

# Dimension Reduction Based on Conditional Multiple Index Density Function

Jun Zhang<sup>a</sup>, Baohua He<sup>b</sup>, Tao Lu<sup>c</sup> and Songqiao Wen<sup>d</sup>

<sup>a</sup>Shenzhen University

<sup>b</sup>Sun Yat-Sen University

<sup>c</sup>University of Nevada

<sup>d</sup>Shenzhen University

**Abstract.** In this paper, a dimension reduction method is proposed by using the first derivative of the conditional density function of response given predictors. To estimate the central subspace, we propose a direct methodology by taking expectation of the product of predictor and kernel function about response, which helps to capture the directions in the conditional density function. The consistency and asymptotic normality of the proposed estimation methodology are investigated. Furthermore, we conduct some simulations to evaluate the performance of our proposed method and compare with existing methods, and a real data set is analyzed for illustration.

## 1 Introduction

Accompanying the advancement of sciences and technologies in various fields such as biology, economics and finance, etc, scientific data has the tendency of growing in both size and complexity. Analysis of high-dimensional data calls for new statistical theories and methodologies. A natural way to analyze high-dimensional data is to first reduce the dimensionality of the original data without losing vital information. To reduce the problem of many covariates to one with a few covariates, sufficient dimension reduction [Cook \(1998\)](#) aims at finding low-dimensional linear combinations of the predictors without loss of information on  $Y|X$ . Suppose  $X$  is a random vector in  $\mathbb{R}^p$  and  $Y$  is a univariate random variable, we need seek a  $p \times d$  matrix  $B$  with  $d \leq p$  satisfying

$$Y \perp\!\!\!\perp X|B^T X, \quad (1.1)$$

where  $\tau$  is the transpose operator on a vector or a matrix and  $\perp\!\!\!\perp$  indicates independence throughout this paper. That is, given  $B^T X$ ,  $Y$  and  $X$  are independent. The space spanned by the column of  $B$ , which is denoted as  $\mathcal{S}(B)$ ,

---

*MSC 2010 subject classifications:* Primary 62G05; 62G08; 62G20

*Keywords and phrases.* Central subspace, Conditional density function, Dimensional reduction, Kernel function.

is defined as a dimension reduction subspace (Cook, 1994, 1998). In fact, we are not concerned about the specific form of  $B$  since any orthogonal transformation of  $B$  from right does not affect the conditional independent property of (1.1). If all the other dimension reduction space include  $\mathcal{S}(B)$  as their subspace, then  $\mathcal{S}(B)$  is a so-called central dimension reduction subspace (CS), and we denote the central dimension reduction subspace (CS) as  $\mathcal{S}_{Y|X}$ . Cook (1998) gives some mild conditions which guarantee the unique existence of  $\mathcal{S}_{Y|X}$ . Without notational confusion, we assume that  $\mathcal{S}_{Y|X}$  coincides with  $\mathcal{S}(B)$  throughout this paper, namely,  $\mathcal{S}_{Y|X} = \mathcal{S}(B)$ . The dimension of  $\mathcal{S}_{Y|X}$  is called the structural dimension and it is denoted as  $d$  in this paper.

