# Estimation of Partial Linear Error-in-Variables Models under Martingale Difference Sequence

**Zhuoxi YU**[1,*], **Dehui WANG**[2], **Na HUANG**[3,1]

1. *School of Management Science and Information Engineering, Jilin University of Finance and Economics, Jilin* 130117, *P. R. China;*
2. *Department of Statistic, College of Mathematics, Jilin University, Jilin* 130021, *P. R. China;*
3. *School of Information Management and Engineering, Shanghai University of Finance and Economics, Shanghai* 200433, *P. R. China*

**Abstract** Consider the partly linear model $Y = x\beta + g(t) + e$ where the explanatory $x$ is erroneously measured, and both $t$ and the response $Y$ are measured exactly, the random error $e$ is a martingale difference sequence. Let $\widetilde{x}$ be a surrogate variable observed instead of the true $x$ in the primary survey data. Assume that in addition to the primary data set containing $N$ observations of $\{(Y_j, \widetilde{x}_j, t_j)_{j=n+1}^{n+N}\}$, the independent validation data containing $n$ observations of $\{(\widetilde{x}_j, x_j, t_j)_{j=1}^{n}\}$ is available. In this paper, a semiparametric method with the primary data is employed to obtain the estimator of $\beta$ and $g(\cdot)$ based on the least squares criterion with the help of validation data. The proposed estimators are proved to be strongly consistent. Finite sample behavior of the estimators is investigated via simulations too.

**Keywords** partial linear error-in-variables models; martingale difference sequence; validation data; strong consistency

**MR(2010) Subject Classification** 62J02; 62G05

## 1. Introduction

Consider the partial linear regression model:

$$Y = x\beta + g(t) + e, \tag{1.1}$$

where $(x, t) \in R \times [0, 1]$ are the nonrandom design points, $\beta$ is an unknown parameter and $g(\cdot)$ is an unknown smooth function defined on $[0, 1]$ and $e = Y - x\beta - g(t)$ are the random errors.

Model (1.1) was first introduced by Engle et al. in [6] and has been widely studied in the literature and the majority of the work was done so far assuming that the random errors are independent and identically distributed (i.i.d), see [4,10,12,20] and the monograph of Hädle et al. in [11]. However, the independence assumption for the errors is not always appropriate in applications, many authors investigated the model (1.1) when the errors are dependent, such

as, when the errors are assumed to be negatively associated random variable, the investigation related to the model (1.1) can be found in [13,14]. Among other works [1,8,9,22]. When the error is a sequence of martingale difference [5,7,16]. Heteroscedasticity also has been applied to the model (1.1) by many researchers [15,25].

In many research settings, the exact measurement of some important variable is difficult, time consuming, or expensive, and can only be performed for a few items in a large scale study. Thus, they are usually replaced by surrogate observations, which are available by some relative simple measuring methods. Some statisticians developed statistical inference techniques based on surrogate data [17,18,23,24]. However, there are few studies for the statistical inference based on surrogate data when the errors are dependent.

For model (1.1), if the $x$ is mismeasured, it is well known that ignoring the measurement error and naively performing a usual regression analysis may lead to incorrect conclusions. Let $\widetilde{x}$ be a surrogate variable observed instead of the true $x$ in the primary survey data. Assume that in addition to the primary data set containing $N$ observations of $\{(Y_j, \widetilde{x}_j, t_j)_{j=n+1}^{n+N}\}$, an independent validation data containing $n$ observations of $\{(x_j, \widetilde{x}_j, t_j)_{j=1}^{n}\}$ is available. Cai [2] gave a consistent estimator in model (1.1) when there are measurement errors in variables with $\rho$-mixing random errors.

The main purpose of this paper is to show how we can obtain a consistent estimator when there are measurement errors in variables with martingale difference errors, that is, the $x$ is mismeasured and the random errors $\{e_i\}$ in (1.1) is a martingale difference sequence with respect to an increase sequence of $\sigma$-fields $\{\mathcal{F}_i\}$; i.e., $e_i$ is $\mathcal{F}_i$-measurable and $E(e_i|\mathcal{F}_{i-1}) = 0$. Further, assume that $E(e_i^2|\mathcal{F}_{i-1}) = \sigma^2$ and $0 < \sigma^2 < \infty$ is unknown.

## 2. Methodology and main results

By employing the validation data, the estimator of the true $x_i$ can be defined as

$$u_n(\nu_i) = \frac{\sum_{j=1}^{n} x_j K_1((\nu_j - \nu_i)/b_n)}{\sum_{j=1}^{n} K_1((\nu_j - \nu_i)/b_n)}, \quad i = n+1, \ldots, n+N, \tag{2.1}$$

where $\nu_j = (\widetilde{x}_j, t_j)$, $j = 1, \ldots, n+N$, $K_1(\cdot)$ is a two-dimensional kernel function and $b_n$ is a bandwith tending to zero as $n \to \infty$. Let

$$\omega_{Ni}(t) = \frac{K_2((t - t_i)/h_N)}{\sum_{i=n+1}^{n+N} K_2((t - t_i)/h_N)}, \quad i = n+1, \ldots, n+N, \tag{2.2}$$

where $K_2(\cdot)$ is also a kernel function and $h_N$ is a bandwith tending to zero as $N \to \infty$.

