tilde{Z}\tilde{X}} (\hat{\beta}_n - \beta) + o_p\left(\frac{1}{\sqrt{n}}\right).$$

*Proof of Lemma 4.2* This proof is given in Qi and Wang [8], we omit the details here. □

**Lemma 4.3** *Under conditions (C1)–(C7), we have*

$$\hat{\beta}_n - \beta = \frac{1}{n} \Gamma_{\tilde{X}}^{-1} \sum_{i=1}^n \{ \delta_i [\tilde{V}_i (\varepsilon_i - e_i^\top \beta) + \Sigma_{uu} \beta] \} + o_p\left(\frac{1}{\sqrt{n}}\right).$$

*Proof of Lemma 4.3* The proof of Lemma 4.3 is similar to the proof of Theorem 1 by Liang et al. [7], we omit the details here. □

*Proof of Theorem 2.1* Here we only consider the asymptotic normality of  $\theta_1$ . The asymptotic result for  $\theta_2$  is obtained similarly. □

For  $\theta_1$ , we have

$$\begin{aligned} \hat{\theta}_1 &= n^{-1} \sum_{i=1}^n \delta_i \varepsilon_i + n^{-1} \sum_{i=1}^n [1 - \delta_i] e_i^\top \beta \\ &\quad + n^{-1} \sum_{i=1}^n [X_i^\top \beta + g(Z_i^\top \alpha)] + n^{-1} \sum_{i=1}^n [1 - \delta_i] V_i^\top (\hat{\beta}_n - \beta) \\ &\quad + n^{-1} \sum_{i=1}^n [1 - \delta_i] [\hat{g}_n(Z_i^\top \hat{\alpha}_n) - g(Z_i^\top \alpha)] \\ &=: \sum_{i=1}^5 I_i + o_p\left(\frac{1}{\sqrt{n}}\right), \end{aligned} \tag{4.3}$$

where

$$\begin{aligned}
 I_1 &= n^{-1} \sum_{i=1}^n \delta_i \varepsilon_i, \\
 I_2 &= n^{-1} \sum_{i=1}^n [1 - \delta_i] e_i^\top \boldsymbol{\beta}, \\
 I_3 &= n^{-1} \sum_{i=1}^n [X_i^\top \boldsymbol{\beta} + g(Z_i^\top \boldsymbol{\alpha})], \\
 I_4 &= n^{-1} \sum_{i=1}^n [1 - \delta_i] V_i^\top (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}), \\
 I_5 &= n^{-1} \sum_{i=1}^n [1 - \delta_i] [\hat{g}_n(Z_i^\top \hat{\boldsymbol{\alpha}}_n) - g(Z_i^\top \boldsymbol{\alpha})].
 \end{aligned}$$

From Taylor expansion and the continuity of  $g'(\cdot)$ , we obtain that

$$\begin{aligned}
 &\hat{g}_n(Z_i^\top \hat{\boldsymbol{\alpha}}_n) - g(Z_i^\top \boldsymbol{\alpha}) \\
 &= [\hat{g}_n(Z_i^\top \boldsymbol{\alpha}) - g(Z_i^\top \boldsymbol{\alpha})] + g'(Z_i^\top \boldsymbol{\alpha}) Z_i^\top (\hat{\boldsymbol{\alpha}}_n - \boldsymbol{\alpha}) + o_p\left(\frac{1}{\sqrt{n}}\right).
 \end{aligned} \tag{4.4}$$