To estimate the central subspace  $\mathcal{S}_{Y|X}$ , there are many useful methods available in the literature which make use of effective tools to reduce high-dimensional variables to equivalent ones comprising only some linear combinations of the original variables. For instance, ordinary least squares (Li and Duan, 1989, OLS), sliced inverse regression (Li, 1991, SIR), principal hessian directions (Li, 1992, pHd), sliced average variance estimation (Cook and Weisberg, 1991, SAVE), directional regression (Li and Wang, 2007, DR), minimum average variance estimation (Xia et al., 2002, MAVE), and outer product of gradient based on conditional density functions (Samarov, 1993, dOPG), density-MAVE (Xia, 2007, dMAVE), score dimension reduction (Wang and Zhu, 2013) etc. These dimension reduction methods are widely used for a large dataset and preserve  $\sqrt{n}$  consistency of the associated estimators, furthermore, the computational cost is not expensive. Once the central subspace is identified, subsequent analysis with the low-dimensional  $B^T X$  will help to construct another regression models or other statistical models based on  $B^T X$ . Specially, when the structural dimension is 1, model (1.1) reduces to the popular single-index models (Liang et al., 2010; Li et al., 2014; Chen et al., 2010; Peng and Huang, 2011; Cui et al., 2009). Moreover, if the estimation structural dimension is 1, 2 or 3, the graphical visualization will gain comprehensive insights about the data, see Cook (1998) with a deep discussion about the graphical methodology. So, the dimension reduction method is widely applied in practice. See, Zhu et al. (2016, ?); Lansangan and Barrios (2017); Luo et al. (2017); Yoshida (2017); Deng and Wang (2017); Sheng and Yin (2016); Zhou and Zhu (2016).

The central subspace satisfying model (1.1) is equivalent to the conditional density function of  $Y|X$  being the same as that of  $Y|B^T X$  for all possible value of  $X$  and  $Y$  if the conditional density function of  $Y$  given  $X$  exists, i.e.,

$$f_{Y|X}(y|x) = f_{Y|B^T X}(y|B^T x). \quad (1.2)$$

Model (1.2) indicates that all the directions can be captured in the conditional density function  $f_{Y|X}(y|x)$ . Then, an estimator of the conditional density function  $f_{Y|X}(y|x)$  is needed to be proposed in the primary step to find the central subspace. Xia (2007) used the “double-kernel” local linear smoothing method (Fan et al., 1996) to estimate the conditional density function  $f_{Y|X}(y|x)$ . They used the fact that the conditional density function  $f_{Y|X}(y|x)$  is asymptotically equivalent to the conditional regression mean function, and the directions defined in model (1.1) will be all captured in this conditional regression mean function.

In this paper, we introduce a simple methodology for dimension reduction that based on the linear condition on  $X$  (Li, 1991) and the existence of the conditional density function of  $Y$  given  $X$ . The linear condition on  $X$  is widely used in the dimension reduction literature, and the existence of conditional density function of  $Y$  given  $X$  is also a mild assumption. The methodology proposed in this paper is to take expectation of the product of  $X$  and the kernel function of  $Y$ . Specifically, we use the idea of “double-kernel” local linear smoothing method (Fan et al., 1996) to estimate the conditional density function  $f_{Y|X}(y|x)$  in the first step. Due to the asymptotically equivalent conditional density function  $f_{Y|X}(y|x)$  and the conditional regression mean function (Fan et al., 1996; Xia, 2007). Together with model assumption (1.2), the conditional regression mean function will be used to capture all the direction of central subspace. Next, the proposed method in this paper is easy to implement by using a spectral decomposition on a kernel matrix. As the bandwidth converges to zero, the kernel matrix will eventually recover the central subspace in the population level. In this paper, we also provide the asymptotic properties of the estimators of the estimated kernel matrix and associated eigenvalues and eigenvectors.

The reminder of the paper is organized as follows. In Section 2, we introduce the rationale of the dimension reduction method at population level, and illustrate its theoretical result. At the sample level, we introduce a direct estimation approach to estimating  $S_{Y|X}$ , and establish the asymptotic properties of the resultant estimators. We demonstrate the methodologies through simulations and an analysis of a real data in Section 3. An analysis of a real data is presented in Section 4. We conclude this paper with a brief discussion in Section 5. All proofs are given in the Appendix.