If $\beta$ is known to be the true parameter in (1.1), then $Y_i - x_i\beta = g(t_i) + e_i$. Hence, the natural estimator of $g(\cdot)$ is

$$g_{nN}(t, \beta) = \sum_{i=n+1}^{n+N} \omega_{Ni}(t)(Y_i - u_n(\nu_i)\beta) =: g_{1N}(t) - g_{2N}(t)\beta. \tag{2.3}$$

The least square estimator, say $\beta_{nN}$, of $\beta$ is defined by

$$\sum_{i=n+1}^{n+N} (Y_i - u_n(\nu_i)\beta - g_{1N}(t_i) + g_{2N}(t_i)\beta)^2 = \min.$$

Solving the equation

$$\sum_{i=n+1}^{n+N} [(u_n(\nu_i) - g_{2N}(t_i))(Y_i - g_{1N}(t_i) - [u_n(\nu_i) - g_{2N}(t_i)]\beta)] = 0,$$

we have

$$\beta_{nN} = \sum_{i=n+1}^{n+N} \overline{u}_n(\nu_i)\overline{Y}_i / S_{nN}^2, \tag{2.4}$$

where

$$\overline{u}_n(\nu_i) = u_n(\nu_i) - \sum_{j=n+1}^{n+N} \omega_{Nj}(t_i)u_n(\nu_j) = u_n(\nu_i) - g_{2N}(t_i),$$

$$\overline{Y}_i = Y_i - \sum_{j=n+1}^{n+N} \omega_{Nj}(t_i)Y_j = Y_i - g_{1N}(t_i),$$

$$S_{nj}^2 = \sum_{i=n+1}^{n+j} [\overline{u}_n(\nu_i)]^2, \quad j = 1, \ldots, N.$$

(2.3), (2.4) together then yield the final estimator of $g(\cdot)$, as follows

$$\widetilde{g}_{nN}(t) =: g_{1N}(t) - g_{2N}(t)\beta_{nN}. \tag{2.5}$$

Now the model (1.1) can be rewritten as

$$Y_i = u_n(\nu_i)\beta + g(t_i) + \varepsilon_i, \quad \varepsilon_i = x_i\beta + e_i - u_n(\nu_i)\beta, \tag{2.6}$$

where $i = n + 1, \ldots, n + N$. $e_i = Y_i - x_i\beta - g(t_i)$ is a martingale difference sequence errors, $\varepsilon_i - E(\varepsilon_i|\mathcal{F}_{i-1}) = e_i$ and $E[(\varepsilon_i - E(\varepsilon_i|\mathcal{F}_{i-1}))^2|\mathcal{F}_{i-1}] = E(e_i^2|\mathcal{F}_{i-1}) = \sigma^2 < +\infty$.

Throughout this paper, $C$ and $C_k, k = 1, 2, \ldots$ will represent positive constants. Let $\nu = (\widetilde{x}, t)$, $\| \nu \| = |\widetilde{x}| + |t|$.

In order to obtain the main results in this section, we shall need the following assumptions:

(A1) $C_1 I(\| \nu \| \leq C) \leq K_1(\nu) \leq C_2, C_3 I(|t| \leq C) \leq K_2(t) \leq C_4$.

(A2) $\lim_{\|\nu\|\to\infty} K_1(\nu) = 0, \lim_{|t|\to\infty} K_2(t) = 0$.

(A3) $n^{-1} \sum_{i=1}^{n} |x_i| \leq C$, $N^{-1} \sum_{i=n+1}^{n+N} |\widetilde{x}_i| \leq C$.

(A4) $\sum_{i=n+1}^{n+N} |\overline{u}_n(\nu_i)| \leq C \sum_{i=n+1}^{n+N} [\overline{u}_n(\nu_i)]^2$.

(A5) $\exists N_1$, as $N > N_1, C_5 \leq \frac{S_{nN}^2}{N} \leq C_6$.

(A6) $g(t)$ satisfies the Lipschitz condition of order 1 on the interval $[0, 1]$.

(A7) $\exists \lambda > 1$, such that $N/n \longrightarrow \lambda$.

**Remark 2.1** By assumption (A1) and (A3), it is easy to obtain that

$$\max_{n+1\leq i,j\leq n+N} |\omega_{Nj}(t_i)| \leq CN^{-1},$$

$$\max_{n+1\leq i\leq n+N} \sum_{j=n+1}^{n+N} |\omega_{Nj}(t_i)| \leq C,$$

$$\max_{n+1\leq i\leq n+N} |\overline{u}_n(\nu_i)| \leq C\sum_{j=1}^{n} \frac{|x_j|}{n} \leq C_7.$$

**Theorem 2.2** ([2])   *Under the conditions (A1)–(A6), if there exists $\alpha > 2$ such that*

$$\sup_{i\geq 1} \|E(|e_i|^\alpha|\mathcal{F}_{i-1})\| \leq C < \infty,$$

*then*

$$\lim_{n,N\to\infty} \beta_{n,N} = \beta \quad a.s.. \tag{2.7}$$

$$\lim_{n,N\to\infty} \sup_{t\in[0,1]} |\widetilde{g}_{nN}(t) - g(t)| = 0 \quad a.s.. \tag{2.8}$$