By Lemma 4.1 and (4.4), it is easy to get

$$\begin{aligned}
 I_5 &= \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] \cdot \frac{\sum_{j=1}^n \delta_j K_{h_2}(Z_j^\top \boldsymbol{\alpha} - Z_i^\top \boldsymbol{\alpha})(\varepsilon_j - e_j^\top \boldsymbol{\beta})}{nf(Z_i^\top \boldsymbol{\alpha})E(\delta|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} \\
 &\quad - (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta})^\top \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] \frac{E(\delta V|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})}{E(\delta|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} \\
 &\quad + (\hat{\boldsymbol{\alpha}}_n - \boldsymbol{\alpha})^\top \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] \left[ g'(Z_i^\top \boldsymbol{\alpha}) Z_i - \frac{E(\delta Z g'(Z^\top \boldsymbol{\alpha})|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})}{E(\delta|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} \right] \\
 &\quad + o_p\left(\frac{1}{\sqrt{n}}\right) + O_p(h_2^2) \\
 &=: I_{51} - I_{52} + I_{53} + o_p\left(\frac{1}{\sqrt{n}}\right) + O_p(h_2^2),
 \end{aligned} \tag{4.5}$$

where

$$\begin{aligned}
 I_{51} &= \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] \cdot \frac{\sum_{j=1}^n \delta_j K_{h_2}(Z_j^\top \boldsymbol{\alpha} - Z_i^\top \boldsymbol{\alpha})(\varepsilon_j - e_j^\top \boldsymbol{\beta})}{nf(Z_i^\top \boldsymbol{\alpha})E(\delta|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})}, \\
 I_{52} &= (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta})^\top \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] \frac{E(\delta V|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})}{E(\delta|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})}, \\
 I_{53} &= (\hat{\boldsymbol{\alpha}}_n - \boldsymbol{\alpha})^\top \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] \left[ g'(Z_i^\top \boldsymbol{\alpha}) Z_i - \frac{E(\delta Z g'(Z^\top \boldsymbol{\alpha})|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})}{E(\delta|Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} \right].
 \end{aligned}$$



We have

$$\begin{aligned}
 I_{51} &= \frac{1}{n} \sum_{j=1}^n \delta_j (\varepsilon_j - e_j^\top \boldsymbol{\beta}) \cdot \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] \frac{K_{h_2}(Z_j^\top \boldsymbol{\alpha} - Z_i^\top \boldsymbol{\alpha})}{f(Z_i^\top \boldsymbol{\alpha}) E(\delta | Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} \\
 &= \frac{1}{n} \sum_{i=1}^n \delta_i (\varepsilon_i - e_i^\top \boldsymbol{\beta}) \cdot \left[ \frac{1}{E(\delta | Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} - 1 \right].
 \end{aligned} \tag{4.6}$$

Combining Lemma 4.2 and calculating directly, we can get

$$\begin{aligned}
 I_{53} &= \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] g'(Z_i^\top \boldsymbol{\alpha}) \tilde{Z}_i^\top \cdot \Gamma_{\tilde{Z}}^{-1} \frac{1}{n} \sum_{i=1}^n \delta_i g'(Z_i^\top \boldsymbol{\alpha}) \tilde{Z}_i (\varepsilon_i - e_i^\top \boldsymbol{\beta}) \\
 &\quad - \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] g'(Z_i^\top \boldsymbol{\alpha}) \tilde{Z}_i^\top \cdot \Gamma_{\tilde{Z}}^{-1} \Gamma_{\tilde{Z}\tilde{X}} (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) + o_p\left(\frac{1}{\sqrt{n}}\right) \\
 &=: I_{531} - I_{532} + o_p\left(\frac{1}{\sqrt{n}}\right),
 \end{aligned} \tag{4.7}$$

where

$$\begin{aligned}
 I_{531} &= \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] g'(Z_i^\top \boldsymbol{\alpha}) \tilde{Z}_i^\top \cdot \Gamma_{\tilde{Z}}^{-1} \frac{1}{n} \sum_{i=1}^n \delta_i g'(Z_i^\top \boldsymbol{\alpha}) \tilde{Z}_i (\varepsilon_i - e_i^\top \boldsymbol{\beta}), \\
 I_{532} &= \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] g'(Z_i^\top \boldsymbol{\alpha}) \tilde{Z}_i^\top \cdot \Gamma_{\tilde{Z}}^{-1} \Gamma_{\tilde{Z}\tilde{X}} (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) + o_p\left(\frac{1}{\sqrt{n}}\right).
 \end{aligned}$$

