



Partial covariate adjusted regression

Damla Şentürk^a, Danh V. Nguyen^{b,*}

^aDepartment of Statistics, Pennsylvania State University, University Park, PA 16802, USA

^bDivision of Biostatistics, University of California, Davis, CA 95616, USA

ARTICLE INFO

Article history:

Received 23 July 2006

Received in revised form

9 March 2008

Accepted 1 April 2008

Available online 22 May 2008

Keywords:

Asymptotic normality

Binning

Confidence intervals

Multiple regression

Varying-coefficient models

ABSTRACT

Covariate adjusted regression (CAR) is a recently proposed adjustment method for regression analysis where both the response and predictors are not directly observed [Şentürk, D., Müller, H.G., 2005. Covariate adjusted regression. *Biometrika* 92, 75–89]. The available data have been distorted by unknown functions of an observable confounding covariate. CAR provides consistent estimators for the coefficients of the regression between the variables of interest, adjusted for the confounder. We develop a broader class of partial covariate adjusted regression (PCAR) models to accommodate both distorted and undistorted (adjusted/unadjusted) predictors. The PCAR model allows for unadjusted predictors, such as age, gender and demographic variables, which are common in the analysis of biomedical and epidemiological data. The available estimation and inference procedures for CAR are shown to be invalid for the proposed PCAR model. We propose new estimators and develop new inference tools for the more general PCAR setting. In particular, we establish the asymptotic normality of the proposed estimators and propose consistent estimators of their asymptotic variances. Finite sample properties of the proposed estimators are investigated using simulation studies and the method is also illustrated with a Pima Indians diabetes data set.

© 2008 Elsevier B.V. All rights reserved.

1. Introduction

Covariate adjusted regression (CAR) has been recently proposed to adjust for the distorting effects of a confounder in a regression setting. It was motivated by a common adjustment method in medical and health related studies. The adjustment entails normalization by anthropometric measurements, such as body mass index (BMI) and/or other measures of body configuration, as confounding variables that affect the primary variables of interest. For example, in a study involving haemodialysis patients, it is of interest to examine the relationship between elevated plasma fibrinogen level (a risk factor for cardiovascular disease in haemodialysis patients) and other predictors, such as serum transferrin protein level (Kaysen et al., 2002; Şentürk and Müller, 2005). However, both primary variables, fibrinogen and transferrin protein levels, are known to depend on BMI, which exerts a confounding effect on the protein measurements. A common approach to adjust for the confounders, like BMI, is to normalize the primary variables of interest by simply dividing (by BMI). Note that this adjustment by division implies that the assumed contamination is of a multiplicative form. Let \tilde{Y} , \tilde{X} , and U denote the observed fibrinogen concentration, serum transferrin level, and confounder BMI, respectively. Using these notations, the adjusted primary variables that are thought to be free from the confounding effect of BMI are

$$Y = \frac{\tilde{Y}}{U} \quad \text{and} \quad X = \frac{\tilde{X}}{U}.$$

* Corresponding author.

E-mail addresses: dsenturk@stat.psu.edu (D. Şentürk), ucdnguyen@ucdavis.edu (D.V. Nguyen).

The basic motivation for the above adjustment is to obtain normalized versions of the observed primary variables by removing the confounder effects, so that the measurements are comparable across patients. Other examples include normalizations by BMI in studies on diabetes, and division of brain volumetric structures by total brain volume in neurological studies (Pinter et al., 2001).

Şentürk and Müller (2005, 2006) proposed a more flexible adjustment, by modeling the confounding through *unknown functions* of the confounder instead of the confounder itself. This reflects the uncertainty encountered in many applications about the precise nature of the commonly assumed multiplicative relation between the confounder and the variables. For the case of p predictors, Şentürk and Müller model the underlying variables as

$$Y = \frac{\tilde{Y}}{\psi(U)}, \quad X_1 = \frac{\tilde{X}_1}{\phi_1(U)}, \dots, X_p = \frac{\tilde{X}_p}{\phi_p(U)},$$

where they are defined to be the parts of the observed variables, $\tilde{Y}, \tilde{X}_1, \dots, \tilde{X}_p$, that are independent of the observable confounder U . In the haemodialysis data example, the latent variables would be defined to be serum protein levels adjusted for BMI. Here, $\phi_1(\cdot), \dots, \phi_p(\cdot)$ and $\psi(\cdot)$ denote unknown smooth contaminating functions of U . CAR gives consistent estimators of the coefficients in the unobserved regression model which can be expressed as

$$Y = \gamma_0 + \sum_{r=1}^p \gamma_r X_r + e,$$

where e is the error term, assumed to be independent of $\{X_r\}_{r=1}^p$ and U . The estimation procedure is based on the observed data: the distorted response, \tilde{Y} , distorted predictors, $\{\tilde{X}_r\}_{r=1}^p$, and confounder U .

The main goal of this paper is to construct estimation and inference procedures needed to allow some of the variables as unadjusted/undistorted predictors, denoted by Z_1, \dots, Z_s . The proposed underlying regression model is then of the form

$$Y = \gamma_0 + \sum_{r=1}^p \gamma_r X_r + \sum_{s=1}^q \delta_s Z_s + e, \quad (1)$$

where $\tilde{X}_r = \phi(U)X_r$ and $\tilde{Y} = \psi(U)Y$ denote the adjusted/distorted predictors and response, respectively, and Z_s denotes the unadjusted predictors. The observed data are $\tilde{Y}, \tilde{X}_r, Z_s$ and U . Furthermore, the confounding covariate, U , is allowed to depend on the unadjusted predictors (Z_s). The flexibility of allowing both distorted and undistorted predictors is needed, particularly in the regression analysis of biomedical data. In many of these applications, which motivated our current work, the researcher is directly interested in the effects of $Z_s =$ age, gender, obesity measures and/or ethnicity on Y . Therefore, the predictors Z_s are unadjusted/undistorted. A specific example is the study of Kaysen et al. (2002) where albumin turnover and protein catabolic rate were adjusted for body surface area (U) via division, while age and gender were among the unadjusted variables. This is an example of model (1) above, where an adjustment method is needed that allows for unadjusted predictors. In this example, model (1) is used to reflect (a) the known dependence of albumin turnover and protein catabolic rate on U and (b) the specific interest in the direct effects of age and gender in the regression relationship. This issue is further discussed in Section 3 in the context of estimation as well as in Section 8.

Under the more general partial covariate adjusted regression (PCAR) setting, formally presented in Section 2, the original CAR estimators (Şentürk and Müller, 2005) for $\{\delta_s\}_{s=1}^q$ are inconsistent and may have an arbitrarily large asymptotic bias as shown in Section 3. We propose alternative estimators that are consistent under this extended CAR setting where the issues of estimation are discussed in Section 3. The proposed PCAR methodology, like CAR, provides consistent estimators not only under multiplicative but also under additive distortion as discussed also in Section 3. The inference procedures developed for the CAR modeling are not valid for the PCAR setting, mainly due to the different dependence structure needed for PCAR. This new structure is explained in detail in Sections 2 and 3. Thus, we develop new theoretical tools for valid inference in the PCAR model. We derive the asymptotic distributions of the proposed estimators, and present them in Section 4. Consistent estimators of the asymptotic variance are also derived in Section 4. Simulation studies to characterize the finite sample properties of the proposed estimators are summarized in Section 6. The method is further illustrated with a Pima Indians diabetes data set given in Section 5. The proofs of the main results are assembled in Section 7, where some technical conditions and auxiliary results are deferred to the Appendix. We conclude with a brief discussion in Section 8.

We note here that an advantage of CAR and PCAR is that, under the identifiability conditions introduced in Section 2, it yields consistent estimates whether the distortion is multiplicative or additive, i.e. $Y = \tilde{Y} - \psi(U)$ and $X_r = \tilde{X}_r - \phi_r(U)$. (A more detailed discussion of the additive distortion case is given in Section 3.) Additive distortions can be handled by the method of nonparametric partial regression; however, there existed no consistent estimation procedure targeting the γ 's under multiplicative distortion of both the response and the predictors. Also, the proposed distortion setting has similarities with measurement error modeling if one views $\psi(U)$ and $\phi_r(U)$ as unobserved errors affecting the response and predictors. There is an extensive literature on additive measurement error modeling, which dates back to Berkson (1950). See Carroll et al. (2006) and references therein for a comprehensive overview. A key difference between the proposed framework and traditional measurement error is that, in the proposed distortion setting, the error is a function of an observable covariate U . Also, since U is observable this information is incorporated into the estimation. While much work is devoted to additive measurement error, work on

multiplicative measurement errors is limited. Relevant work includes Hwang (1986) and Iturria et al. (1999) who proposed estimation procedures targeting the regression coefficients under multiplicative measurement error in the predictors. However, the case of multiplicative measurement errors that affect both the predictors and the response has not been considered previously to our knowledge. Therefore, the method and theory presented here may potentially be of interest to the area of measurement error modeling, despite the differences in the methodological motivation.

2. PCAR models

We consider the underlying (unobserved) regression model

$$Y_{ni} = \gamma_0 + \sum_{r=1}^p \gamma_r X_{nri} + \sum_{s=1}^q \delta_s Z_{nsi} + e_{ni} = \boldsymbol{\chi}_{ni}^T \boldsymbol{\alpha} + e_{ni}, \quad (2)$$

where Y_{ni} , e_{ni} , $\boldsymbol{\chi}_{ni} = (1, X_{n1i}, \dots, X_{npi}, Z_{n1i}, \dots, Z_{nqi})^T$ and $\boldsymbol{\alpha} = (\gamma_0, \dots, \gamma_p, \delta_1, \dots, \delta_q)^T$ are the response, error, $p + q$ predictors and unknown regression coefficients, respectively. The error variable e has mean zero, variance σ^2 and a finite moment that is higher than the 4th moment. The goal is estimation and inference for the parameter vector $\boldsymbol{\alpha}$ of the unobserved regression model (2). Estimation is based on available distorted predictor and response data, namely $\{\tilde{Y}_{ni}, \tilde{\boldsymbol{\chi}}_{ni}, U_{ni}\}_{i=1}^n$, where

$$\begin{aligned} \tilde{Y}_{ni} &= \psi(U_{ni})Y_{ni} \quad \text{and} \quad \tilde{\boldsymbol{\chi}}_{ni} = (1, \tilde{X}_{n1i}, \dots, \tilde{X}_{npi}, Z_{n1i}, \dots, Z_{nqi})^T \\ &= \{1, \phi_1(U_{ni})X_{n1i}, \dots, \phi_p(U_{ni})X_{npi}, Z_{n1i}, \dots, Z_{nqi}\}^T. \end{aligned} \quad (3)$$

The unknown distorting functions $\{\psi(\cdot), \phi_r(\cdot)\}_{r=1}^p$ are assumed to be smooth functions of the confounder, U .

Some constraints on the unknown smooth distortion functions are needed for the identifiability of the estimation problem. A set of reasonable constraints for $\psi(\cdot)$ and $\{\phi_r(\cdot)\}$ is implied by the natural assumption that the mean distorting effect should correspond to no distortion (Şentürk and Müller, 2005), i.e. the means of adjusted variables are the same as the means of the observed variables, $E(\tilde{X}_r) = E(X_r)$ and $E(\tilde{Y}) = E(Y)$. These conditions directly imply that

$$E\{\psi(U)\} = 1 \quad \text{and} \quad E\{\phi_r(U)\} = 1. \quad (4)$$

We consider the following dependence structure. The underlying predictors X_r and the undistorted predictors Z_s are allowed to be dependent. The error, e , is assumed to be mutually independent of X_r , Z_s , and U . We depart from the original CAR model (Şentürk and Müller, 2005), where U is independent of all the latent predictors, by allowing U to depend on Z_s , while still being independent of X_r . We will elaborate further on this important difference of the proposed setting from CAR at the end of Section 3. This is an important flexibility of the proposed method, since the common confounder correlates with all the observed variables in these distortion settings. This is consistent with the assumption that the observed predictors \tilde{X}_r and Z_s are dependent on the confounder U , and that the latent variable X_r is defined to be the part of \tilde{X}_r that is independent of U .