## 3. Simulation study

In this section, we carry out some simulations to show the finite sample performance of the proposed method. The surrogate variable $\widetilde{x}$ was generated such that $\widetilde{x} = 1.15x + \delta\epsilon$, where the design points $x_i = \Phi^{-1}(i/n + 1)$ and $\epsilon$ is $N(0,1)$ distributed, and $\delta$ is the standard deviation of the measurement error. Results for $\delta = 0.2$ and $\delta = 0.4$ are reported. Simulations were run with validation and primary data size $(n, N) = (30, 90), (60, 180), (100, 300)$ and $(n, N) = (30, 150), (60, 300), (100, 500)$. The kernel function $K_1(\cdot)$ was taken to be the product kernel $K_1(x_1, x_2) = K_0(x_1)K_0(x_2)$, where $K_0(x)$ and the kernel function $K_2(\cdot)$ were taken to be the Gaussian kernels and the bandwidth $b_n = n^{-1/5}$, $h_N = N^{-1/5}$.

On the other hand, the design points $t_i = i/N$, the function $g(\cdot)$ is chosen to be $g(t) = \sin 2\pi t$. Since $\{e_i, \mathcal{F}_i, i \in Z\}$ is a sequence of martingale differences, we first take the martingale sequence $\{\xi_i, \sigma(\xi_i), i \geq 1\}$, and let $e_i = \xi_{i+1} - \xi_i$. We generate the first random number $\xi_1 \sim N(0,1)$, then take the following $\xi_2, \ldots, \xi_{n+1}$ according to the conditional distribution $\xi_{i+1}|\xi_i \sim N(\xi_i, 0.2^2), i = 1, \ldots, n$.

For each sample size $(n, N)$ and selected value of $\delta$, we calculated the estimator $\beta_{nN}$ of $\beta = 0.5$, which is given by (2.4). In each case the number of simulated realization is 500. We calculate the sample means and the average square error (ASE) of $\beta_{nN}$, The simulation results are summarized in Table 1. Moreover, we also display a set of the corresponding nonparametric component estimator $\widetilde{g}_{nN}(t_i)$ given by (2.5) in Figures 1 and 2.

| | $(n, N)$ | (30,90) | (60,180) | (100,300) | (30,150) | (60,300) | (100,500) |
|---|---|---|---|---|---|---|---|
| $\delta = 0.2$ | Mean($\beta_{nN}$) | 0.4426 | 0.4749 | 0.5194 | 0.4365 | 0.4647 | 0.5158 |
| | $\beta_{nN}$ASE | 0.0684 | 0.0231 | 0.0185 | 0.0748 | 0.0439 | 0.0318 |
| $\delta = 0.4$ | Mean($\beta_{nN}$) | 0.4377 | 0.4678 | 0.5226 | 0.4351 | 0.4632 | 0.5213 |
| | $\beta_{nN}$ASE | 0.0691 | 0.0380 | 0.0335 | 0.0788 | 0.0535 | 0.0353 |

Table 1  Sample means and ASE for $\beta_{nN}$ with $\beta = 0.5$

Moreover, we also display a set of the corresponding nonparametric component estimator $\widetilde{g}_{nN}(t_i)$ given by (2.5). Figures 1 and 2 show the curves of $g(x)$ and its semiparametric estimator $\widetilde{g}_{nN}(x)$ under the two types of models, respectively. The solid line is the regression function $g(x)$, the broken line is the estimator $\widetilde{g}_{nN}(x)$.
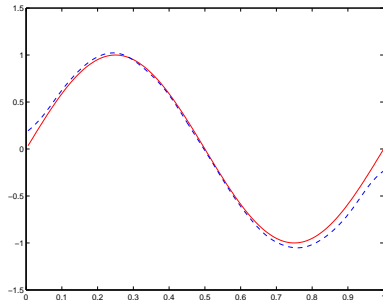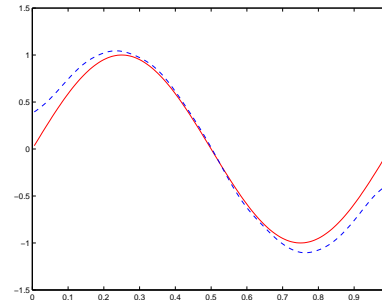


Figure 1 $(n, N) = (60, 180), \delta = 0.2$      Figure 2 $(n, N) = (60, 180), \delta = 0.4$

## 4. The proofs

In order to prove our results, we need the following lemmas.