## 2 Methodology Development

### 2.1 The Population level

We introduce the proposed dimension reduction method in the population level. When the conditional density function of  $Y$  given  $X$  exists, model (1) is equivalent to the conditional density function of  $Y|X$  being the same as that of  $Y|B^T X$  for all possible values of  $(x, y)$  over the support of  $(X, Y)$ , ie,  $f_{Y|X}(y|x) = f_{Y|B^T X}(y|B^T x)$ . It implies that

$$\partial f_{Y|X}(y|x)/\partial x = B[\partial f_{Y|B^T X}(y|B^T x)/\partial(B^T x)]. \quad (2.1)$$

In other words, all the basis directions in  $\mathbf{S}_{Y|X}$  can be captured by the first derivation of the conditional density function. Thus, to recover  $\mathbf{S}_{Y|X}$ , we need to estimate  $\partial f_{Y|X}(y|x)/\partial x$  or  $\mathbf{E}[\partial f_{Y|X}(y|X)/\partial X]$ . However, we could not estimate this derivative directly by nonparametric smoothing when the dimensions of  $X$  are high. To avoid using nonparametric smoothing method in high-dimensional case, if  $X$  is further assumed temporarily to be normally distributed with mean zero and identical covariance matrix, a direct application of [Stein \(1981\)](#)'s Lemma 4 and combining with equations (1.2) and (2.1) entail that

$$\begin{aligned} \mathbf{E}[X f_{Y|X}(y|X)] &= \mathbf{E}[\partial f_{Y|X}(y|X)/\partial X] \\ &= B\mathbf{E}[\partial f_{Y|B^T X}(y|B^T X)/\partial(B^T X)]. \end{aligned} \quad (2.2)$$

It is seen that  $\mathbf{E}[X f_{Y|X}(y|X)]$  can be used to seek the directions of  $\mathbf{S}_{Y|X}$ . However, in the argument (2.2), the density  $f_{Y|X}(y|x)$  is still unknown and needs to be estimated. [Xia \(2007\)](#) used the idea of "double-kernel" smoothing method studied in [Fan et al. \(1996\)](#) to construct an estimator of conditional density function  $f_{Y|X}(y|x)$ . We denote  $K(\cdot)$  as a symmetric density function and  $h$  is the bandwidth,  $h > 0$  and  $K_h(x) = h^{-1}K(x/h)$ . If  $h \rightarrow 0$  and  $n \rightarrow \infty$ , we have

$$\begin{aligned} \mathbf{E}(K_h(Y - y)|X = x) &= \mathbf{E}(K_h(Y - y)|B^T X = B^T x) \\ &\rightarrow f_{Y|B^T X}(y|B^T x) \end{aligned} \quad (2.3)$$

Combining with (2.2) and (2.3), if  $X$  follows a normal distribution with identical covariance matrix, as  $h \rightarrow 0$  and  $n \rightarrow \infty$ , we have

$$\begin{aligned} \mathbf{E}(X K_h(Y - y)) &= \mathbf{E}[X \mathbf{E}(K_h(Y - y)|X)] \\ &\rightarrow \mathbf{E}(X f_{Y|B^T X}(y|B^T x)) = B[\partial f_{Y|B^T X}(y|B^T x)/\partial(B^T x)]. \end{aligned} \quad (2.4)$$

Equation (2.4) indicates all the directions can be captured by the  $\mathbf{E}[XK_h(Y - y)]$  under normality assumption when  $h \rightarrow 0$  as  $n \rightarrow \infty$ . In fact, we avoid to estimate  $f_{Y|X}(y|x)$ ,  $\partial f_{Y|X}(y|x)/\partial x$  and  $\mathbf{E}[\partial f_{Y|X}(y|X)/\partial X]$  by using non-parametric smoothing methods, while a simple moment-based estimator of  $\mathbf{E}(XK_h(Y - y))$  helps to infer the central subspace  $\mathbf{S}_{Y|X}$ . However, the normality assumption for  $X$  is relatively restrictive. Without loss of generality, we assume  $\mathbf{E}X = 0$  and relax the normality assumption of  $X$  to the widely used linearity condition (Li, 1991) in the following