The assumption that the underlying predictors, $\{X_r\}_{r=1}^p$, and response, Y , are independent of the contaminating variable U is a fundamental assumption for the estimation procedure. It defines the proposed contamination setting through defining the unobserved, underlying variables. This independence assumption cannot be checked in practice since X_r and Y are unobservable. Instead, the question of more relevance in practice is whether the independence conditions help define interpretable latent variables of interest from their observable counterparts. In the haemodialysis data example, the latent variables are defined to be serum protein levels adjusted for BMI, which are commonly used in medical studies.

We refer to the model described by (2)–(4) as the PCAR model, since only a partial set of the predictors are adjusted for the confounder. Note also that the CAR model is a special case of the PCAR model.

For the estimation, note that it follows from (2) and the mutual independence of $\{e$ and $U\}$, $\{e$ and $(X_r, Z_s)\}$, and $\{U$ and $X_r\}$, for $r = 1, \dots, p$, $s = 1, \dots, q$, that the regression of \tilde{Y} on $\tilde{\boldsymbol{\chi}} = (1, \tilde{X}_1, \dots, \tilde{X}_p, Z_1, \dots, Z_q)^T$ leads to a fully observable varying coefficient model (Cleveland et al., 1991; Hastie and Tibshirani, 1993),

$$\tilde{Y}_{ni} = \beta_0(U_{ni}) + \sum_{r=1}^p \beta_r(U_{ni})\tilde{X}_{nri} + \sum_{s=1}^q \eta_s(U_{ni})Z_{nsi} + \varepsilon_{ni}. \quad (5)$$

where

$$\beta_0(U_{ni}) = \gamma_0\psi(U_{ni}), \quad \beta_r(U_{ni}) = \gamma_r \frac{\psi(U_{ni})}{\phi_r(U_{ni})}, \quad \eta_s(U_{ni}) = \delta_s\psi(U_{ni}), \quad (6)$$

and $\varepsilon(U_{ni}) = \psi(U_{ni})e_{ni}$. For an extensive overview of the estimation procedures proposed for varying coefficient models, see Wu and Yu (2002). Note that in (6), the varying coefficient functions, $\{\beta_r(\cdot)\}_{r=1}^p$, are proportional to the quotient of the original distorting functions, $\{\psi(\cdot)/\phi_r(\cdot)\}$; both the intercept function, $\beta_0(\cdot)$, and the functions $\{\eta_s(\cdot)\}_{s=1}^q$ are proportional to $\psi(\cdot)$. The constants of proportionality are precisely the underlying regression parameters, $\{\gamma_r, \delta_s\}$, of interest. These connections allow estimation of the underlying model through the varying coefficient functions. In Section 3 below, we describe an estimation procedure which targets $\{\gamma_r, \delta_s\}$ and mitigates the effects of the distorting functions $\{\psi(\cdot), \phi_r(\cdot)\}$.

3. Estimation procedure

The estimation of the regression coefficients, $\gamma_0, \{\gamma_r\}_{r=1}^p$ and $\{\delta_s\}_{s=1}^q$, in the underlying regression model $E(Y) = \gamma_0 + \sum_{r=1}^p \gamma_r X_r + \sum_{s=1}^q \delta_s Z_s$ is a two-step procedure. The first step involves estimation of the varying coefficient functions in model (5), namely $\beta_0(\cdot), \{\beta_r(\cdot)\}_{r=1}^p$ and $\{\eta_s(\cdot)\}_{s=1}^q$ using a binning approach. These varying functions are estimable because $\tilde{Y}, \tilde{X}_r, Z_s$, and U are all observable. The underlying regression coefficients are targeted in the second step, with weighted averages of the estimated $\beta_0(\cdot), \beta_r(\cdot)$ and $\eta_s(\cdot)$ for γ_0, γ_r and δ_s , respectively. The estimation makes use of the relations between the varying coefficient functions and the regression coefficients given by (6) and the identifiability conditions (4), as will be described next.

The binning approach for the estimation of the varying coefficient functions involves dividing the support of U into m equidistant bins and then fitting linear regressions of \tilde{Y} on $\tilde{\mathbf{X}}$ using the data falling within each bin. The observed data are the collection of n samples: $\{\tilde{Y}_{ni}, \tilde{\mathbf{X}}_{ni}, U_{ni}\}_{i=1}^n$. It is assumed that the confounding covariate, U , is bounded below and above, $a \leq U \leq b$, where $a < b$ are real numbers. In practice a and b would be taken to be $\min_i U_{ni}$ and $\max_i U_{ni}$, respectively. The estimation procedure initially divides the interval $[a, b]$ into m equidistant intervals, denoted B_{n1}, \dots, B_{nm} and referred to as bins. Let L_{nj} be the number of U_{ni} 's falling into bin j . Furthermore, let $(U'_{nj,k}, \tilde{X}'_{nrjk}, Z'_{nsjk}, \tilde{Y}'_{nj,k}, X'_{nrjk}, Y'_{nj,k})$ be the k th data element in the j th bin, B_{nj} , where $k = 1, \dots, L_{nj}$. Data elements in any given bin are marked by a prime.

After the initial binning of the data, a linear regression is fitted to the data observed within each bin $B_{nj}, j = 1, \dots, m$. The least squares estimator of the multiple regression of the data in the j th bin is

$$(\hat{\beta}_{n0j}, \dots, \hat{\beta}_{nrj}, \hat{\eta}_{n1j}, \dots, \hat{\eta}_{nqj})^T = (\tilde{\mathbf{X}}_{nj}^T \tilde{\mathbf{X}}_{nj})^{-1} \tilde{\mathbf{X}}_{nj}^T \tilde{\mathbf{Y}}_{nj} \tag{7}$$

where the response vector is $\tilde{\mathbf{Y}}_{nj} = (\tilde{Y}'_{nj1}, \dots, \tilde{Y}'_{njL_{nj}})^T$ and $\tilde{\mathbf{X}}_{nj} = (\tilde{X}'_{nrj1}, \dots, \tilde{X}'_{nrjL_{nj}})^T$ is the $L_{nj} \times (p + q + 1)$ data matrix in bin j , with the k th observation $\tilde{X}'_{nrjk} = (1, \tilde{X}'_{n1jk}, \dots, \tilde{X}'_{nrjk}, Z'_{n1jk}, \dots, Z'_{nqjk})^T$. The estimated regression coefficients in each bin (7) are the estimators of the varying coefficient functions.

In the second step of the estimation procedure, the estimators of the targeted regression parameters, $\gamma_0, \{\gamma_r\}_{r=1}^p$ and $\{\delta_s\}_{s=1}^q$, are obtained as weighted averages of the estimators $\{\hat{\beta}_{n0j}, \dots, \hat{\beta}_{nrj}, \hat{\eta}_{n1j}, \dots, \hat{\eta}_{nqj}\}_{j=1}^m$ from the m bins. The proposed PCAR estimators for $\gamma_0, \{\gamma_r\}_{r=1}^p$ and $\{\delta_s\}_{s=1}^q$ are

$$\hat{\gamma}_{n0} = \sum_{j=1}^m \frac{L_{nj}}{n} \hat{\beta}_{n0j}, \quad \hat{\gamma}_{nr} = \frac{1}{\bar{X}_{nr}} \sum_{j=1}^m \frac{L_{nj}}{n} \hat{\beta}_{nrj} \bar{X}'_{nrj} \quad \text{and} \quad \hat{\delta}_{ns} = \sum_{j=1}^m \frac{L_{nj}}{n} \hat{\eta}_{nsj} \tag{8}$$

where $\bar{X}_{nr} = n^{-1} \sum_{i=1}^n \tilde{X}_{nri}$ and $\bar{X}'_{nrj} = L_{nj}^{-1} \sum_{k=1}^{L_{nj}} \tilde{X}'_{nrjk}$. The weights in (8) depend on the number of data points in each bin, namely L_{nj} for $j = 1, \dots, m$. Note that the estimators proposed in (8) are method of moments estimators targeting $E\{\beta_0(U)\}, E\{\beta_r(U)\tilde{X}_r\}/E(\tilde{X}_r)$ and $E\{\eta_s(U)\}$, respectively. It follows then that they are consistent for the underlying parameters, formally stated in Section 4, since $E\{\beta_0(U)\} = \gamma_0, E\{\beta_r(U)\tilde{X}_r\}/E(\tilde{X}_r) = \gamma_r$ and $E\{\eta_s(U)\} = \delta_s$, by the relations in (6) and the identifiability conditions.

We note that the estimators $\hat{\gamma}_{n0}$ and $\hat{\gamma}_{nr}$ have the same form as the CAR estimators (Şentürk and Müller, 2005), whereas $\hat{\delta}_{ns}$ are different. Furthermore, a straightforward application of the CAR algorithm yield inconsistent estimators for δ_s under the more general PCAR model. To see this, denote the original CAR estimators for δ_s by $\{\delta_{ns}^*\}_{s=1}^q$. It follows from Şentürk and Müller (2005), that

$$\hat{\delta}_{ns}^* = \frac{1}{\bar{Z}_{ns}} \sum_{j=1}^m \frac{L_{nj}}{n} \hat{\eta}_{nsj} \bar{Z}'_{nsj},$$

where $\bar{Z}_{ns} = n^{-1} \sum_{i=1}^n Z_{nsi}$ and $\bar{Z}'_{nsj} = L_{nj}^{-1} \sum_{k=1}^{L_{nj}} Z'_{nsjk}$. The estimators $\hat{\delta}_{ns}^*$ do not target δ_s , instead they target $E\{\eta_s(U)Z_s\}/E(Z_s) = \delta_s E\{\psi(U)Z_s\}/E(Z_s) = \delta_s C_s$, where $C_s \equiv E\{\psi(U)Z_s\}/E(Z_s) = [\text{cov}\{\psi(U), Z_s\}/E(Z_s)] + 1$ can get arbitrary large as $E(Z_s)$ approaches zero. Note here that since Z_s is correlated with U , one can argue that a latent variable can also be created for Z_s that would be independent of U , similar to the construction of X_r and Y . An argument can be made that the dependence structure suggests $Z_s = \alpha(U, Z_s^*)$, for some latent variable Z_s^* that is independent of U , which can be $Z_{nsi} = \phi_s^*(U_{ni})Z_{nsi}^*$. In this case, C_s further simplifies to equal $\text{cov}\{\psi(U), \phi_s^*(U)\} + 1$. (This insight was provided by a reviewer.) From this last form of C_s , it follows immediately that if (1) the response is not distorted by U , i.e. $\psi(U) = 1$, or if (2) Z_s is independent of U , i.e. $\phi_s^*(U) = 1, C_s = 1$, in which case the CAR estimator for δ_s would be consistent. For the latter mentioned case (2) of assuming Z_s is independent of U , this is not realistic in most data applications, since U is the common confounder that correlates with all the observed variables in these distortion settings. In the data examples given in the Introduction, the undistorted predictor age is not independent of the confounders like BMI, body weight or height.

For the former mentioned case (1) when the response is undistorted, simpler adjustments given by Hwang (1986) and Iturria et al. (1999), for multiplicative measurement error only in the predictors, would also be applicable. Hwang (1986) proposes a consistent estimator for the regression coefficients by estimating and adjusting for the bias of the regular least squares estimator. The estimation assumes that consistent estimates of the moments of the measurement error are available. Iturria et al. (1999)

proposes two estimation methods, where the first considers specific distributional forms for the measurement error and the second also models the distribution of the unobserved predictors. The two approaches of Hwang and Iturria are similar to PCAR in assuming that the error is independent of the unobserved predictors and that it has a mean of one. The difference of PCAR from the previous two approaches is that no knowledge of the distributional forms or the specific moments are assumed. Instead, information from the observed covariate U , of which the measurement error $\phi_r(U)$ is a function of, is utilized in the proposed estimation procedure.

The CAR estimators are biased for only the case of undistorted predictors. In other words, if all variables are considered as distorted then CAR provides consistent estimators. This is the difference between CAR and PCAR. While the PCAR set-up allows researchers the choice to consider a subset of predictors as distorted, the CAR set-up requires that all predictors are considered as distorted in order to yield consistent estimates.