**Lemma 4.1** *Let $\{\xi_i, \mathcal{F}_i, i \geq 1\}$ be a martingale difference sequence $|\xi_i| \leq b$ a.s. $i = 1, 2, \ldots$. Then, for $\forall \varepsilon > 0$, we have*

$$P(|\sum_{i=1}^{n} \xi_i| > \varepsilon) \leq 2 \exp\{-\frac{\varepsilon^2}{2(2\overline{\sigma}_n^2 + b\varepsilon)}\},$$

*where $\overline{\sigma}_n^2 = \sum_{i=1}^{n} \| E(\xi_i^2 | \mathcal{F}_{i-1}) \|$.*

**Proof of Lemma 4.1** The proof of Lemma 4.1 is similar to the proof of Inference 3.1 by Li and Liu in [16]. □

**Lemma 4.2** *Let $\{S_n = \sum_{i=1}^{n} \xi_i, n \geq 1\}$ be a random sequence, and $\{\mathcal{F}_n, n \geq 1\}$ be a non-decreasing $\sigma$-fields such that $S_n \in \mathcal{F}_n, n \geq 1$. If*

    *(i) $\sum_{i=1}^{n} P(|\xi_i| \geq C | \mathcal{F}_{i-1}) < \infty$;*

    *(ii) $\sum_{i=1}^{n} E(\xi_i I(|\xi_i| \leq C) | \mathcal{F}_{i-1}) < \infty$;*

    *(iii) $\sum_{i=1}^{n} \{E(\xi_i^2 I(|\xi_i| \leq C) | \mathcal{F}_{i-1}) - (E(\xi_i I(|\xi_i| \leq C) | \mathcal{F}_{i-1}))^2\} < \infty$,*

*then $\{S_n, n \geq 1\}$ converges a.s..*

**Proof of Theorem 2.2** It is easy to see

$$\beta_{n,N} - \beta$$

$$= S_{n,N}^{-2} \sum_{i=n+1}^{n+N} \overline{u}_n(\nu_i) \Big[ (u_n(\nu_i)\beta + g(t_i) + \varepsilon_i) - \sum_{j=n+1}^{n+N} \omega_{Nj}(t_i)(u_n(\nu_j)\beta + g(t_j) + \varepsilon_j) \Big] - \beta$$

$$= S_{n,N}^{-2} \sum_{i=n+1}^{n+N} \overline{u}_n(\nu_i)\varepsilon_i - S_{n,N}^{-2} \sum_{i=n+1}^{n+N} \Big[ \overline{u}_n(\nu_i) \Big( \sum_{j=n+1}^{n+N} \omega_{Nj}(t_i)\varepsilon_j \Big) \Big] +$$

$$S_{n,N}^{-2} \sum_{i=n+1}^{n+N} \overline{u}_n(\nu_i) \Big( g(t_i) - \sum_{j=n+1}^{n+N} \omega_{Nj}(t_i)g(t_j) \Big) +$$

$$S_{n,N}^{-2} \sum_{i=n+1}^{n+N} \overline{u}_n(\nu_i) \Big( u_n(\nu_i)\beta - \sum_{j=n+1}^{n+N} \omega_{Nj}(t_i)u_n(\nu_j)\beta \Big) - \beta$$

$$= S_{n,N}^{-2} \sum_{j=n+1}^{n+N} \Big( \overline{u}_n(\nu_j) - \sum_{i=n+1}^{n+N} \omega_{Nj}(t_i)\overline{u}_n(\nu_i) \Big) \varepsilon_j + S_{n,N}^{-2} \sum_{i=n+1}^{n+N} \overline{u}_n(\nu_i) \Big( g(t_i) - \sum_{j=n+1}^{n+N} \omega_{Nj}(t_i)g(t_j) \Big)$$

$$= S_{n,N}^{-2} \sum_{j=n+1}^{n+N} \Big( \overline{u}_n(\nu_j) - \sum_{i=n+1}^{n+N} \omega_{Nj}(t_i)\overline{u}_n(\nu_i) \Big) (\varepsilon_j - E(\varepsilon_j|\mathcal{F}_{j-1})) +$$

$$S_{n,N}^{-2} \sum_{j=n+1}^{n+N} \Big( \overline{u}_n(\nu_j) - \sum_{i=n+1}^{n+N} \omega_{Nj}(t_i)\overline{u}_n(\nu_i) \Big) (x_j - u_n(\nu_j))\beta +$$

$$S_{n,N}^{-2} \sum_{i=n+1}^{n+N} \overline{u}_n(\nu_i) \Big( g(t_i) - \sum_{j=n+1}^{n+N} \omega_{Nj}(t_i)g(t_j) \Big)$$

$$=: I_1 + I_2 + I_3.$$

By the uniform continuity of the function $g(\cdot)$ on the interval $[0, 1]$, (2.2), and the conditions (A1), (A2) and (A4), similarly to the proof of Theorem 2 in [2], we have, as $N \longrightarrow \infty$,

$$I_3 = S_{n,N}^{-2} \sum_{i=n+1}^{n+N} \overline{u}_n(\nu_i) \Big( g(t_i) - \sum_{j=n+1}^{n+N} \omega_{Nj}(t_i)g(t_j) \Big) = \mathrm{o}(1). \tag{4.1}$$