$$\mathbf{E}[X|B^\tau X] = P_B^\tau(\Sigma_X)X, \quad (2.5)$$

where  $P_B(\Sigma_X) = B(B^\tau \Sigma_X B)^{-1}B^\tau \Sigma_X$ ,  $\Sigma_X = \text{Var}(X)$ . As a consequence, under the linearity condition of  $X$ ,  $\mathbf{E}[XK_h(Y - y)]$  still could be used to seek  $\mathbf{S}_{Y|X}$  when  $h \rightarrow 0$  and  $n \rightarrow \infty$ . Define

$$\mathbf{D}_h(y) = \mathbf{E}[XK_h(Y - y)], \quad (2.6)$$

$$\mathbf{D}_B(y) = P_B^\tau \mathbf{E}[Xf_{Y|B^\tau X}(y|B^\tau X)], \quad (2.7)$$

First of all, we list some conditions for our asymptotic results.

- (A1) The conditional density functions  $f_{Y|X}(y|x)$  and  $f_{Y|B^\tau X}(y|B^\tau x)$  have continuous and bounded second order derivatives with respect to  $x$ . Furthermore,  $\mathbf{D}_B(y) < \infty$  for all possible values of  $y$  on the support of  $Y$ , and  $\mathbf{D}_B < \infty$ .
- (A2)  $\mathbf{E}X^\tau X < \infty$ ,  $\mathbf{D}_h(y) < \infty$ ,  $\mathbf{D}_h < \infty$  for all possible values of  $y$  on the support of  $Y$  as  $h \rightarrow 0$ .
- (A3) The kernel function  $K(\cdot)$  is symmetric about 0 and has a compact support. Moreover,

$$\int K(u)du = 1, \quad \int uK(u)du = 0, \quad \int u^2 K(u) < \infty.$$

- (A4) The density function  $g(\cdot)$  of  $Y$  has bounded second order derivatives. Moreover,  $\mathbf{E}(X^\tau X)^4 < \infty$ ,  $\mathbf{E}[\mathbf{D}_B(Y)g(Y)] < \infty$ .

**Theorem 2.1.** *If  $\mathbf{E}X = 0$  and  $\Sigma_X$  is a positive definite matrix, moreover,  $X$  satisfies the linearity condition (2.5). Under the conditions (A1)-(A4), as  $h \rightarrow 0$ ,  $n \rightarrow \infty$ , we have*

- (1)  $\Sigma_X^{-1} \mathbf{D}_h(y) \rightarrow \Sigma_X^{-1} \mathbf{D}_B(y) \subseteq \mathbf{S}_{Y|X}$ .
- (2) Let  $\mathbf{D}_h = \mathbf{E}[\mathbf{D}_h(\tilde{Y})\mathbf{D}_h^\tau(\tilde{Y})]$ ,  $\mathbf{D}_B = \mathbf{E}[\mathbf{D}_B(\tilde{Y})\mathbf{D}_B^\tau(\tilde{Y})]$ , where  $\tilde{Y}$  is an independent copy of  $Y$ . The kernel matrix  $\mathbf{V}_h$ , defined as  $\mathbf{V}_h = \Sigma_X^{-1} \mathbf{D}_h \Sigma_X^{-1}$ , satisfies  $\mathbf{V}_h \rightarrow \mathbf{V}_B = \Sigma_X^{-1} \mathbf{D}_B \Sigma_X^{-1} \subseteq \mathbf{S}_{Y|X}$ .

Theorem 2.1 indicates that  $\mathbf{V}_B$  estimates the central subspace  $\mathcal{S}_{Y|X}$  in the population level. If we apply a spectral decomposition on the kernel matrix  $\mathbf{V}_B$  to get  $\{\beta_1, \dots, \beta_k\}$ , which are the eigenvectors of kernel matrix  $\mathbf{V}_B$  corresponding to its largest  $k$  nonzero eigenvalues, then the space  $S(\beta_1, \dots, \beta_k)$  spanned by  $\{\beta_1, \dots, \beta_k\}$  will recover  $\mathcal{S}_{Y|X}$ . Theorem 2.1 entails a direct estimation approach by using the kernel matrix  $\mathbf{V}_h$  to estimate  $\mathcal{S}_{Y|X}$  when  $h \rightarrow 0$  as  $n \rightarrow \infty$ . In next subsection, we introduce the estimation procedures of  $\mathbf{D}_h(y)$  and  $\mathbf{V}_h$  and present the asymptotic properties of these estimators.