As the PCAR set-up allows researchers to consider a subset of predictors as distorted, an important issue is how to determine whether a predictor should be considered as distorted or undistorted. Note that this distinction cannot be made using a statistical/analytical approach via studying the dependence or the relations between the observed variables and U . This is because all observed variables are dependent on U in the model (correlate with U) *whether they are distorted or undistorted*. Hence, this decision instead should be made by considering the duality between (a) the decision/assumption on undistorted (distorted) variables and (b) the specification of the underlying model; more specifically, the consideration of the predictor choice in the underlying model.

More precisely, determining to include a predictor as distorted or undistorted corresponds to two completely different underlying models of interest, one involving Z_s and the other involving Z_s^* , respectively. Specification of the underlying model will depend on the specific interest of the researcher and the specific context of the application. For example, in the data analysis presented in Section 5, the goal is to uncover the relation between a diabetes marker and diastolic blood pressure adjusted for BMI. Hence, the diabetes marker and diastolic blood pressure are considered as adjusted/distorted variables. On the other hand, age and triceps skin fold thickness are considered unadjusted/undistorted, as we are interested in the direct effects of age and triceps skin fold thickness on the BMI-adjusted diabetes markers.

We also note here that the PCAR estimators given in (8) are consistent for the parameters of the underlying model (2) also under additive distortion. More precisely, consider the simple case of one distorted and one undistorted predictor. The regression model in (2) simplifies to $Y = \gamma_0 + \gamma_1 X + \delta_1 Z + e$, where the multiplicative error structure is replaced by the additive error structure, given by $\tilde{Y} = Y + \psi_a(U)$ and $\tilde{X} = X + \phi_a(U)$. The proposed PCAR estimators given in (8) are consistent even when the multiplicative error is replaced by additive error as described above. The above additive error model leads to the following specific varying coefficient model: $E(\tilde{Y}|\tilde{X}, Z, U) = \beta_0(U) + \beta_1(U)\tilde{X} + \eta_1(U)Z$, where $\beta_0(U) = \gamma_0 - \gamma_1 \phi_a(U) + \psi_a(U)$, $\beta_1(U) = \gamma_1$ and $\eta_1(U) = \delta_1$. The PCAR estimators given in (8), namely $\hat{\gamma}_0$, $\hat{\gamma}_1$ and $\hat{\delta}_1$, target

$$E\{\beta_0(U)\}, \quad E\{\beta_1(U)\tilde{X}\}/E(\tilde{X}), \quad \text{and} \quad E\{\eta_1(U)\},$$

respectively. This holds regardless of the specific error structure, whether it be additive or multiplicative. Furthermore, under the additive distortion model, we have that

$$E\{\beta_0(U)\} = \gamma_0, \quad E\{\beta_1(U)\tilde{X}\}/E(\tilde{X}) = \gamma_1 \quad \text{and} \quad E\{\eta_1(U)\} = \delta_1.$$

This follows since $E\{\psi_a(U)\} = E\{\phi_a(U)\} = 0$ in the additive distortion model, under the identifiability condition of no average distortion, i.e. $E(\tilde{Y}) = E(Y)$ and $E(\tilde{X}) = E(X)$. Thus, the PCAR estimators proposed in (8) are consistent for parameters of the underlying model also under additive distortion structure.

4. Asymptotic properties

We present the asymptotic distribution of the estimators $\hat{\gamma}_{n0}$, $\hat{\gamma}_{nr}$ and $\hat{\delta}_{ns}$ in (8) when the number of subjects n tends to infinity. As in typical smoothing applications, the number of bins $m = m(n)$ is required to satisfy $m \rightarrow \infty$, $n/(m \log n) \rightarrow \infty$ and $m/\sqrt{n} \rightarrow \infty$ as $n \rightarrow \infty$. We denote convergence in distribution by $\xrightarrow{\mathcal{D}}$ and convergence in probability by \xrightarrow{P} . Also let

$$(\tilde{\gamma}_{n0j}, \dots, \tilde{\gamma}_{npj}, \tilde{\delta}_{n1j}, \dots, \tilde{\delta}_{nqj})^T = (\mathbf{X}_{nj}^T \mathbf{X}'_{nj})^{-1} \mathbf{X}_{nj}^T \mathbf{Y}'_{nj} \tag{9}$$

denote the least squares estimators of the multiple regression of the unobserved data falling into B_{nj} , where the vectors \mathbf{X}'_{nj} and \mathbf{Y}'_{nj} are defined the same way as $\tilde{\mathbf{X}}'_{nj}$ and $\tilde{\mathbf{Y}}'_{nj}$, with \mathbf{X}'_{nrjk} and \mathbf{Y}'_{nj} replacing $\tilde{\mathbf{X}}'_{nrjk}$ and $\tilde{\mathbf{Y}}'_{nj}$, respectively. This quantity is not estimable, but will be used in the proof of the main results.

For the PCAR estimators given in (8) to be well defined, the least squares estimators given in (7) must exist for each bin B_{nj} , i.e. $\det(\tilde{\mathbf{X}}_{nj}^T \tilde{\mathbf{X}}'_{nj}) \neq 0$. Correspondingly, the estimators in (9) will exist under the condition that $\det(\mathbf{X}_{nj}^T \mathbf{X}'_{nj}) \neq 0$. The following theorems are given under event E_n , explicitly defined in the Appendix, which summarizes the above outlined conditions for the PCAR estimators to be well defined.

For the following theorems, we define the following notations: $\lambda_{\psi} = E\{\psi^2(U)\}$, $\lambda_{\phi} = E\{\phi^2(U)\}$, $\lambda_{\psi\phi_r} = E\{\psi(U)\phi_r(U)\}$, $\sigma_{\psi}^2 = \text{var}\{\psi(U)\}$, $m_{r,k} = E(X_r^k)$, $\boldsymbol{\chi}^T = (1, X_1, \dots, X_p, Z_1, \dots, Z_q)$, $\boldsymbol{\Gamma} = E(\boldsymbol{\chi}\boldsymbol{\chi}^T|U)$, $\boldsymbol{\mathcal{M}} = \boldsymbol{\Gamma}^{-1} \boldsymbol{\chi}\boldsymbol{\chi}^T \boldsymbol{\Gamma}^{-1}$, $\tilde{\boldsymbol{\Theta}}_{nj} = L_{nj}^{-1} \tilde{\boldsymbol{\chi}}_{nj}^T \tilde{\boldsymbol{\chi}}'_{nj}$ and $\omega_{n,\ell,k} = (\tilde{\boldsymbol{\Theta}}_{nj})_{\ell 1}^{-1} + (\tilde{\boldsymbol{\Theta}}_{nj})_{\ell 2}^{-1} \tilde{\chi}'_{n1jk} + \dots + (\tilde{\boldsymbol{\Theta}}_{nj})_{\ell, p+q+1}^{-1} Z'_{nqjk}$ for $\ell = 1, \dots, p + s + 1$ and $k = 1, \dots, L_{nj}$.

Theorem 1. Under the technical conditions (C1)–(C7) in Section 6, on event E_n with $\text{pr}(E_n) \rightarrow 1$ as $n \rightarrow \infty$,

$$\sqrt{n}(\hat{\gamma}_{nr} - \gamma_r) \xrightarrow{\mathcal{D}} \mathbb{N}(0, \sigma_{\gamma_r}^2), \quad 0 \leq r \leq p,$$

$$\sqrt{n}(\hat{\delta}_{ns} - \delta_s) \xrightarrow{\mathcal{D}} \mathbb{N}(0, \sigma_{\delta_s}^2), \quad 1 \leq s \leq q,$$

where

$$\sigma_{\gamma_0}^2 = \gamma_0^2 \sigma_{\psi}^2 + \sigma^2 E\{\psi^2(U) \cdot \mathcal{M}_{11}\},$$

$$\sigma_{\gamma_r}^2 = \frac{\gamma_r^2 (\lambda_{\psi} m_{r,2} - m_{r,1}^2) + \sigma^2 m_{r,1}^2 E\{\psi^2(U) \cdot \mathcal{M}_{r+1,r+1}\} - 2\gamma_r^2 (\lambda_{\psi} \phi_r m_{r,2} - m_{r,1}^2) - \gamma_r^2 \text{var}(\tilde{X}_r)}{m_{r,1}^2},$$

$$\sigma_{\delta_s}^2 = \delta_s^2 \sigma_{\psi}^2 + \sigma^2 E\{\psi^2(U) \cdot \mathcal{M}_{p+s+1,p+s+1}\} \text{ for } 1 \leq r \leq p \text{ and } 1 \leq s \leq q.$$

Theorem 1 establishes the asymptotic normality of the proposed PCAR estimators. The following theorem provides consistent estimators of the asymptotic variances given in Theorem 1.

Theorem 2. Under the technical conditions (C1)–(C7) in Section 6, on event E_n with $\text{pr}(E_n) \rightarrow 1$ as $n \rightarrow \infty$,

$$\hat{\sigma}_{n\gamma_r}^2 \xrightarrow{P} \sigma_{\gamma_r}^2, \quad 0 \leq r \leq p,$$

$$\hat{\sigma}_{n\delta_s}^2 \xrightarrow{P} \sigma_{\delta_s}^2, \quad 1 \leq s \leq q,$$

where

$$\hat{\sigma}_{n\gamma_0}^2 = \sum_{j=1}^m \frac{L_{nj}}{n} \hat{\beta}_{n0j}^2 - \hat{\gamma}_{n0}^2 + \hat{\sigma}_n^2 \frac{n^{-1} \sum_{j=1}^m \hat{\beta}_{n0j}^2 \sum_{k=1}^{L_{nj}} \omega_{n,1,k}^2}{\sum_{j=1}^m n^{-1} L_{nj} \hat{\beta}_{n0j}^2},$$

$$\hat{\sigma}_{n\gamma_r}^2 = \frac{n^{-1} \sum_{j=1}^m \hat{\beta}_{nrj}^2 \sum_{k=1}^{L_{nj}} \tilde{X}_{nrjk}^2 + \hat{\gamma}_{nr}^2 \tilde{X}_{nr}^2 - 2\hat{\gamma}_{nr} n^{-1} \sum_{j=1}^m \hat{\beta}_{nrj} \sum_{k=1}^{L_{nj}} \tilde{X}_{nrjk} + \hat{\gamma}_{nr}^2 s_{\tilde{X}_r}^2}{\tilde{X}_{nr}^2}$$

$$+ \hat{\sigma}_n^2 \frac{n^{-1} \sum_{j=1}^m \hat{\beta}_{n0j}^2 \tilde{X}_{nrj}^2 \sum_{k=1}^{L_{nj}} \omega_{n,r+1,k}^2}{\tilde{X}_{nr}^2 \sum_{j=1}^m n^{-1} L_{nj} \hat{\beta}_{n0j}^2},$$

$$\hat{\sigma}_{n\delta_s}^2 = \sum_{j=1}^m \frac{L_{nj}}{n} \hat{\eta}_{nsj}^2 - \hat{\delta}_{ns}^2 + \hat{\sigma}_n^2 \frac{n^{-1} \sum_{j=1}^m \hat{\beta}_{n0j}^2 \sum_{k=1}^{L_{nj}} \omega_{n,p+s+1,k}^2}{\sum_{j=1}^m n^{-1} L_{nj} \hat{\beta}_{n0j}^2},$$

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{j=1}^m \sum_{k=1}^{L_{nj}} \left(\tilde{Y}'_{nj} - \hat{\beta}_{n0j} - \sum_{r=1}^p \hat{\beta}_{nrj} \tilde{X}'_{nrjk} - \sum_{s=1}^q \hat{\eta}_{nsj} Z'_{nsjk} \right)^2$$

and

$$s_{\tilde{X}_r}^2 = \frac{1}{n-1} \sum_{i=1}^n (\tilde{X}_{nri} - \tilde{X}_{nr})^2.$$

Normalizing by the above consistent variance estimators, it holds that

$$\frac{\sqrt{n}}{\hat{\sigma}_{n\gamma_r}} (\hat{\gamma}_{nr} - \gamma_r) \xrightarrow{\mathcal{D}} \mathbb{N}(0, 1), \quad 0 \leq r \leq p \text{ and } \frac{\sqrt{n}}{\hat{\sigma}_{n\delta_s}} (\hat{\delta}_{ns} - \delta_s) \xrightarrow{\mathcal{D}} \mathbb{N}(0, 1), \quad 1 \leq s \leq q.$$

Therefore, the approximate $(1 - \alpha)100\%$ asymptotic confidence intervals for γ_r and δ_s have the endpoints

$$\hat{\gamma}_{nr} \pm z_{\alpha/2} \frac{\hat{\sigma}_{n\gamma_r}}{\sqrt{n}} \text{ and } \hat{\delta}_{ns} \pm z_{\alpha/2} \frac{\hat{\sigma}_{n\delta_s}}{\sqrt{n}}, \tag{10}$$

where $z_{\alpha/2}$ is the $(1 - \alpha/2)$ th quantile of the standard Gaussian distribution.