Clearly, $I_1$ can be decomposed as

$$I_1 = S_{n,N}^{-2} \sum_{j=n+1}^{n+N} \Big( \overline{u}_n(\nu_j) - \sum_{i=n+1}^{n+N} \omega_{Nj}(t_i)\overline{u}_n(\nu_i) \Big) (\varepsilon_j - E(\varepsilon_j|\mathcal{F}_{j-1}))$$

$$= S_{n,N}^{-2} \sum_{j=n+1}^{n+N} \Big( \overline{u}_n(\nu_j) - \sum_{i=n+1}^{n+N} \omega_{Nj}(t_i)\overline{u}_n(\nu_i) \Big) e_j$$

$$= \sum_{j=n+1}^{n+N} S_{n,N}^{-2}\overline{u}_n(\nu_j)e_j - \sum_{j=n+1}^{n+N} \Big[ S_{n,N}^{-2} \Big( \sum_{i=n+1}^{n+N} \omega_{Nj}(t_i)\overline{u}_n(\nu_i) \Big) \Big] e_j$$

$$=: I_{11} - I_{12}.$$

$$I_{11} = \sum_{j=n+1}^{n+N} S_{n,N}^{-2}\overline{u}_n(\nu_j)e_j =: \sum_{j=n+1}^{n+N} b_j e_j.$$

By assumptions (A1), (A3) and (A5), it follows that

$$\sum_{j=n+1}^{n+N} b_j^2 = \sum_{j=n+1}^{n+N} \frac{\overline{u}_n^2(\nu_j)}{S_{n,N}^4} = \frac{1}{S_{n,N}^2} = \mathrm{O}(N^{-1}),$$

$$\max_{n+1 \leq j \leq n+N} |b_j| = \max_{n+1 \leq j \leq n+N} \Big| \frac{\overline{u}_n(\nu_j)}{S_{n,N}^2} \Big| \leq \frac{C_7}{S_{n,N}^2} = \mathrm{O}(N^{-1}).$$

For $\mu_N = \frac{1}{\log\log N}$, let $\delta_N = \sqrt{\frac{\mu_N}{2C_8}}$, $e'_{Nj} = e_j I(|e_j| \leq \delta_N^2 j^{\frac{1}{2}})$, $e''_{Nj} = e_j - e'_{Nj}$, $e_{Nj} =$

$b_j(e'_{Nj} - E(e'_{Nj}|\mathcal{F}_{j-1}))$, so $\{e_{Nj}, \mathcal{F}_j, n+1 \leq j \leq n+N\}$ is a martingale difference sequence,

$$\max_{n+1 \leq j \leq n+N} |e_{Nj}| \leq 2 \max_{n+1 \leq j \leq n+N} |b_j| \delta_N^2 (n+N)^{\frac{1}{2}} = \mathrm{O}(N^{-\frac{1}{2}}(\log\log N)^{-1}),$$

$$\sum_{j=n+1}^{n+N} \| E(e_{Nj}^2|\mathcal{F}_{j-1}) \| \leq \sum_{j=n+1}^{n+N} b_j^2 \| E(e_{Nj}^{'2}|\mathcal{F}_{j-1}) \| \leq \frac{\sigma^2}{S_{nN}^2} = \mathrm{O}(N^{-1}).$$

By Lemma 4.1, we have

$$P\Big(\Big|\sum_{j=n+1}^{n+N} e_{Nj}\Big| \geq \mu_N\Big) \leq 2 \exp\Big\{ \frac{-\mu_N^2}{\frac{4\sigma^2}{S_{nN}^2} + \frac{2\mu_N^2(n+N)^{\frac{1}{2}}}{S_{nN}^2}} \Big\}$$

$$\leq 2 S_{nN}^{-4}(n+N)^{\frac{1}{2}} = \mathrm{O}(N^{-\frac{3}{2}}).$$

For $\forall \mu > 0$, $\exists N$, such that $\mu_N < \mu$, then

$$P\Big(\Big|\sum_{j=n+1}^{n+N} e_{Nj}\Big| \geq \mu\Big) \leq P\Big(\Big|\sum_{j=n+1}^{n+N} e_{Nj}\Big| \geq \mu_N\Big) \leq \mathrm{O}(N^{-\frac{3}{2}}),$$

by the Borel-Cantelli Lemma, it follows that

$$\sum_{j=n+1}^{n+N} e_{Nj} \longrightarrow 0 \quad \text{a.s.,} \quad n, N \longrightarrow \infty. \tag{4.2}$$

Since $e''_{Nj} = e_j - e'_{Nj} = e_j I(|e_j| > \delta_N^2 j^{\frac{1}{2}}) = e_j I(e_j > \delta_N^2 j^{\frac{1}{2}}) + e_j I(e_j < -\delta_N^2 j^{\frac{1}{2}})$, let $e_{Nj}^{''+} = e_j I(e_j > \delta_N^2 j^{\frac{1}{2}}), e_{Nj}^{''-} = e_j I(e_j < -\delta_N^2 j^{\frac{1}{2}})$.