By a straightforward calculation, it follows that

$$\begin{aligned}
 I_4 - I_{52} - I_{532} &= \left[ \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] \tilde{V}_i^\top - \frac{1}{n} \sum_{i=1}^n [1 - \delta_i] g'(Z_i^\top \boldsymbol{\alpha}) \tilde{Z}_i^\top \Gamma_{\tilde{Z}}^{-1} \Gamma_{\tilde{Z}\tilde{X}} \right] \\
 &\quad \cdot (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) + o_p\left(\frac{1}{\sqrt{n}}\right), \\
 &= (E[1 - \delta] \tilde{V}^\top - E[(1 - \delta) g'(Z^\top \boldsymbol{\alpha}) \tilde{Z}^\top] \Gamma_{\tilde{Z}}^{-1} \Gamma_{\tilde{Z}\tilde{X}}) \\
 &\quad \cdot (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) + o_p\left(\frac{1}{\sqrt{n}}\right),
 \end{aligned} \tag{4.8}$$

Furthermore, it is easy to get

$$\begin{aligned}
 I_1 + I_2 + I_{51} &= \frac{1}{n} \sum_{i=1}^n \frac{\delta_i \varepsilon_i}{E(\delta | Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} + \frac{1}{n} \sum_{i=1}^n \left[ 1 - \frac{\delta_i}{E(\delta | Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} \right] e_i^\top \boldsymbol{\beta}.
 \end{aligned} \tag{4.9}$$

Combining (4.3), (4.5), (4.6), (4.7), (4.8), (4.9), and Lemma 4.3, one can get

$$\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta} = n^{-1} \sum_{i=1}^n \frac{\delta_i}{E(\delta | Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} \varepsilon_i + n^{-1} \sum_{i=1}^n \left[ 1 - \frac{\delta_i}{E(\delta | Z^\top \boldsymbol{\alpha} = Z_i^\top \boldsymbol{\alpha})} \right] e_i^\top \boldsymbol{\beta}$$

$$\begin{aligned}
 &+ n^{-1} \sum_{i=1}^n [X_i^T \boldsymbol{\beta} + g(Z_i^T \boldsymbol{\alpha}) - \boldsymbol{\theta}] \\
 &+ E[g'(Z^T \boldsymbol{\alpha}) \tilde{Z}^T] \cdot \Gamma_{\tilde{Z}}^{-1} \frac{1}{n} \sum_{i=1}^n \delta_i g'(Z_i^T \boldsymbol{\alpha}) \tilde{Z}_i (\varepsilon_i - e_i^T \boldsymbol{\beta}) \\
 &+ E[\tilde{V}^T - g'(Z^T \boldsymbol{\alpha}) \tilde{Z}^T \Gamma_{\tilde{Z}}^{-1} \Gamma_{\tilde{Z}\tilde{X}}] \cdot \frac{1}{n} \Gamma_{\tilde{X}}^{-1} \sum_{i=1}^n \{\delta_i [\tilde{V}_i (\varepsilon_i - e_i^T \boldsymbol{\beta}) + \Sigma_{uu} \boldsymbol{\beta}]\} \\
 &+ o_p\left(\frac{1}{\sqrt{n}}\right) + O_p(h_2^2). \tag{4.10}
 \end{aligned}$$

This, together with the central limit theorem, proves Theorem 2.1 for  $\hat{\boldsymbol{\theta}}_1$ .

**Acknowledgements**

The authors thank the two referees and editor(s) for carefully reading the paper and for their valuable suggestions and comments which greatly improved the paper.

**Funding**

This work is supported by Philosophy and Social Sciences Planning Project of Guangdong Province during the ‘13th Five-Year’ Plan Period (No. GD18CYJ08, GD18XGL26), National Social Science Foundation of China (18CTQ032) and Guangdong Polytechnic of Science and Technology Research Project (No. XJPY2018006, XJMS2018006).

**Availability of data and materials**

The data sets analyzed in the current study can be generated by Monte Carlo experiments.