Remark. These proposed variance estimators are motivated by the identifiability conditions, the definition of the smooth varying coefficients functions given in (6), Lemmas 3 and 4(a). Using the consistency of $\hat{\beta}_{nrj}$ and $\hat{\eta}_{nsj}$ for the values of the functions β_r and η_s at the midpoint of the j th bin and the definitions of \tilde{Y}'_{nj} and \tilde{X}'_{nrjk} , we target the quantities $\sigma^2 \lambda_\psi$, $\gamma_0^2 \lambda_\psi$, $\delta_s^2 \lambda_\psi$, $\gamma_r^2 \lambda_\psi m_{r,2}$ and $\gamma_r^2 \lambda_\psi \phi_r m_{r,2}$ with the estimators $\hat{\sigma}_n^2$, $\sum_{j=1}^m n^{-1} L_{nj} \hat{\beta}_{n0j}^2$, $\sum_{j=1}^m n^{-1} L_{nj} \hat{\eta}_{nsj}^2$, $n^{-1} \sum_{j=1}^m \hat{\beta}_{nrj}^2 \sum_{k=1}^{L_{nj}} \tilde{X}'_{nrjk}{}^2$ and $n^{-1} \hat{\gamma}_{nr} \sum_{j=1}^m \hat{\beta}_{nrj} \sum_{k=1}^{L_{nj}} \tilde{X}'_{nrjk}{}^2$, respectively. Furthermore, relying mainly on Lemmas 3 and 4(a), we target $\gamma_0^2 E\{\psi^2(U) \cdot \mathcal{M}_{11}\}$, $\gamma_0^2 m_{r,1}^2 E\{\psi^2(U) \cdot \mathcal{M}_{r+1,r+1}\}$ and $\gamma_0^2 E\{\psi^2(U) \cdot \mathcal{M}_{p+s+1,p+s+1}\}$ with $n^{-1} \sum_{j=1}^m \hat{\beta}_{n0j}^2 \sum_{k=1}^{L_{nj}} \omega_{n,1,k}^2$, $n^{-1} \sum_{j=1}^m \hat{\beta}_{nrj}^2 \sum_{k=1}^{L_{nj}} \omega_{n,r+1,k}^2$ and $n^{-1} \sum_{j=1}^m \hat{\beta}_{n0j}^2 \sum_{k=1}^{L_{nj}} \omega_{n,p+s+1,k}^2$, respectively.

5. Application to the Pima Indians diabetes data

We illustrate the proposed PCAR methodology with an application to the Pima Indians diabetes data set, available at <http://www.ics.uci.edu/~mllearn>. Obesity is an important contributing factor to diabetes and has been widely studied in the Pima Indians population (Smith et al., 1988; Knowler et al., 1991; Hanson et al., 1998). One-half of adult Pima Indians have diabetes and 95% of those with diabetes are overweight (National Institute of Diabetes and Digestive and Kidney Diseases, <http://diabetes.niddk.nih.gov>). The available data come from a larger database, where the subgroup used consists of $n = 524$ females at least 21 years old and of Pima Indian heritage. (The population lives near Phoenix, AZ, U.S.A.) An oral glucose tolerance test is one of the diagnostic tests for type II diabetes. The goal is to uncover the underlying, BMI adjusted, regression relation, $PGC = \gamma_0 + \gamma_1 DBP + \delta_1 Age + \delta_2 TSFT + e$, based on the observed plasma glucose concentration (PGC; from a oral glucose tolerance test), diastolic blood pressure (DBP), triceps skin fold thickness (TSFT), age and BMI. We chose to adjust only the main relation of interest, namely the one between plasma glucose concentration (the response) and diastolic blood pressure for BMI, and included age and triceps skin fold thickness as unadjusted predictors as they are commonly accounted factors in studies on diabetes.

Table 1 gives the regression coefficient estimates for $(\gamma_0, \gamma_1, \delta_1, \delta_2)$ using the proposed PCAR method, CAR method, the ordinary least squares (OLS) estimates from regressing the observed PGC on (DBP, Age, TSFT) without adjusting for the confounder BMI, and adjustment via division, i.e. regressing \widehat{PGC}/BMI on $(\widehat{DBP}/BMI, Age, TSFT)$. The approximate 95% asymptotic confidence intervals for the regression parameters obtained through all three methods are also displayed. The approximate confidence intervals for PCAR estimates were obtained as proposed in (10).

The implementation of the binning algorithm allows for merging of sparsely populated bins. Bin widths were chosen such that there are at least $(p + q + 1)$ points, enough to fit the linear regression with $(p + q)$ predictors in each bin. If there were bins with less than $(p + q + 1)$ elements, such bins were randomly merged with neighboring bins. The merging algorithm is randomized to avoid the introduction of any additional bias. It starts by merging bins with no points. If there are more than one such bin, it randomly picks one and merges it with its neighbor of smallest number of points. After merging all the bins with no points, the bins with one point and eventually bins with $p + q$ points are merged. For this example with $n = 524$ (after the removal of outliers), the average number of points per bin was 15, yielding a total of 34 bins after merging. Note that CAR estimates have been shown to be sufficiently robust regarding different choices of m , under the rate conditions given in Section 4 (Şentürk and Müller, 2006). We have found this property to hold for the proposed PCAR estimates as well, where the range of m values yielding robust estimates for different sample sizes is given explicitly in the next section.

Note that coefficients obtained by adjustment via division are quite different from the other three methods applied. In this adjustment the coefficient of DBP becomes quite pronounced compared to the other two predictors. This is most likely due to the pseudo dependence created between \widehat{PGC}/BMI and \widehat{DBP}/BMI via division by the common variable BMI. This is an example of the misleading conclusions that adjustment by division may suggest. In other words, if the original contamination is not exactly multiplication by the confounder (BMI in this example), then normalization by division may create further confounding, or “coupling” (as defined in Archie, 1981), creating a pseudo dependence that does not exist in the original data.

Even though OLS estimates for blood pressure and age are different from the PCAR and CAR estimates, all are found statistically significant at the usual 5% level. Thus, diastolic blood pressure and age are still important predictors of PGC even after adjusted for BMI. However, using OLS, TSFT is a significant predictor of PGC, but it is not significant using PCAR and CAR at the 5% significance

Table 1

Parameter estimates for the regression model $PGC = \gamma_0 + \gamma_1 DBP + \delta_1 Age + \delta_2 TSFT + e$, obtained by least squares regression of $\tilde{Y} = \widehat{PGC}$ (plasma glucose concentration) on $\tilde{X}_1 = \widehat{DBP}$ (diastolic blood pressure), $Z_1 = Age$ and $Z_2 = TSFT$ (triceps skin fold thickness), and by adjustment through division, i.e. regressing \widehat{PGC}/BMI on \widehat{DBP}/BMI , Age and TSFT, and alternatively by PCAR and CAR, for $n = 524$ subjects

Coeff.	Least sq. reg.		Adj. by division		PCAR		CAR	
	Est.	95% CI	Est.	95% CI	Est.	95% CI	Est.	95% CI
γ_0	70.6	(54.9, 86.2)	2.2	(1.7, 2.8)	75.7	(51.2, 100.1)	75.7	(53.4, 97.9)
DBP	0.25	(0.02, 0.48)	0.72	(0.54, 0.90)	0.42	(0.13, 0.71)	0.42	(0.14, 0.70)
Age	0.64	(0.38, 0.89)	0.01	(0.00, 0.02)	0.47	(0.17, 0.78)	0.48	(0.18, 0.78)
TSFT	0.42	(0.16, 0.68)	-0.02	(-0.03, -0.01)	0.22	(-0.22, 0.66)	0.08	(-0.36, 0.51)

Confidence intervals at the 95% level were obtained by the standard t -statistic for least squares regression and adjustment through division, by the proposed asymptotic intervals (10) for PCAR and by asymptotic intervals given in Şentürk and Müller (2006) for CAR.

level. This result is not too surprising, since both TSFT and BMI are indicators of obesity. They are positively correlated (Pearson correlation 0.67). Thus, adjusting for one, the other becomes an insignificant factor for predicting plasma glucose concentration. We note that even though estimation via CAR leads to the same conclusion as PCAR on the significance of the predictors for this analysis, the estimates from these two methods are different for TSFT. This is again to be expected, since CAR estimates are shown to be biased for the undistorted predictors.

6. Numerical studies

To examine the numerical properties of the estimators, we implemented the following simulation studies. The underlying multiple regression model is

$$Y = \gamma_0 + \gamma_1 X_1 + \gamma_2 X_2 + \delta Z + e, \tag{11}$$

where the parameters of interest are $(\gamma_0, \gamma_1, \gamma_2, \delta)^T = (4, -1, 0.3, 3)$. The error variable is $e \sim N(0, 0.5)$, and the confounder variable U is generated from a uniform distribution on $[2, 6]$. We considered the joint distribution of the predictors to be multivariate normal: $(X_1, X_2, Z)^T \sim N_3(\mu, \Sigma)$, with a general covariance structure

$$\Sigma = \begin{bmatrix} 0.490 & 0.168 & 0.280 \\ 0.168 & 1.440 & -0.360 \\ 0.280 & -0.360 & 1.000 \end{bmatrix}.$$

The mean vector is $\mu = (0.7, 1.2, |U| - 3.5)^T$, so that the undistorted predictor Z is dependent on U . To simulate the distorted (observed) data, we consider the following distorting functions, $\psi(U) = (U+3)/7$, $\phi_1(U) = (U+1)^2/26.3333$, and $\phi_2(U) = (U+10)/14$, satisfying the identifiability constraints that $E\{\psi(U)\} = 1$ and $E\{\phi_r(U)\} = 1$. The distorted response and predictors are $\tilde{Y} = \psi(U)Y$, $\tilde{X}_1 = \phi_1(U)X_1$, and $\tilde{X}_2 = \phi_2(U)X_2$. Under this simulation setting, we examine (1) the confidence interval coverage levels based on the asymptotic results and (2) the finite sample bias of the estimators, as well as comparing CAR and PCAR estimators in terms of variance and MSE.

We conducted 1000 Monte Carlo simulation runs for sample sizes $n = 100, 150, 350, 800$, and 1400 to study the approximate asymptotic confidence intervals given in (10). For the sample sizes $n = 100, 150, 350, 800$ and 1400, the total number of bins formed were $m = 16, 27, 32, 50$ and 70. Table 2 summarizes the coverage and interval lengths, averaged over the 1000 simulation runs, for the approximate 95% asymptotic confidence intervals for the parameter vector $(\gamma_0, \gamma_1, \gamma_2, \delta)^T = (4, -1, 0.3, 3)$. The numerical study indicates that the estimated non-coverage percentages are close to the target value of 0.05, as the sample size n increases. The estimated interval lengths are decreasing as n increases, as expected.

We also examined the bias, variance and mean squared error (MSE) of the proposed estimators in comparison to the CAR estimators. For example, the estimated (absolute bias, variance, MSE) values for PCAR estimators at the smallest sample size $n = 100$ are (0.0112, 0.2223, 0.2224), (0.0120, 0.1392, 0.1393), (0.0028, 0.0421, 0.0421) and (0.0167, 0.0651, 0.0654) for $\hat{\gamma}_0, \hat{\gamma}_1, \hat{\gamma}_2$ and $\hat{\delta}$, respectively. These values are averages over 1000 Monte Carlo runs. The results are similar for other sample sizes, where the variance seems to be the dominating factor contributing to the MSE. The estimated (bias, variance, MSE) values for the CAR estimator for δ (even though their asymptotic distributions are different, the three other point estimates for γ_0, γ_1 and γ_2 are the same for the two methods), $\hat{\delta}^*$, are (1.294, 1.694, 3.368) at the same sample size of $n = 100$. The multiplicative bias factor of the CAR estimate for δ , shown to be $C_s = E\{\psi(U)Z_s\}/E(Z_s)$ in Section 3, is equal to 1.416 for this simulation set-up. As expected, the CAR estimate $\hat{\delta}^*$ is off target for $\delta = 3$ (with a mean of 4.294 at $n = 100$, and 4.146 at $n = 1400$). In addition to being biased, note that the CAR estimator $\hat{\delta}^*$ has substantially larger variance relative to the PCAR estimator.