For a fixed constant $d > 0$, we have

$$P(|e_{Nj}^{''+}| > d|\mathcal{F}_{j-1}) = P(e_{Nj} > \max\{d, \delta_N^2 j^{\frac{1}{2}}\}|\mathcal{F}_{j-1})$$

$$\leq P(e_{Nj} > \delta_N^2 j^{\frac{1}{2}}|\mathcal{F}_{j-1}) \leq P(|e_{Nj}|^2 > \delta_N^4 j|\mathcal{F}_{j-1}).$$

Then

$$\sum_{j=n+1}^{n+N} P(|e_{Nj}^{''+}| > d|\mathcal{F}_{j-1}) \leq \sum_{j=n+1}^{n+N} P(|e_{Nj}|^2 > \delta_N^4 j|\mathcal{F}_{j-1})$$

$$\leq \sum_{j=n+1}^{n+N} \frac{\sup_{j \geq 1} \| E|e_j|^\alpha|\mathcal{F}_{j-1} \|}{\delta_N^{2\alpha} j^{\frac{\alpha}{2}}} < \infty. \tag{4.3}$$

Let

$$(e_{Nj}^{''+})_d =: e_{Nj}^{''+} I(|e_{Nj}^{''+}| < d) = e_j I(e_j > \delta_N^2 j^{\frac{1}{2}}) I(\delta_N^2 j^{\frac{1}{2}} < e_j \leq d) = e_j I(\delta_N^2 j^{\frac{1}{2}} < e_j \leq d).$$

Then

$$\sum_{j=n+1}^{n+N} E\{(e_{Nj}^{''+})_d|\mathcal{F}_{j-1}\} \leq \sum_{j=n+1}^{n+N} dE(I(\delta_N^2 j^{\frac{1}{2}} < e_j \leq d)|\mathcal{F}_{j-1})$$

$$\leq \sum_{j=n+1}^{n+N} dE(I(|e_j| > \delta_N^2 j^{\frac{1}{2}})|\mathcal{F}_{j-1})$$

$$\leq d \sum_{j=n+1}^{n+N} P(|e_j|^2 > \delta_N^4 j|\mathcal{F}_{j-1}) < \infty. \tag{4.4}$$

Note that

$$\sum_{j=n+1}^{n+N} \{E\{(e_{Nj}^{''+})_d^2|\mathcal{F}_{j-1}\} - E^2\{(e_{Nj}^{''+})_d|\mathcal{F}_{j-1}\}\}$$

$$\leq \sum_{j=n+1}^{n+N} \{E\{(e_{Nj}^{''+})_d^2|\mathcal{F}_{j-1}\} = \sum_{j=n+1}^{n+N} E\{e_j^2 I(\delta_N^2 j^{\frac{1}{2}} < e_j \leq d)|\mathcal{F}_{j-1}\}$$

$$\leq d^2 \sum_{j=n+1}^{n+N} E(I(|e_j| > \delta_N^2 j^{\frac{1}{2}})|\mathcal{F}_{j-1}) \leq d^2 \sum_{j=n+1}^{n+N} P(|e_j|^2 > \delta_N^4 j|\mathcal{F}_{j-1}) < \infty. \qquad (4.5)$$

Using Lemma 4.2, together with (4.3), (4.4) and (4.5) gives

$$\sum_{j=n+1}^{n+N} |e_{Nj}^{''+}| < \infty \quad \text{a.s..} \qquad (4.6)$$

Similarly to the arguments used in the proof of (4.6), we can prove that $\sum_{j=n+1}^{n+N} |e_{Nj}^{''-}| < \infty$ a.s.. So we have

$$\sum_{j=n+1}^{n+N} |e_{Nj}^{''}| < \infty \quad \text{a.s..} \qquad (4.7)$$

From (4.7),

$$|\sum_{j=n+1}^{n+N} b_j e_{Nj}^{''}| \leq (\max_{n+1 \leq j \leq n+N} |b_j|)(\sum_{j=n+1}^{n+N} |e_{Nj}^{''}|) = \mathrm{O}(S_{nN}^{-2}), \quad n, N \longrightarrow \infty. \qquad (4.8)$$

$$|\sum_{j=n+1}^{n+N} b_j E(e_{Nj}^{''}|\mathcal{F}_{j-1})| \leq (\max_{n+1 \leq j \leq n+N} |b_j|)\Big(\sum_{j=n+1}^{n+N} E(|e_j|I(|e_j| \geq \delta_N^2 j^{\frac{1}{2}})|\mathcal{F}_{j-1})\Big)$$

$$\leq CS_{nN}^{-2} \sum_{j=n+1}^{n+N} \delta_N^{-2} j^{-\frac{1}{2}} E(e_j^2 I(|e_j| \geq \delta_N^2 j^{\frac{1}{2}})|\mathcal{F}_{j-1})$$

$$\leq CS_{nN}^{-2} \sum_{j=n+1}^{n+N} \delta_N^{-2} j^{-\frac{1}{2}} \leq CS_{nN}^{-2} N \log\log N/(n+1)^{\frac{1}{2}}$$