## 2.2 Estimation Procedures and Asymptotic Results

Suppose that  $\{(X_i, Y_i)\}_{i=1}^n$  are i.i.d sample from the model (1.2). Denote  $K_h(\cdot) = h^{-1}K(\cdot/h)$ . The estimators of  $\mathbf{D}_h(y)$  and  $\Sigma_X$  are defined as

$$\hat{\mathbf{D}}_h(y) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}) K_h(Y_i - y), \quad (2.8)$$

$$\hat{\Sigma}_X = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^\tau, \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i. \quad (2.9)$$

Using (2.8) and (2.9), the moment estimators of  $\mathbf{D}_h$  and  $\mathbf{V}_h$  in Theorem 2.1 are proposed as

$$\hat{\mathbf{D}}_h = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{D}}_h(Y_i) \hat{\mathbf{D}}_h^\tau(Y_i), \quad (2.10)$$

$$\hat{\mathbf{V}}_h = \hat{\Sigma}_X^{-1} \hat{\mathbf{D}}_h \hat{\Sigma}_X^{-1}. \quad (2.11)$$

**Theorem 2.2.** *Under the conditions of Theorem 2.1, if  $h \rightarrow 0$  and  $nh \rightarrow \infty$  as  $n \rightarrow \infty$ , we have*

$$\hat{\mathbf{V}}_h \xrightarrow{P} \mathbf{V}_B.$$

The consistency of  $\hat{\mathbf{V}}_h$  to  $\mathbf{V}_B$  enables us to estimate  $\mathcal{S}_{Y|X}$  by using the first  $k$  eigenvectors of the estimated kernel matrix  $\hat{\mathbf{V}}_h$  associated with its  $k$  largest nonzero eigenvalues.

Let  $\text{Vech}(A) = (a_{11}, \dots, a_{p1}, a_{22}, \dots, a_{p2}, a_{33}, \dots, a_{pp})^\tau$  be a  $p(p+1)/2$  dimension vector for any symmetric  $p \times p$  dimensional matrix  $A = (a_{ij})_{p \times p}$ . In the following, we use notation  $A^{\otimes 2} = AA^\tau$  for any matrix or vector  $A$ .

Define

$$\begin{aligned} \mathbf{T}_h(x, y) &= (x^{\otimes 2} - \Sigma_X) \Sigma_X^{-1} \mathbf{D}_B + \mathbf{D}_B \Sigma_X^{-1} (x^{\otimes 2} - \Sigma_X) \quad (2.12) \\ &+ \mathbf{E}\{\mathbf{D}_B(y)g(y)\} x^\tau + x \mathbf{E}\{\mathbf{D}_B^\tau(y)g(y)\} \\ &+ \mathbf{E}(\hat{\mathbf{D}}_h^0 | (X_1, Y_1) = (x, y)) - \mathbf{E} \hat{\mathbf{D}}_h^0, \end{aligned}$$

where  $\hat{D}_h^0$  is defined as (A.2) in Appendix. Next, we present the asymptotic distribution of the estimated kernel matrix  $\hat{V}_h$ .