As stated above, for sample sizes $n = 100, 150, 350, 800$ and 1400, the total number of bins formed were $m = 16, 27, 32, 50$ and 70, respectively. We carried out additional simulation studies to examine the affect of the number of bins m . The results suggest that the estimators are robust, based on estimated MSE, when m is chosen in the intervals $[13, 18], [18, 27], [25, 45], [35, 65]$ and $[60, 90]$ corresponding to sample sizes $n = 100, 150, 350, 800$ and 1400. In a given application, the above intervals can give rough guidelines on how the choice of m may change with sample size, although a sensitivity analysis for the choice of m specific to the data would also be informative.

Table 2

Estimated coverage (in percent) and average interval lengths for the approximate 95% confidence intervals formed for the parameters of the regression model (11), corresponding to sample sizes $n = 100, 150, 350, 800$, and 1400

n	γ_0		γ_1		γ_2		δ	
	Coverage	Length	Coverage	Length	Coverage	Length	Coverage	Length
100	87.5	1.15	90.4	0.99	90.8	0.52	92.2	0.68
150	88.6	0.92	92.3	0.75	91.3	0.41	92.4	0.54
350	90.0	0.46	93.9	0.36	91.2	0.19	93.0	0.26
800	93.7	0.26	94.0	0.19	92.1	0.10	93.2	0.15
1400	94.3	0.19	94.3	0.14	94.7	0.08	94.8	0.11

Results are based on 1000 Monte Carlo simulation runs for each sample size.

Finally, we note that even though the smallest sample size at which the proposed asymptotic confidence intervals attain reasonable coverage in our simulation study is $n = 100$ (the coverage is around 5% off the targeted level for $n = 100$), the proposed PCAR point estimators still yield reasonable (bias, variance and MSE) values: (0.0515, 0.5748, 0.5775), (0.0026, 0.3552, 0.3553), (0.0256, 0.0916, 0.0922) and (0.0205, 0.3328, 0.3332) for $\gamma_0 = 4, \gamma_1 = -1, \gamma_2 = 0.3, \delta = 3$, respectively, at $n = 50$. These values are estimated under the current simulation set-up with three predictor variables for $n = 50$ and $m = 8$. For the case of simple linear regression, roughly the same number of bins can be attained with $n = 30$, and hence $n = 30$ would be the smallest sample size where CAR and PCAR give reasonable point estimates. For sample sizes smaller than 30, systematic localization via binning may not be fully feasible so the PCAR estimates should be taken with caution. In this case a very rough localization by stratification (by U) into 2–3 groups for crude comparisons is possible.

7. Proofs of the main results

We provide the major steps of the proofs of the main results (Theorems 1 and 2) here and defer the auxiliary results for these proofs to the Appendix, where they are listed as Lemmas 1–4. We introduce the following technical conditions:

- (C1) The covariate U is bounded below and above, $-\infty < a \leq U \leq b < \infty$ for real numbers $a < b$. The density $f(u)$ of U satisfies $\inf_{a \leq u \leq b} f(u) > c_1 > 0, \sup_{a \leq u \leq b} f(u) < c_2 < \infty$ for real c_1, c_2 , and is uniformly Lipschitz continuous, i.e., there exists a real number M such that $\sup_{a \leq u \leq b} |f(u+c) - f(u)| \leq M|c|$ for any real number c .
- (C2) The variables (e, U, X_r) are mutually independent for $r = 1, \dots, p$. In addition, (e, Z_s) are assumed to be independent.
- (C3) For the predictors, $\sup_{1 \leq i \leq n, 1 \leq r \leq p, 1 \leq s \leq q} \{|X_{nri}|, |Z_{nsi}|\} \leq B$ for some bound $B \in \mathbb{R}$. In addition, the predictors X_r satisfy the condition that $E(X_r) \neq 0$.
- (C4) Contamination functions $\psi(\cdot)$ and $\phi_r(\cdot), 1 \leq r \leq p$, are twice continuously differentiable, satisfying $E\psi(U) = 1, E\phi_r(U) = 1$, and $\phi_r(\cdot) > 0, 1 \leq r \leq p$.
- (C5) The matrices $\Gamma_{nj}, j = 1, \dots, m$ are nonsingular, i.e. $\rho = |\inf_j \det(\Gamma_{nj})| > 0$, where $\Gamma_{nj} = E(L_{nj}^{-1} \chi_{nj}^T \chi_{nj}^* | U_{nj}^*)$, $\chi_{nj}^* = (\chi_{nj1}^*, \dots, \chi_{njL_{nj}}^*)^T$ is a $L_{nj} \times (p + q + 1)$ undistorted data matrix in bin j , and $\chi_{nj}^* = (1, X'_{n1jk}, \dots, X'_{npjk}, Z'_{n1jk}, \dots, Z'_{nqjk})^T$ denotes the k th observation.

The technical conditions above are similar to those introduced in Şentürk and Müller (2006), except for the new independence structure outlined in (C2), the boundedness of the undistorted predictors Z_s in (C3), and the bin dependent limiting matrices Γ_{nj} in (C5), resulting from the dependence structure between Z_s and U . Bounded covariates are standard in asymptotic theory for least squares regression, as are conditions (C2) and (C5) (see Lai et al., 1979). The identifiability conditions stated in (C4) are equivalent to $E(\tilde{Y}|X) = E(Y|X)$ and $E(\tilde{X}_r|X_r) = X_r$.

In the proofs of the main results, the following notations will be utilized.

- 1. $A \square B$: The Hadamard product of two matrices, A and B , of the same dimension. The matrix $A \square B$ is also of the same dimension with (i, j) th element equal to the product of the (i, j) th elements of matrices A and B .
- 2. $\mathbf{1}_{a \times b}$: A matrix of size $a \times b$ with all entries equal to one.
- 3. $\hat{\theta}_{nj} = (\hat{\beta}_{n0j}, \hat{\beta}_{n1j}, \dots, \hat{\beta}_{npj}, \hat{\eta}_{n1j}, \dots, \hat{\eta}_{nqj})^T$.
- 4. $\tilde{\theta}_{nj} = (\tilde{\gamma}_{n0j} \psi(U_{nj}^*), \tilde{\gamma}_{n1j} \psi(U_{nj}^*) / \phi_1(U_{nj}^*), \dots, \tilde{\gamma}_{npj} \psi(U_{nj}^*) / \phi_p(U_{nj}^*), \tilde{\delta}_{n1j} \psi(U_{nj}^*), \dots, \tilde{\delta}_{nqj} \psi(U_{nj}^*))^T$.
- 5. $\theta_{nj} = (\gamma_0 \psi(U_{nj}^*), \gamma_1 \psi(U_{nj}^*) / \phi_1(U_{nj}^*), \dots, \gamma_p \psi(U_{nj}^*) / \phi_p(U_{nj}^*), \delta_1 \psi(U_{nj}^*), \dots, \delta_q \psi(U_{nj}^*))^T$.
- 6. We use $\chi_{nj(i)}$ to denote the matrix χ_{nj}^* and $L_{nj(i)}$ to denote the number of points in the j th bin such that $U_{ni} \in B_{nj}$, and $\kappa_{rk(i)} \equiv \{(L_{nj(i)}^{-1} \chi_{nj(i)}^T \chi_{nj(i)}^*)^{-1} \chi_{nj(i)}^T\}_{rk(i)}$ is the (r, k) th element of the matrix $\{(L_{nj}^{-1} \chi_{nj}^T \chi_{nj}^*)^{-1} \chi_{nj}^T\}$ for $1 \leq r \leq p + q + 1$, where $U_{ni} = U'_{nj k}$ is the k th element in the ordered sample $(U'_{nj1}, \dots, U'_{njL_{nj}}) \in B_{nj}$.
- 7. $\Theta_{nj} = L_{nj}^{-1} \chi_{nj}^T \chi_{nj}^*$.

Proof of Theorem 1. From Lemma 4(b), we have that

$$\sup_j |(L_{nj}^{-1} \tilde{\chi}_{nj}^T \tilde{Y}_{nj}) - \{\Delta \square (L_{nj}^{-1} \chi_{nj}^T Y_{nj}^*)\}| = O_p(m^{-1}) \mathbf{1}_{(p+q+1) \times 1}, \tag{12}$$

where $\Delta = \{\psi(U_{nj}^*), \psi(U_{nj}^*) / \phi_1(U_{nj}^*), \dots, \psi(U_{nj}^*) / \phi_p(U_{nj}^*), \psi(U_{nj}^*), \dots, \psi(U_{nj}^*)\}^T$.

Lemma 3 together with (12) implies that, on event E_n ,

$$\sup_j |\hat{\theta}_{nj} - \tilde{\theta}_{nj}| = O_p(m^{-1}) \mathbf{1}_{(p+q+1) \times 1}. \tag{13}$$

We first consider the case of $r = 0$ and show that $\sqrt{n}(\hat{\gamma}_{n0} - \gamma_0)$ is asymptotically normal. Using Lemma 4, (13), and some algebra, $\sqrt{n}(\hat{\gamma}_{n0} - \gamma_0) = \sqrt{n}(\sum_{j=1}^m L_{nj}n^{-1}\hat{\beta}_{n0j} - \gamma_0)$ can be expressed as

$$\sum_{j=1}^m \sum_{k=1}^{L_{nj}} \left[\frac{\gamma_0 \psi(U'_{nj}k)}{\sqrt{n}} + \frac{\psi(U'_{nj}k)e'_{nj}k}{\sqrt{n}} \{(L_{nj}^{-1} \mathbf{X}'_{nj} \mathbf{X}_{nj})^{-1} \mathbf{X}'_{nj} \mathbf{1}_k\} - \sqrt{n} \gamma_0 + O_p(\sqrt{n}/m) \right].$$

Since the above sum is over all bins indexed by j , and over all points within the bins indexed by k , it is equal to the sum over all data points indexed by i , summed up in a random order. Thus, the above expression for $\sqrt{n}(\hat{\gamma}_{n0} - \gamma_0)$ can be further simplified to

$$\sum_{\substack{i=1 \\ j,k}}^t \left[\frac{\gamma_0 \psi(U_{ni})}{\sqrt{n}} + \frac{\psi(U_{ni})e_{ni} \kappa_{1k(i)}}{\sqrt{n}} - \frac{\gamma_0}{\sqrt{n}} \right] + O_p(\sqrt{n}/m) \equiv \sum_{i=1}^t W_{n0i} + O_p(\sqrt{n}/m). \tag{14}$$

Therefore, $\sqrt{n}(\hat{\gamma}_{n0} - \gamma_0)$ is asymptotically equivalent to $S_{n0t} \equiv \sum_{i=1}^t W_{n0i}$ because the second term $O_p(\sqrt{n}/m)$ is negligible when $m/\sqrt{n} \rightarrow \infty$ as $n \rightarrow \infty$. Next, let F_{n0t} be the σ -field generated by $\{e_{n1}, \dots, e_{nt}, U_{n1}, \dots, U_{nt}, L_{nj(1)}, \dots, L_{nj(t)}, \mathbf{X}'_{nj(1)}, \dots, \mathbf{X}'_{nj(t)}\}$. Then $\{S_{n0t}, F_{n0t}, 1 \leq t \leq n\}$ is a mean zero martingale for $n \geq 1$, since $E(S_{n0t}) = 0$, $E(S_{n0,t+1} | F_{n0t}) = S_{n0t}$, and S_{n0t} is adapted to F_{n0t} . Furthermore, note that the σ -fields are nested, that is, $F_{n0t} \subseteq F_{n0,t+1}$ for all $t \leq n$. Hence, it follows from Lemma 1 that $S_{n0n} \rightarrow \mathbb{N}(0, \sigma_{\gamma_0}^2)$ in distribution (McLeish, 1974, Theorem 2.3 and subsequent discussion). This establishes the asymptotic normality of $\sqrt{n}(\hat{\gamma}_{n0} - \gamma_0)$.