$$= \mathrm{O}(N^{-\frac{1}{2}} \log\log N). \qquad (4.9)$$

Combining (4.2), (4.7), (4.8) and (4.9), we have

$$|I_{11}| = |\sum_{j=n+1}^{n+N} b_j(e_{Nj}' + e_{Nj}^{''})|$$

$$= |\sum_{j=n+1}^{n+N} e_{Nj} + \sum_{j=n+1}^{n+N} b_j e_{Nj}^{''} + \sum_{j=n+1}^{n+N} b_j E(e_{Nj}' + e_{Nj}^{''}|\mathcal{F}_{j-1}) - \sum_{j=n+1}^{n+N} b_j E(e_{Nj}^{''}|\mathcal{F}_{j-1})|$$

$$\leq |\sum_{j=n+1}^{n+N} e_{Nj}| + |\sum_{j=n+1}^{n+N} b_j e_{Nj}^{''}| + |\sum_{j=n+1}^{n+N} b_j E(e_{Nj}^{''}|\mathcal{F}_{j-1})| \longrightarrow 0 \quad \text{a.s.,} \quad n, N \longrightarrow \infty.$$

It follows that

$$I_{11} \longrightarrow 0 \quad \text{a.s.} \quad n, N \longrightarrow \infty. \qquad (4.10)$$

Since

$$I_{12} = \sum_{j=n+1}^{n+N} \left[ S_{n,N}^{-2} \left( \sum_{i=n+1}^{n+N} \omega_{Nj}(t_i) \overline{u}_n(\nu_i) \right) \right] e_j =: \sum_{j=n+1}^{n+N} \overline{b}_j e_j,$$

from Remark 2.1, we have

$$\max_{n+1 \leq j \leq n+N} |\overline{b}_j| \leq S_{n,N}^{-2} \max_{n+1 \leq i \leq n+N} \sum_{j=n+1}^{n+N} |\omega_{Nj}(t_i)| \max_{n+1 \leq i \leq n+N} |\overline{u}_n(\nu_i)| = \mathrm{O}(N^{-1}). \qquad (4.11)$$

$$\sum_{j=n+1}^{n+N} \overline{b}_j^2 \leq S_{n,N}^{-4} \sum_{j=n+1}^{n+N} \max_{n+1 \leq i,j \leq n+N} |\omega_{Nj}(t_i)|^2 \sum_{i=n+1}^{n+N} [\overline{u}_n(\nu_i)]^2 = \mathrm{O}(N^{-1}). \qquad (4.12)$$

Similarly to the proof of (4.10), we can prove that

$$I_{12} \longrightarrow 0 \quad \text{a.s.} \quad n, N \longrightarrow \infty. \qquad (4.13)$$

With the definition of (2.1) and the conditions of (A1), (A2), (A3), (A5), similarly to the proof of Theorem 2 in [2], we have, as $n, N \longrightarrow \infty$,

$$I_2 = S_{n,N}^{-2} \sum_{j=n+1}^{n+N} \left( \overline{u}_n(\nu_j) - \sum_{i=n+1}^{n+N} \omega_{Nj}(t_i) \overline{u}_n(\nu_i) \right) (x_j - u_n(\nu_j)) \beta = \mathrm{o}(1). \qquad (4.14)$$

(4.1), (4.10), (4.13) and (4.14) together imply (2.7) of Theorem 2.2.

Observe that

$$\widetilde{g}_{nN}(t) - g(t) = (\beta - \beta_{nN}) g_{2N}(t) + \left( \sum_{j=n+1}^{n+N} \omega_{Nj}(t) g(t_j) - g(t) \right) + \sum_{j=n+1}^{n+N} \omega_{Nj}(t) \varepsilon_j$$

$$=: A_{nN}(t) + B_{nN}(t) + C_{nN}(t).$$

By (2.7) and $\sup_{0 \leq t \leq 1} |g_{2N}(t)| \leq C$, it follows that

$$\sup_{0 \leq t \leq 1} A_{nN}(t) \longrightarrow 0 \quad \text{a.s.,} \quad n, N \longrightarrow \infty. \qquad (4.15)$$

Using the reasoning similar to the proof of (4.1), we have

$$\sup_{0 \leq t \leq 1} B_{nN}(t) \longrightarrow 0, \quad \text{as} \quad N \longrightarrow \infty. \qquad (4.16)$$

Since $E(\varepsilon_j | \mathcal{F}_{j-1}) = (x_j - u_n(v_j)) \beta$, noting that

$$\sum_{j=n+1}^{n+N} \omega_{Nj}(t) \varepsilon_j = \sum_{j=n+1}^{n+N} \omega_{Nj}(t) (\varepsilon_j - E(\varepsilon_j | \mathcal{F}_{j-1})) + \sum_{j=n+1}^{n+N} \omega_{Nj}(t) (x_j - u_n(v_j)) \beta$$

$$= \sum_{j=n+1}^{n+N} \omega_{Nj}(t) e_j + \sum_{j=n+1}^{n+N} \omega_{Nj}(t) (x_j - u_n(v_j)) \beta,$$

by the same arguments as in the proofs of (4.10) and (4.14), we can prove

$$\sum_{j=n+1}^{n+N} \omega_{Nj}(t) e_j \longrightarrow 0 \quad \text{a.s.,} \quad n, N \longrightarrow \infty. \qquad (4.17)$$

$$\sum_{j=n+1}^{n+N} \omega_{Nj}(t) (x_j - u_n(v_j)) \beta \longrightarrow 0, \quad \text{as} \quad n, N \longrightarrow \infty. \qquad (4.18)$$

(4.17) and (4.18) show that

$$\sup_{0 \le t \le 1} C_{nN}(t) \longrightarrow 0 \quad \text{a.s.,} \quad n, N \longrightarrow \infty. \tag{4.19}$$

Hence, Theorem 2.2 is proved. $\square$

**Acknowledgements** We thank the referees for their time and comments.

# References

[1] G. ANEIROS, A. QUINTELA. *Modified cross-validation in semiparametric regression models with dependent errors.* Comm. Statist. Theory Methods, 2001, **30**(2): 289–307.