**Theorem 2.3.** *Under the conditions of Theorem 2.1,  $nh^4 \rightarrow 0$  and  $nh^2 \rightarrow \infty$  as  $n \rightarrow \infty$ , moreover,  $a^\tau \text{Cov}(\text{Vech}(\Sigma_X^{-1} \mathbf{T}_h(X, Y) \Sigma_X^{-1})) a \rightarrow \sigma_a^2 > 0$  for any  $a \in \mathbb{R}^{\frac{p(p+1)}{2}}$  and  $a \neq 0$ , we have*

$$\sqrt{n}(\hat{V}_h - \mathbf{V}_B) \xrightarrow{L} \mathcal{H}. \quad (2.13)$$

where  $a^\tau \text{Vech}(\mathcal{H})$  follows a normal distribution  $N(0, \sigma_a^2)$  for any  $a \in \mathbb{R}^{\frac{p(p+1)}{2}}$  and  $a \neq 0$ .

Let  $\lambda(A)$  stands for the vector of ordered eigenvalues of  $A$ , i.e, denoted as  $\lambda(A) = (\lambda_1(A), \dots, \lambda_p(A))^\tau$  satisfying  $\lambda_1(A) \geq \dots \geq \lambda_p(A)$ . Denote  $\lambda_1(\mathbf{V}_B) > \dots > \lambda_l(\mathbf{V}_B)$  are the distinct eigenvalues of  $\mathbf{V}_B$  with the multiplicity of  $\lambda_i(\mathbf{V}_B)$  being  $m_i, i = 1, \dots, l$ , and  $m_1 + m_2 + \dots + m_l = p$ .

The orthogonal matrices  $Q$  makes

$$Q^\tau \mathbf{V}_B Q = \begin{bmatrix} \lambda_1(\mathbf{V}_B) I_{m_1} & 0 & \dots & 0 \\ 0 & \lambda_2(\mathbf{V}_B) I_{m_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_l(\mathbf{V}_B) I_{m_l} \end{bmatrix},$$

where  $\tilde{\mathcal{H}} = (\tilde{\mathcal{H}}_{i,j})$  denotes the partitioning of  $Q^\tau \mathcal{H} Q$  in blocks of order  $m_i \times m_j$ . Theorem 2.3 entails to derive the asymptotic distribution of the eigenvalues of the estimated kernel matrix  $\hat{V}_h$ . By application of Theorem 3.1 and Theorem 3.2 in Eaton and Tyler (1991), we have the following asymptotic result.

**Theorem 2.4.** *Under the conditions of Theorem 2.3,  $\lambda_1(\mathbf{V}_B) > \lambda_2(\mathbf{V}_B) > \dots > \lambda_l(\mathbf{V}_B)$  is the distinct eigenvalues of  $\mathbf{V}_B$  with the multiplicity of  $\lambda_i(\mathbf{V}_B)$  being  $m_i, i = 1, \dots, l$  and  $m_1 + m_2 + \dots + m_l = p$ , we have*

$$\sqrt{n}(\lambda(\hat{V}_h) - \lambda(\mathbf{V}_B)) \xrightarrow{L} \begin{bmatrix} \lambda(\tilde{\mathcal{H}}_{11}) \\ \lambda(\tilde{\mathcal{H}}_{22}) \\ \vdots \\ \lambda(\tilde{\mathcal{H}}_{ll}) \end{bmatrix}.$$

From Theorem 2.4, if the first  $k$ th of the kernel matrix  $\mathbf{V}_B$  are nonzero, the first  $k$ th nonzero eigenvalues  $\{\hat{b}_1, \dots, \hat{b}_k\}$  of the estimated kernel matrix  $\hat{V}_h$  are root- $n$  consistent, and the first  $d$ th eigenvectors with  $d \leq k$  corresponding to nonzero eigenvalues helps to infer the central subspace  $\mathcal{S}_{Y|X}$ .