We proceed next to establish the asymptotic normality of $\sqrt{n}(\hat{\gamma}_{nr} - \gamma_r)$ for $r = 1, \dots, p$. Let $\hat{v}_{nr} = \sum_{j=1}^m L_{nj}n^{-1}\hat{\beta}_{nrj} \bar{X}'_{nrj}$ and $\bar{v}_{nr} = \sum_{j=1}^m L_{nj}n^{-1}\bar{X}'_{nrj}$ and note that $\hat{\gamma}_{nr} = \hat{v}_{nr}/\bar{v}_{nr}$. We first show that

$$\sqrt{n} \begin{pmatrix} \hat{v}_{nr} - \gamma_r E(X_r) \\ \bar{v}_{nr} - E(X_r) \end{pmatrix} \xrightarrow{\mathcal{D}} \mathbb{N}_2(\mathbf{0}, \Sigma_r). \tag{15}$$

For (15) to hold, by the Cramér–Wold device, it is enough to show the asymptotic normality of

$$\sqrt{n}[a(\hat{v}_{nr} - \gamma_r E(X_r)) + b(\bar{v}_{nr} - E(X_r))] \tag{16}$$

for real a, b . The asymptotic normality of $\sqrt{n}(\hat{\gamma}_{nr} - \gamma_r)$ ($1 \leq r \leq p$) will follow from (15) by applying the δ -method with $\hat{\gamma}_{nr} = \hat{v}_{nr}/\bar{v}_{nr}$. Again applying Lemma 4 together with (13) and some simple algebra, we can express \hat{v}_{nr} and \bar{v}_{nr} as

$$\hat{v}_{nr} = \sum_{j=1}^m \sum_{k=1}^{L_{nj}} \left[\frac{\gamma_r \psi(U'_{nj}k) \mathbf{X}'_{nrjk}}{n} + \frac{\bar{X}'_{nrj} \psi(U'_{nj}k) e'_{nj}k}{n} \{(L_{nj}^{-1} \mathbf{X}'_{nj} \mathbf{X}_{nj})^{-1} \mathbf{X}'_{nj} \mathbf{1}_{rk}\} \right] + O_p(m^{-1})$$

and

$$\bar{v}_{nr} = \sum_{j=1}^m \sum_{k=1}^{L_{nj}} \frac{1}{n} \phi_r(U'_{nj}k) \mathbf{X}'_{nrjk} + O_p(m^{-1}).$$

Thus, using similar simplifications as was done for the case of $r = 0$ in (14), the linear combination (16), namely $\sqrt{n}[a(\hat{v}_{nr} - \gamma_r E(X_r)) + b(\bar{v}_{nr} - E(X_r))]$, can be expressed as

$$\sum_{\substack{i=1 \\ j,k}}^n \left[a \frac{\gamma_r}{\sqrt{n}} \psi(U_{ni}) \mathbf{X}_{nri} + a \frac{\bar{X}'_{nrj(i)}}{\sqrt{n}} \psi(U_{ni}) e_{ni} \kappa_{rk(i)} - a \frac{\gamma_r}{\sqrt{n}} E(X_r) + \frac{b}{\sqrt{n}} \phi_r(U_{ni}) \mathbf{X}_{nri} - b \frac{E(X_r)}{\sqrt{n}} \right] + O_p(\sqrt{n}/m) \equiv \sum_{i=1}^t W_{nri} + O_p(\sqrt{n}/m).$$

The second term $O_p(\sqrt{n}/m)$ is asymptotically negligible and it is straightforward to verify that $\{S_{nrt} \equiv \sum_{i=1}^t W_{nri}, F_{n0t}, 1 \leq t \leq n\}$ is a mean zero martingale for $n \geq 1$. Analogous to the case of $r = 0$, described in more details earlier, it follows from Lemma 2 that $S_{nrm} \xrightarrow{\mathcal{D}} \mathbb{N}(0, (a, b) \Sigma_r (a, b)^T)$. Finally, a direct application of the δ -method gives $\sqrt{n}(\hat{\gamma}_{nr} - \gamma_r) \xrightarrow{\mathcal{D}} \mathbb{N}(0, \sigma_{\gamma_r}^2)$ for $1 \leq r \leq p$, where $\sigma_{\gamma_r}^2$ is explicitly given in Theorem 1.

The asymptotic normality of $\sqrt{n}(\hat{\delta}_{ns} - \delta_s) \rightarrow \mathbb{N}(0, \sigma_{\delta_s}^2)$ follows similarly to the case of $\sqrt{n}(\hat{\gamma}_{n0} - \gamma_0)$, since they have similar forms in (13). (See also definition/notation 4.) The asymptotic variance $\sigma_{\delta_s}^2$ which has a similar form as $\sigma_{\gamma_0}^2$, is given explicitly in Theorem 1. This completes the proof of Theorem 1. \square

Proof of Theorem 2. The following relation holds on event A_n

$$\sup_j |(\tilde{\gamma}_{n0j} - \gamma_0, \dots, \tilde{\gamma}_{npj} - \gamma_p, \tilde{\delta}_{n1j} - \delta_1, \dots, \tilde{\delta}_{nqj} - \delta_q)^\top| = o_p(1) \mathbf{1}_{(p+q+1) \times 1}. \tag{17}$$

It follows from Lemma 4(a) and (b). Utilizing (17) together with (13) gives

$$\sup_j |\hat{\theta}_{nj} - \theta_{nj}| = o_p(1) \mathbf{1}_{(p+q+1) \times 1}. \tag{18}$$

By the Law of Large Numbers, (18), and boundedness considerations

$$\begin{aligned} \hat{\sigma}_n^2 &= \frac{1}{n} \sum_{j=1}^m \sum_{k=1}^{L_{nj}} \left(\tilde{Y}'_{nj k} - \hat{\beta}_{n0j} - \sum_{r=1}^p \hat{\beta}_{nrj} \tilde{X}'_{nrjk} - \sum_{s=1}^q \hat{\eta}_{nsj} Z'_{nsjk} \right)^2 \\ &= \frac{1}{n} \sum_{j=1}^m \sum_{k=1}^{L_{nj}} \{ \psi(U_{nj}^*) e'_{nj k} + o_p(1) \}^2 = \frac{1}{n} \sum_{i=1}^n \psi(U_{ni}^2) e_{ni}^2 + o_p(1) = \sigma^2 \lambda_\psi + o_p(1). \end{aligned}$$

It has been shown in Şentürk and Müller (2006) that

$$\begin{aligned} A_1 &\equiv \frac{1}{n} \sum_{j=1}^m \hat{\beta}_{nrj}^2 \sum_{k=1}^{L_{nj}} \tilde{X}_{nrjk}^2 = \gamma_r^2 \lambda_\psi m_{r,2} + o_p(1), \\ A_2 &\equiv \frac{1}{n} \sum_{j=1}^m \hat{\beta}_{nrj} \sum_{k=1}^{L_{nj}} \tilde{X}_{nrjk}^2 = \gamma_r \lambda_\psi \phi_r m_{r,2} + o_p(1), \\ A_3 &\equiv \sum_{j=1}^m \frac{L_{nj}}{n} \hat{\beta}_{n0j}^2 = \gamma_0^2 \lambda_\psi + o_p(1), \end{aligned}$$

and it can be shown similarly that $A_4 \equiv \sum_{j=1}^m n^{-1} L_{nj} \hat{\eta}_{nsj}^2 = \delta_s^2 \lambda_\psi + o_p(1)$. Also, using Lemmas 3 and 4(a), (b) and the Law of Large Numbers, we have

$$\begin{aligned} A_5 &\equiv \frac{1}{n} \sum_{j=1}^m \hat{\beta}_{n0j}^2 \sum_{k=1}^{L_{nj}} \omega_{n,1,k}^2 \xrightarrow{P} \gamma_0^2 E\{\psi^2(U) \cdot \mathcal{M}_{11}\}, \\ A_6 &\equiv \frac{1}{n} \sum_{j=1}^m \hat{\beta}_{n0j} \tilde{X}_{nrj}^2 \sum_{k=1}^{L_{nj}} \omega_{n,r+1,k}^2 \xrightarrow{P} \gamma_0^2 \{E(X_r)\}^2 E\{\psi^2(U) \cdot \mathcal{M}_{r+1,r+1}\}, \end{aligned}$$

and

$$A_7 \equiv \frac{1}{n} \sum_{j=1}^m \hat{\beta}_{n0j}^2 \sum_{k=1}^{L_{nj}} \omega_{n,p+s+1,k}^2 \xrightarrow{P} \gamma_0^2 E\{\psi^2(U) \cdot \mathcal{M}_{p+s+1,p+s+1}\}.$$

The estimators of the asymptotic variances given in Theorem 2, in terms of the above quantities, are: (1) $\hat{\sigma}_{n\gamma_0}^2 = A_3 - \hat{\gamma}_{n0}^2 + \hat{\sigma}_n^2 A_5 / A_3$, (2) $\hat{\sigma}_{n\gamma_r}^2 = (A_1 + \hat{\gamma}_{nr}^2 \tilde{X}_{nr}^2 - 2\hat{\gamma}_{nr} \tilde{X}_{nr} A_2 + \hat{\gamma}_{nr}^2 s_{\tilde{X}_r}^2) / \tilde{X}_{nr}^2 + (\hat{\sigma}_n^2 A_6) / (\tilde{X}_{nr}^2 A_3)$, and (3) $\hat{\sigma}_{n\delta_s}^2 = A_4 - \hat{\delta}_{ns}^2 + (\hat{\sigma}_n^2 A_7) / A_3$. Thus, the first part of Theorem 2 follows by noting that $\hat{\gamma}_{n0} \xrightarrow{P} \gamma_0$, $\hat{\gamma}_{nr} \xrightarrow{P} \gamma_r$, $\hat{\delta}_{ns} \xrightarrow{P} \delta_s$, $s_{\tilde{X}_r}^2 \xrightarrow{P} \text{var}(\tilde{X}_r)$ and $\tilde{X}_{nr} \xrightarrow{P} E(X_r)$. Asymptotic confidence intervals given in the second part of the Theorem follow immediately from Theorem 1 and Slutsky's theorem using the consistent variance estimators. \square

8. Discussion

In this work we extend covariate adjusted regression (CAR) models to partial covariate adjusted regression (PCAR) models that allow for the specification of the effects of unadjusted predictors. Asymptotic normality of the proposed estimators are derived.

The PCAR (and CAR) estimation approach was designed to estimate the underlying regression relationship directly, bypassing the estimation of the exact distortion forms. Although the current approach leads to estimators that are simple to implement with known asymptotic properties and good finite sample performance, it does not provide direct estimates of the distorting functions. If the primary interest is in the distorting functions then alternative approaches are needed. A potential alternative approach would involve considering refined estimators for the varying coefficient functions to be used in targeting the distorting functions. However, nontrivial work is required and this remains largely an open problem.

Another interesting and relevant issue, brought to light by a reviewer, is an alternative analysis of the diabetes data given in Section 5. The analysis proceeds by considering the observed variable $Z_s = \text{Age}$ as $Z_s = \phi_r^*(U)Z_s^*$, where $U = \text{BMI}$ (body mass index) and Z_s^* is a latent variable that can be interpreted as “core biological wear”. This implies that the underlying model of interest would involve the latent construct of core biological wear instead of age (defined to be the length of time that a person has lived, which is observable). Such an application is an interesting extension of CAR in the presence of demographic variables, such as age or gender, and may be of key interest to the area of latent variable modeling in general. However, if one is interested in the direct effect of age on the response Y , then the specific model of interest would include age as a predictor. In this case, the PCAR methodology developed here would be appropriate, as the application of CAR estimation would result in inconsistent estimates. In either case, the specific area of application and relevant body of scientific literature can provide guidance to the researcher in choosing the relevant model. For example, if a researcher is unsure of whether to enter age into the model as actually age or as “core biological wear”, then the statistician can provide guidance on focusing/clarifying the specific aims of the research hypotheses with the researcher and the available relevant literature. After the specific research hypothesis is defined, which may involve illiciting the relationship (between the predictor and response) of interest to the researcher, then a suitable/reasonable PCAR model can be chosen/entertained. Hence, once the hypothesis (relationship of interest) is identified, then the (underlying) model can be specified (containing predictor Z_s or Z_s^*) and the appropriate estimators (PCAR or CAR) can be applied accordingly.