**Competing interests**

The authors declare that they have no competing interests.

**Authors’ contributions**

The authors contributed equally to the writing of this paper. All authors read and approved the final manuscript.

**Author details**

<sup>1</sup>Guangdong Polytechnic of Science and Technology, Zhuhai, P.R. China. <sup>2</sup>School of Economics, Liaoning University, Shenyang, P.R. China.

**Publisher’s Note**

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Received: 4 September 2019 Accepted: 23 January 2020 Published online: 30 January 2020

**References**

1. Chen, X., Cui, H.J.: Empirical likelihood for partially linear single-index errors-in-variables model. *Commun. Stat., Theory Methods* **38**(15), 2498–2514 (2009)
2. Cheng, P.E.: Nonparametric estimation of mean functionals with data missing at random. *J. Am. Stat. Assoc.* **89**, 81–87 (1994)
3. Fan, J.Q., Gijbels, I.: *Local Polynomial Modelling and Its Applications*. Chapman & Hall, London (1996)
4. He, X.M., Liang, H.: Quantile regression estimates for a class of linear and partially linear errors-in-variables models. *Stat. Sin.* **10**, 129–140 (2000)
5. Liang, H., Härdle, W., Carroll, R.J.: Estimation in a semiparametric partially linear errors-in-variables model. *Ann. Stat.* **27**, 1519–1535 (1999)
6. Liang, H., Wang, N.: Partially linear single-index measurement error models. *Stat. Sin.* **15**, 99–116 (2005)
7. Liang, H., Wang, S.J., Carroll, R.J.: Partially linear models with missing response variables and error-prone covariates. *Biometrika* **94**, 185–198 (2007)
8. Qi, X., Wang, D.H.: Estimation in a partially linear single-index model with missing response variables and error-prone covariates. *J. Inequal. Appl.* **2016**, 11 (2016). <https://doi.org/10.1186/s13660-015-0941-8>
9. Wang, Q.H., Linton, O., Härdle, W.: Semiparametric regression analysis with missing response at random. *J. Am. Stat. Assoc.* **99**, 334–345 (2004)
10. Wei, C.H.: Estimation in partially linear varying-coefficient errors-in-variables models with missing responses (Chinese ed.). *Acta Math. Sci.* **30**, 1042–1054 (2010)
11. Wei, C.H., Jia, X.J., Hu, H.S.: Statistical inference on partially linear additive models with missing response variables and error-prone covariates. *Commun. Stat., Theory Methods* **44**, 872–883 (2015)
12. Wei, C.H., Luo, Y.B., Wu, X.Z.: Empirical likelihood for partially linear additive errors-in-variables models. *Stat. Pap.* **53**(2), 485–496 (2012)
13. Wei, C.H., Mei, C.L.: Empirical likelihood for partially linear varying-coefficient models with missing response variables and error-prone covariates. *J. Korean Stat. Soc.* **41**, 97–103 (2012)

14. Xia, Y.C., Härdle, W.: Semi-parametric estimation of partially linear single-index models. *J. Multivar. Anal.* **97**, 1162–1184 (2006)
15. You, J.H., Chen, G.M.: Estimation of a semiparametric varying-coefficient partially linear errors-in-variables model. *J. Multivar. Anal.* **97**(2), 324–341 (2006)
16. You, J.H., Zhou, Y., Chen, G.M.: Corrected local polynomial estimation in varying-coefficient models with measurement errors. *Can. J. Stat.* **34**(3), 391–410 (2006)
17. Zhao, P.X., Xue, L.G.: Variable selection for varying coefficient models with measurement errors. *Metrika* **74**, 231–245 (2011)

**Submit your manuscript to a SpringerOpen<sup>®</sup> journal and benefit from:**

- ▶ Convenient online submission
- ▶ Rigorous peer review
- ▶ Open access: articles freely available online
- ▶ High visibility within the field
- ▶ Retaining the copyright to your article

---

Submit your next manuscript at ▶ [springeropen.com](https://www.springeropen.com)

---