**Theorem 2.5.** *Let  $e$  be any unit length vector which is orthogonal to  $\mathbf{S}_{Y|X}$ , and suppose that  $e^\tau \text{Cov}(\Sigma_X^{-1} \mathbf{T}_h(X, Y) \Sigma_X^{-1} b_j) e \rightarrow e^\tau W_j e > 0$  as  $h \rightarrow 0$ ,  $j = 1, \dots, k$ . Under the conditions of Theorem 2.3, we have*

$$\sqrt{n} e^\tau \hat{b}_j \xrightarrow{L} N(0, e^\tau W_j e).$$

### 3 Simulations Studies

In this section, we conduct some simulations to evaluate the performance of our proposed method. To evaluate the estimation accuracy, we use a measure between two subspace of  $\mathbb{R}^p$ . Let  $B$  is a  $p \times k$  matrix spanning  $\mathbf{S}_{Y|X}$ , and  $\hat{B}$  is a  $p \times k$  matrix to estimate  $B$ . The measure between  $\mathbf{S}_{Y|X}$  and its estimator  $\hat{\mathbf{S}}_{Y|X}$  is defined as

$$\text{dist}(\mathbf{S}_{Y|X}, \hat{\mathbf{S}}_{Y|X}) = \|P_B - P_{\hat{B}}\|,$$

where  $P_B$  and  $P_{\hat{B}}$  is the projection operator in the standard inner product of  $B$  and  $\hat{B}$ , and  $\|\cdot\|$  is the Euclidean matrix norm. The smaller value of  $\text{dist}(\mathbf{S}_{Y|X}, \hat{\mathbf{S}}_{Y|X})$ , the better performance of  $\hat{\mathbf{S}}_{Y|X}$ . A more detailed discussion of the distance measure can be referred to Li et al. (2005).

We compare our proposed method with some useful methods, SIR (Li, 1991), SAVE (Cook and Weisberg, 1991), pHd (Li, 1992), DR (Li and Wang, 2007), rMAVE (Xia et al., 2002), dMAVE (Xia, 2007). For SIR, SAVE, and DR, we consider two cases with the slice number  $H = 5$  and  $H = 10$ . For our method, we use the kernel function  $K(t) = (15/16)(1-t^2)^2 I(t^2 < 1)$  and bandwidth  $h = n^{-1/3}$  to satisfy the conditions in Theorems. The dimension of predictor  $X$  is chosen as  $p = 5, 10, 20$  for Example 1 and Example 2.

**Example 1** The following three models are used:

$$Y = \frac{1}{2} \exp(X_1 - 1) \varepsilon; \quad (3.1)$$

$$Y = \sin \left\{ \frac{1}{\sqrt{2}} X_1 - \frac{1}{\sqrt{2}} X_2 + \varepsilon \right\}; \quad (3.2)$$

$$Y = \cos \left\{ \exp\left(\frac{1}{\sqrt{2}} X_1 - \frac{1}{\sqrt{2}} X_2 + \varepsilon\right) \right\}; \quad (3.3)$$

In these models, the structure dimension of  $\mathbf{S}_{Y|X}$  is one. The variables  $X_{ij}$  independently follows  $t(8)$  ( $t$ -distribution with 8 freedom-degree) for  $i = 1, \dots, n$ ,  $j = 1, \dots, p$ , and the error  $\varepsilon \sim N(0, 1)$  satisfying  $\varepsilon \perp\!\!\!\perp X$ . Among these dimension reduction methods, we present estimation error for each model with 300 replications. In Table 1, we present the mean of distance

measure  $\text{dist}(\mathbf{S}_{Y|X}, \hat{\mathbf{S}}_{Y|X}) = \|P_B - P_{\hat{B}}\|$  between these methods and the true central subspace and also the standard error.

From Table 1, we see that our method outperforms SAVE, pHd and rMAVE in this three models. As the numbers of predictor  $p$  increases, our method still could detect the underlying true dimension reduction subspace and has better performance than SAVE and rMAVR. Both SIR with slice number  $H = 5$ ,  $H = 10$  and dMAVE have comparable performance with our method, and the latter has a slightly better performance than DR in the three models.