[2] Guanghui CAI. *Estimation of partial linear error-in-variables models for $\rho^-$-mixing dependence data.* J. Math. Chem., 2008, **43**(1): 375–385.

[3] R. J. CARROLL, M. P. WAND. *Semiparametric estimation in logistic measurment error models.* J. R. Statist. Soc., 1991, **53**: 653–653.

[4] H. CHEN. *Convergence rates for parametric components in a partly linear model.* Ann. Statist., 1988, **16**(1): 136–146.

[5] Xia CHEN, Hengjian CUI. *Empirical likelihood inference for partial linear models under martingale difference sequence.* Statist. Probab. Lett., 2008, **78**(17): 2895–2901.

[6] R. ENGLE, C. GRANGER, J. RICE, et al. *Nonparametric estimators of the relation between weather and electricity sales.* J. Amer. Statist. Assoc., 1986, **81**: 310–320.

[7] Guoliang FAN, Hanying LIANG. *Empirical likelihood inference for semiparametric model with linear process errors.* J. Korean Statist. Soc., 2010, **39**(1): 55–65.

[8] Guoliang FAN, Hanying LIANG, Hongxia XU. *Empirical likelihood for a heteroscedastic partial linear model.* Comm. Statist. Theory Methods, 2011, **40**(8): 1396–1417.

[9] Jiti GAO, V. V. ANH. *Semiparametric regression under long-range dependent errors.* J. Statist. Plann. Inference, 1999, **80**(1-2): 37–57.

[10] Jiti GAO, Xiru CHEN, Lincheng ZHAO. *Asymptotic normality of a class of estimators in partial linear models.* Acta Math. Sinica, 1994, **37**(2): 256–268.

[11] W. HÄDLE, H. LIANG, J. GAO. *Partial Linear Models.* Heidelberg: Physica-Verlag, 2000.

[12] N. HECKMAN. *Spline smoothing in partly linear models.* J. Roy. Statist. Soc., 1986, **48**: 244–248.

[13] Hanying LIANG, Guoliang FAN. *Berry-Esseen type bounds of estimators in a semiparametric model with linear process errors.* J. Multivariate Anal., 2009, **100**(1): 1–15.

[14] Hanying LIANG, V. MAMMITZSCH, J. STEINEBACH. *On a semiparametric regression model whose errors form a linear process with negatively associated innovations.* Statistics, 2006, **40**(3): 207–226.

[15] Hanying LIANG, B. Y. JING. *Asmptotic normality in partial linear models based on dependent errors.* J. Statist. Plann. Infer., 2009, **139**: 1357–1371.

[16] Guoliang LI, Luqin LIU. *Strong consistency of a class of estimators in partial linear model under martingle difference sequence.* Acta Meth. Sinica., 2007, **27**(5): 788–801.

[17] M. S. PEPE. *Inference using surrogate outcome data and a validation sample.* Biometrika, 1992, **79**(2): 355–365.

[18] M. S. PEPE, T. R. FLEMING. *A general nonparametric method for dealing with errors in missing or surrogate covariate data.* J. Amer. Statist. Assoc., 1992, **86**: 108–113.

[19] J. H. SEPANSKI, L. F. LEE. *Semiparametric estimation of nonlinear errors-in-variables models with validation study.* J. Nonparametr. Statist., 1995, **4**(4): 365–394.

[20] P. SPECKMAN. *Kernel smoothing in partial linear models.* J. Roy. Statist. Soc. Ser. B, 1988, **50**(3): 413–436.

[21] L. A. STEFANSKI, R. J. CARROLL. *Covariate measurement error in generalized linear models.* Biometrika, 1985, **72**: 583–592.

[22] Xiaoqian SUN, Jinhong YOU, Gemai CHEN. *Convergence rates of estimators in partial linear regression models with $MA(\infty)$ error process.* Comm. Statist. Theory Methods, 2002, **31**(12): 2251–2273.

[23] Qihua WANG. *Estimation of partial linear error-in-variables models with validation data.* J. Multivariate Anal., 1999, **69**(1): 30–64.

[24] Qihua WANG, J. N. K. RAO. *Empirical likelihood-based inferenc in linear errors-in-covariables models with validation data.* Biometrika, 2002, **89**(2): 345–357.

[25] Jinhong YOU, Gemai CHEN, Yong ZHOU. *Statistical inference of partially linear regression models with heteroscedastic errors.* J. Multivariate Anal., 2007, **98**(8): 1539–1557.