Finally, we note the following regarding the implementation of the binning procedure, in the context of the data analysis. For the theory, we assumed that the support of U is the interval $[a, b]$. To be able to bin the data with respect to $U = \text{BMI}$ and compute the estimators, we take a and b to be the min/max of the data. For any given population under study, a reasonable range can be inferred to define the limits a and b . For adults, we can reasonably set the limits to 14 and 65 BMI, for instance. In our data, the observed min and max are 18.2 and 49.7 and the intervals/bins are between these observed limits. However, if one is applying the binning using the limits 14 and 65, for instance, our estimator weighted by L_j/n will assign zero weight to bins/intervals with no data (as it should), e.g. bins with U between 14 to 18.2 and 49.7 to 65 BMI. Thus, starting the binning at the min/max of the data are approximately equivalent to starting at $[a, b]$.

Acknowledgments

We are grateful to the reviewers for many detailed suggestions which substantially improved the paper as well as to the Editor for careful review. This work is supported by Grant Number UL1 RR024146 from the National Center for Research Resources (NCRR), a component of the National Institutes of Health (NIH) and NIH Roadmap for Medical Research, NIEHS Grant P01-ES011269-06, NIH Grants UL1RR024922 and RL1AG032115, and National Institute of Child Health and Human Development Grant HD036071.

Appendix

In this section we provide the additional lemmas and their proofs utilized earlier for the main results. We begin by formally defining the events under which the two main theorems of Section 4 are given. Summarizing the existence conditions for the PCAR estimators, define the events

$$\begin{aligned} \tilde{A}_n &= \left\{ \omega \in \Omega : \inf_j |\det(L_{nj}^{-1} \tilde{\chi}_{nj}^T \tilde{\chi}_{nj}')| > \zeta \text{ and } \min_j L_{nj} > p + q \right\}, \\ A_n &= \left\{ \omega \in \Omega : \inf_j |\det(L_{nj}^{-1} \chi_{nj}^T \chi_{nj}')| > \zeta \text{ and } \min_j L_{nj} > p + q \right\}, \end{aligned} \tag{19}$$

where $\zeta = \min\{\rho/2, [\inf_j \{\phi_1^2(U_{nj}^*), \dots, \phi_p^2(U_{nj}^*)\}]^p \rho/2\}$, ρ is as defined in (C5), $U_{nj}^* = L_{nj}^{-1} \sum_{k=1}^{L_{nj}} U'_{nj,k}$ is the average of the U 's in B_{nj} , and (Ω, \mathcal{F}, P) is the underlying probability space. The estimators in (8) and (9) are well defined on events \tilde{A}_n and A_n , respectively. The event E_n in Theorems 1 and 2 is defined to be the intersection of A_n and \tilde{A}_n , $E_n = A_n \cap \tilde{A}_n$. It is shown following Lemma 4 that $\text{pr}(E_n) \rightarrow 1$ as $n \rightarrow \infty$.

We next introduce some additional technical conditions that are needed for the proof of Lemma 4 given below:

- (C6) The functions $h_1(u) = \int x g_1(x, u) dx$ and $h_2(u) = \int x g_2(x, u) dx$ are uniformly Lipschitz, where $g_1(\cdot, \cdot)$ and $g_2(\cdot, \cdot)$ are the joint density functions of (χ, U) and $(\chi e, U)$, respectively.
- (C7) The error term satisfies $E|e^\tau| < \infty$ for $\tau > 4$.

Lemma 1. Under the technical conditions (C1)–(C7), on event A_n (19), the martingale differences W_{n0t} satisfy the conditions

$$(a) \sum_{t=1}^n E\{W_{n0t}^2 I(|W_{n0t}| > \varepsilon)\} \rightarrow 0 \text{ for all } \varepsilon > 0,$$

$$(b) \Delta_{n0}^2 = \sum_{t=1}^n W_{n0t}^2 \xrightarrow{P} \sigma_{\gamma_0}^2 \text{ for } \sigma_{\gamma_0}^2 > 0.$$

Proof. Let $W_{n0t} = w_{n0t}v_{n0t}$, where $w_{n0t} = 1/\sqrt{n}$, $v_{n0t} = \gamma_0\psi(U_{nt}) + \psi(U_{nt})e_{nt}\kappa_{1k(t)} - \gamma_0 \equiv \alpha_{1nt} + \alpha_{2nt}e_{nt}$, $\alpha_{1nt} = \gamma_0\psi(U_{nt}) - \gamma_0$, $\alpha_{2nt} = \psi(U_{nt})\kappa_{1k(t)}$, and $E(v_{n0t}) = 0$. Using (C1), (C3) and (C4), it holds on event A_n that $\sup_{1 \leq t \leq n} |\alpha_{1nt}| < c_1$ and $\sup_{1 \leq t \leq n} |\alpha_{2nt}| < c_2$ for some $c_1, c_2 > 0$. Thus, it holds for $\varepsilon > 0$ that

$$\begin{aligned} \sum_{t=1}^n E\{W_{n0t}^2 I(|W_{n0t}| > \varepsilon)\} &= \sum_{t=1}^n \int x^2 I(|x| > \varepsilon) dF_{w_{n0t}v_{n0t}}(x) \\ &= \sum_{t=1}^n \int x^2 I(|x| > \varepsilon/|w_{n0t}|) w_{n0t}^2 dF_{v_{n0t}}(x) = n^{-1} \sum_{t=1}^n \int x^2 I(|x| > \sqrt{n}\varepsilon) dF_{v_{n0t}}(x) \\ &\leq n^{-1} \sum_{t=1}^n \{E(v_{n0t}^4)\}^{1/2} \{\text{pr}(v_{n0t}^2 > n\varepsilon^2)\}^{1/2}. \end{aligned}$$

Now, $E(v_{n0t}^4)$ is bounded uniformly in n and t , since e_{nt} has finite fourth moment by (C7). Also note that $\text{pr}(v_{n0t}^2 > n\varepsilon^2) = \text{pr}((\alpha_{1nt} + \alpha_{2nt}e_{nt})^2 > n\varepsilon^2) \leq \text{pr}(\alpha_{1nt}^2 + \alpha_{2nt}^2 e_{nt}^2 + 2|\alpha_{1nt}\alpha_{2nt}e_{nt}| > n\varepsilon^2) \leq \text{pr}(c_1^2 + c_2^2 e_{nt}^2 + 2c_1 c_2 |e_{nt}| > n\varepsilon^2)$. Lemma 1(a) follows, since $\text{pr}(c_1^2 + c_2^2 e_{nt}^2 + 2c_1 c_2 |e_{nt}| > n\varepsilon^2) \rightarrow 0$ uniformly in n and t , e_{nt}^2 and $|e_{nt}|$ being i.i.d. with finite fourth moments.

Next, consider the term Δ_{n0}^2 given in Lemma 1(b). It is equal to

$$\begin{aligned} \Delta_{n0}^2 &= \gamma_0^2 \left\{ n^{-1} \sum_t \psi^2(U_{nt}) \right\} + \gamma_0^2 - 2\gamma_0^2 \left\{ n^{-1} \sum_t \psi(U_{nt}) \right\} + 2\gamma_0 n^{-1} \sum_t \psi^2(U_{nt})e_{nt}\kappa_{1k(t)} \\ &\quad - 2\gamma_0 n^{-1} \sum_t \psi(U_{nt})e_{nt}\kappa_{1k(t)} + n^{-1} \sum_{j=1}^m \sum_{k=1}^{L_{nj}} \{(\Theta_{nj})_{11}^{-1} e'_{nj} \psi(U'_{nj})\} \\ &\quad + (\Theta_{nj})_{12}^{-1} e'_{nj} \psi(U'_{nj}) X'_{n1jk} + \dots + (\Theta_{nj})_{1,p+q+1}^{-1} e'_{nj} \psi(U'_{nj}) Z'_{nqjk} \}^2 \equiv T_1 + \dots + T_6. \end{aligned}$$

Using Law of Large Numbers, it holds that $T_1 + T_2 + T_3 \xrightarrow{P} \gamma_0^2 \text{var}\{\psi(U)\}$. Since T_4 and T_5 have expected values zero and variances $O(n^{-1})$, they are both $O_p(n^{-1/2})$. By Lemma 4(a) and the Law of Large Numbers, term T_6 is equal to

$$\begin{aligned} &\sigma^2 E\{\psi^2(U)\} \{(\Gamma^{-1})_{11}^2 + (\Gamma^{-1})_{12}^2 X_1^2 + \dots + (\Gamma^{-1})_{1,p+q+1}^2 Z_q^2\} \\ &\quad + \{2(\Gamma^{-1})_{11}(\Gamma^{-1})_{12} X_1 + \dots + 2(\Gamma^{-1})_{11}(\Gamma^{-1})_{1,p+q+1} Z_q\} \\ &\quad + \{2(\Gamma^{-1})_{12}(\Gamma^{-1})_{13} X_1 X_2 + \dots + 2(\Gamma^{-1})_{12}(\Gamma^{-1})_{1,p+q+1} X_1 Z_q\} + \dots \\ &\quad + \{2(\Gamma^{-1})_{1,p+q}(\Gamma^{-1})_{1,p+q+1} Z_{q-1} Z_q\} + o_p(1) \\ &= \sigma^2 E\{\psi^2(U)\} \mathcal{M}_{11} + o_p(1), \end{aligned}$$

where Γ and \mathcal{M} are as defined prior to Theorem 1. Thus, $\Delta_{n0}^2 \xrightarrow{P} \gamma_0^2 \sigma_{\psi}^2 + \sigma^2 E\{\psi^2(U)\} \mathcal{M}_{11} \equiv \sigma_{\gamma_0}^2$ and Lemma 1(b) follows. \square

Lemma 2. Under the technical conditions (C1)–(C7), on event A_n (19), the martingale differences W_{nrt} satisfy the conditions

$$(a) \sum_{t=1}^n E\{W_{nrt}^2 I(|W_{nrt}| > \varepsilon)\} \rightarrow 0 \text{ for all } \varepsilon > 0,$$

$$(b) \Delta_{nr}^2 = \sum_{t=1}^n W_{nrt}^2 \xrightarrow{P} (a, b) \Sigma_r(a, b) \Gamma \text{ for } (a, b) \Sigma_r(a, b) \Gamma > 0.$$

Proof. Part (a) of Lemma 2 follows in a similar fashion as part (a) of Lemma 1. Therefore, we focus on the proof of part (b). The term Δ_{nr}^2 in Lemma 2(b) is equal to

$$\begin{aligned} \Delta_{nr}^2 &= a^2 \gamma_r^2 \left\{ n^{-1} \sum_t \psi^2(U_{nt}) X_{nrt}^2 \right\} + a^2 \gamma_r^2 \{E(X_r)\}^2 + b^2 \left\{ n^{-1} \sum_t \phi_r^2(U_{nt}) X_{nrt}^2 \right\} \\ &\quad + b^2 m_{r,1}^2 - 2a^2 \gamma_r^2 m_{r,1} \left\{ n^{-1} \sum_t \psi(U_{nt}) X_{nrt} \right\} + 2ab \gamma_r m_{r,1}^2 \\ &\quad + 2ab \gamma_r \left\{ n^{-1} \sum_t \psi(U_{nt}) \phi_r(U_{nt}) X_{nrt}^2 \right\} - 2b^2 m_{r,1} \left\{ n^{-1} \sum_t \phi_r(U_{nt}) X_{nrt} \right\} \\ &\quad - 2ab \gamma_r m_{r,1} \left\{ n^{-1} \sum_t \psi(U_{nt}) X_{nrt} \right\} - 2ab \gamma_r m_{r,1} \left\{ n^{-1} \sum_t \phi_r(U_{nt}) X_{nrt} \right\} \\ &\quad + 2a^2 \gamma_r n^{-1} \sum_t \psi^2(U_{nt}) e_{nt} \bar{X}'_{nrj(t)} X_{nrt} \kappa_{rk(t)} - 2a^2 \gamma_r E(X_r) n^{-1} \sum_t \psi(U_{nt}) e_{nt} \bar{X}'_{nrj(t)} \kappa_{rk(t)} \\ &\quad + 2ab n^{-1} \sum_t \psi(U_{nt}) \phi_r(U_{nt}) e_{nt} \bar{X}'_{nrj(t)} X_{nrt} \kappa_{rk(t)} - 2ab E(X_r) n^{-1} \sum_t \psi(U_{nt}) e_{nt} \bar{X}'_{nrj(t)} \kappa_{rk(t)} \\ &\quad + a^2 n^{-1} \sum_{j=1}^m \sum_{k=1}^{l_{nj}} \{(\Theta_{nj})_{r1}^{-1} \bar{X}'_{nrj} e'_{nj} \psi(U'_{nj}) + (\Theta_{nj})_{r2}^{-1} \bar{X}'_{nrj} e'_{nj} \psi(U'_{nj}) X'_{n1jk} + \dots \\ &\quad + (\Theta_{nj})_{r,p+q+1}^{-1} \bar{X}'_{nrj} e'_{nj} \psi(U'_{nj}) Z'_{nqjk}\}^2 \equiv T_1 + \dots + T_{15}. \end{aligned}$$

Using the Law of Large Numbers, it holds that $T_1 + \dots + T_{10} \xrightarrow{p} a^2 \gamma_r^2 [m_{r,1}^2 \sigma_\psi^2 + \text{var}(X_r) \lambda_\psi] + 2ab \gamma_r [\lambda_\psi \phi_r m_{r,2} - m_{r,1}^2] + b^2 \text{var}(\tilde{X}_r)$. Since T_{11}, T_{12}, T_{13} and T_{14} have expected values zero and variances $O(n^{-1})$, they are all $O_p(n^{-1/2})$. By Lemma 4(a) and the Law of Large Numbers, term T_{15} is equal to

$$\begin{aligned} &a^2 \sigma^2 m_{r,1}^2 E\{\psi^2(U) [(\Gamma^{-1})_{r+1,1}^2 + (\Gamma^{-1})_{r+1,2}^2 X_1^2 + \dots + (\Gamma^{-1})_{r+1,p+q+1}^2 Z_q^2 \\ &\quad + 2(\Gamma^{-1})_{r+1,1} (\Gamma^{-1})_{r+1,2} X_1 + \dots + 2(\Gamma^{-1})_{r+1,1} (\Gamma^{-1})_{r+1,p+q+1} Z_q] \\ &\quad + 2(\Gamma^{-1})_{r+1,2} (\Gamma^{-1})_{r+1,3} X_1 X_2 + \dots + 2(\Gamma^{-1})_{r+1,2} (\Gamma^{-1})_{r+1,p+q+1} X_1 Z_q] + \dots \\ &\quad + \{2(\Gamma^{-1})_{r+1,p+q} (\Gamma^{-1})_{r+1,p+q+1} Z_{q-1} Z_q\} + o_p(1) \\ &= a^2 \sigma^2 m_{r,1}^2 E\{\psi^2(U) \mathcal{M}_{r+1,r+1}\} + o_p(1). \end{aligned}$$

Thus

$$\Delta_{nr}^2 \xrightarrow{p} (a, b) \Sigma_r (a, b)^T = (a, b) \begin{bmatrix} \Sigma_{r11} & \Sigma_{r12} \\ \Sigma_{r12} & \Sigma_{r22} \end{bmatrix} (a, b)^T,$$

where $\Sigma_{r11} = \gamma_r^2 [m_{r,1}^2 \sigma_\psi^2 + \text{var}(X_r) \lambda_\psi] + \sigma^2 m_{r,1}^2 E\{\psi^2(U) \mathcal{M}_{r+1,r+1}\}$, $\Sigma_{r12} = \gamma_r [\lambda_\psi \phi_r m_{r,2} - m_{r,1}^2]$, and $\Sigma_{r22} = \text{var}(\tilde{X}_r)$. Hence Lemma 2(b) follows. \square

Lemma 3. Under the technical conditions (C1)–(C6), it holds on event E_n that,

$$\sup_j |(\tilde{\Theta}_{nj})^{-1} - \{\Phi_{nj} \square (\Theta_{nj})^{-1}\}| = O(m^{-1}) \mathbf{1}_{(p+q+1) \times (p+q+1)},$$

where

$$\Phi_{nj} = \begin{bmatrix} \Phi_{nj}^{11} & \Phi_{nj}^{21T} \\ \Phi_{nj}^{21} & \mathbf{1}_{q \times q} \end{bmatrix}, \quad \Phi_{nj}^{21} = \begin{bmatrix} 1 & 1/\phi_1(U_{nj}^*) & \dots & 1/\phi_p(U_{nj}^*) \\ \vdots & \vdots & & \vdots \\ 1 & 1/\phi_1(U_{nj}^*) & \dots & 1/\phi_p(U_{nj}^*) \end{bmatrix}_{p+1 \times q},$$

and

$$\Phi_{nj}^{11} = \begin{bmatrix} 1 & 1/\phi_1(U_{nj}^*) & \dots & 1/\phi_p(U_{nj}^*) \\ 1/\phi_1(U_{nj}^*) & 1/\phi_1^2(U_{nj}^*) & \dots & 1/(\phi_p(U_{nj}^*)\phi_1(U_{nj}^*)) \\ \vdots & & \ddots & \\ 1/\phi_p(U_{nj}^*) & 1/(\phi_p(U_{nj}^*)\phi_1(U_{nj}^*)) & \dots & 1/\phi_p^2(U_{nj}^*) \end{bmatrix}_{p+1 \times p+1}.$$

The proof follows from Lemma 3 of Şentürk and Müller (2006) by substituting 1 in place of $\phi_{p+1}, \dots, \phi_{p+q}$.

Lemma 4. Under the technical conditions (C1)–(C7), for a sequence r_n such that $r_n = O_p(\sqrt{(m \log n)/n})$, on event A_n

- (a) $\sup_j |(\Theta_{nj})^{-1} - \Gamma_{nj}^{-1}| = O_p(r_n) \mathbf{1}_{(p+q+1) \times (p+q+1)}$,
- (b) $\sup_j |L_{nj}^{-1} \mathbf{z}_{nj}^T e'_{nj}| = O_p(r_n) \mathbf{1}_{(p+q+1) \times 1}$,

where Γ_{nj} is assumed to be nonsingular by (C5), and $e'_{nj} = (e'_{nj1}, \dots, e'_{njL_{nj}})^T$.

The proof is similar to the proof of Lemma 4 given in Şentürk and Müller (2006). However a key difference is that the limiting term in part (a), Γ_{nj}^{-1} , contains expectations taken conditional on U . The conditioning on U does not disappear because of the dependence between Z_s and U in the case of PCAR.

Proof that $\text{pr}(E_n) \rightarrow 1$. The formula given in (32) in Şentürk and Müller (2006) can be extended to $\sup_j |\det(\Theta_{nj}) - \det(\Gamma_{nj})| = O_p(r_n)$, where r_n is as defined in Lemma 4. This implies, on event A_n , that $\text{pr}(\inf_j \det(\Theta_{nj}) > \zeta) \rightarrow 1$ as $n \rightarrow \infty$, where $\zeta = \min\{\rho/2, [\inf_j (\phi_1^2(U_{nj}^*), \dots, \phi_p^2(U_{nj}^*))]^\rho \rho/2\}$ and ρ is as defined in (C5). Similarly, it can be shown that $\text{pr}(\min_j L_{nj} \leq p + q) \rightarrow 0$ as $n \rightarrow \infty$, where $p + q$ denotes the number of predictors. Thus, $\text{pr}(A) \rightarrow 1$ as $n \rightarrow \infty$. Furthermore, Lemma 3 implies that

$$\sup_j |\det(\tilde{\Theta}_{nj}) - \phi_1^2(U_{nj}^*) \dots \phi_p^2(U_{nj}^*) \det(\Theta_{nj})| = O_p(m^{-1}).$$

This shows that $\text{pr}(\inf_j \det(\tilde{\Theta}_{nj}) > \zeta) \rightarrow 1$ as $n \rightarrow \infty$, which implies $\text{pr}(\tilde{A}_n) \rightarrow 1$ as $n \rightarrow \infty$. Thus $\text{pr}(E_n) \rightarrow 1$ as $n \rightarrow \infty$. □

References

Archie, J.P., 1981. Mathematical coupling of data: a common source of error. *Ann. Surgery* 193, 296–303.
 Berkson, J., 1950. Are there two regressions? *J. Amer. Statist. Assoc.* 45, 164–180.
 Carroll, R.J., Ruppert, D., Stefanski, L.A., Crainiceanu, C., 2006. *Measurement Error in Nonlinear Models: A Modern Perspective*. second ed.. Chapman & Hall, Boca Raton, FL.
 Cleveland, W.S., Grosse, E., Shyu, W.M., 1991. Local regression models. In: Chambers, J.M., Hastie, T.J. (Eds.), *Statistical Models in S*. Wadsworth & Brooks, Pacific Grove, pp. 309–376.
 Hanson, R.L., Ehm, M.G., Pettitt, D.J., Prochazka, M., Thompson, D.B., Timberlake, D., Foroud, T., Kobes, S., Baier, L., Burns, D.L., Almasy, L., Blangero, J., Garvey, W.T., Bennett, P.H., Knowler, W.C., 1998. An autosomal genomic scan for loci linked to type II diabetes mellitus and body-mass index in Pima Indians. *Amer. J. Hum. Genet.* 63, 1130–1138.
 Hastie, T., Tibshirani, R., 1993. Varying coefficient models. *J. Roy. Statist. Soc. B* 55, 757–796.
 Hwang, J.T., 1986. Multiplicative errors-in-variables models with applications to recent data released by the U.S. department of energy. *J. Amer. Statist. Assoc.* 81, 680–688.
 Iturria, S., Carroll, R.J., Firth, D., 1999. Polynomial regression and estimating functions in the presence of multiplicative measurement error. *J. Roy. Statist. Soc. B* 61, 547–561.
 The Hemo Study Group, Kaysen, G.A., Dubin, J.A., Müller, H.G., Mitch, W.E., Rosales, L.M., Levin, N.W., 2002. Relationship among inflammation nutrition and physiologic mechanisms establishing albumin levels in hemodialysis patients. *Kidney Int.* 61, 2240–2249.
 Knowler, W.C., Pettitt, D.J., Saad, M.F., Charles, M.A., Nelson, R.G., Howard, B.V., Bogardus, C., Bennett, P.H., 1991. Obesity in the Pima Indians: its magnitude and relationship with diabetes. *Amer. J. Clin. Nutr.* 53, 1543S–1551S.
 Lai, T.L., Robbins, H., Wei, C.Z., 1979. Strong consistency of least-squares estimates in multiple regression 2. *J. Multivariate Anal.* 9, 343–361.
 McLeish, D.L., 1974. Dependent central limit theorems and invariance principles. *Ann. Statist.* 2, 620–628.
 Pinter, J.D., Brown, W.E., Eliez, S., Schmitt, J.E., Capone, G.T., Reiss, A.L., 2001. Amygdala and hippocampal volumes in children with Down syndrome: a high-resolution MRI study. *Neurology* 56, 972–974.
 Şentürk, D., Müller, H.G., 2005. Covariate adjusted regression. *Biometrika* 92, 75–89.
 Şentürk, D., Müller, H.G., 2006. Inference for covariate adjusted regression via varying coefficient models. *Ann. Statist.* 34, 654–679.
 Smith, J.W., Everhart, J.E., Dickson, W.C., Knowler, W.C., Johannes, R.S., 1988. Using the ADAP learning algorithm to forecast the onset of diabetes mellitus. In: *Proceedings of the Symposium on Computer Applications and Medical Care*. pp. 261–265.
 Wu, C.O., Yu, K.F., 2002. Nonparametric varying coefficient models for the analysis of longitudinal data. *Int. Statist. Rev.* 70, 373–393.