**Example 2** The following two models are used:

$$Y = \text{sign}(\beta_1^T X + \varepsilon_1) \times \log(|6 + \beta_2^T X + \varepsilon_2|), \quad (3.4)$$

$$Y = \text{sign}(2\beta_1^T X + \varepsilon_1) \times \log(|4 + 2\beta_2^T X + \varepsilon_2|), \quad (3.5)$$

where  $\text{sign}(\cdot)$  is the sign function. In this example, it is seen that that  $\mathbf{S}_{Y|X}$  is spanned by  $(\beta_1, \beta_2)$ .

For model (3.4), the first five elements of  $\beta_1$  are  $(1, 1, 1, 1, 0)/2$  and the rest are all zeros. The first five elements of  $\beta_2$  are  $(-1, 2, -2, 1, 1)/\sqrt{11}$  and the rest are all zeros. The distribution of  $X$  is  $N(0, I_p)$ , where  $I_p$  is a  $p \times p$  identical matrix, and unobservable noise  $\varepsilon_1$  and  $\varepsilon_2$  follow from  $N(0, 0.5^2)$ . Moreover,  $X$  is independent with  $(\varepsilon_1, \varepsilon_2)$ , and  $\varepsilon_1$  is independent with  $\varepsilon_2$ .

For model (3.5), the first five elements of  $\beta_1$  are  $(1, 1, 1, 1, 0)/2$  and the rest are all zeros. The first five elements of  $\beta_2$  are  $(1, -1, 1, -1, 0)/2$ , and the rest are all zeros too. The distribution of  $X$  is also  $N(0, I_p)$  and independent with  $(\varepsilon_1, \varepsilon_2)$ , and unobservable noise  $\varepsilon_1$  and  $\varepsilon_2$  are independent and both are  $N(0, 1)$ . The models investigated here are similar to [Chen and Li \(1998\)](#) and [Xia \(2007\)](#). The model (3.5) is the same as the model in Example 4.1 in [Xia \(2007\)](#). We consider the models (3.4)-(3.5) with different dimension reduction methods with sample size  $n = 300$  in this example. The slice numbers of SIR, SAVE and DR are also chosen as  $H = 5$  and  $H = 10$ . The simulation results of mean and standard error with 300 replications are reported in Table 2.

From Table 2, the mean of  $\text{dist}(\mathbf{S}_{Y|X}, \hat{\mathbf{S}}_{Y|X})$  by our method is better than SAVE, pHd, rMAVE. As the numbers of predictor  $p$  increase, our method still has comparable performance with SIR and DR with the slice number  $H = 5$ ,  $H = 10$ . For model (3.5), when we change different parameter values of  $(\beta_1, \beta_2)$  and reduce the variance of  $(\varepsilon_1, \varepsilon_2)$ , the performance of our method is better than dMAVE. For model (3.4), dMAVE could not detect the underlying dimension reduction subspace when the number of predictors is relatively small  $p = 5$ , while our method has stable performance in both cases even when the number of predictors  $p$  increases to 20.



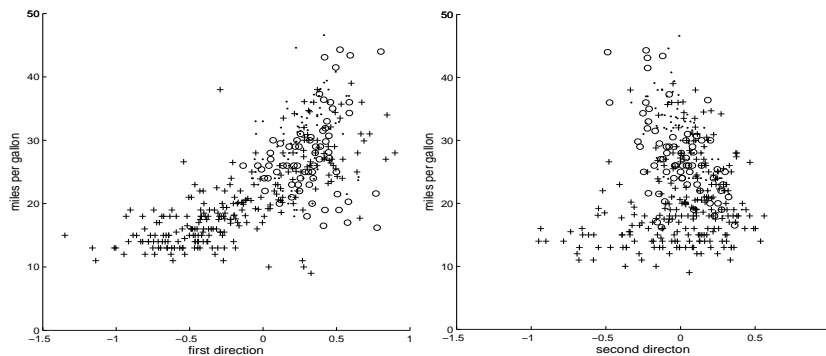
Technology. There are 406 observations with eight variables:  $Y$ -miles per gallon,  $X_1$ -number of cylinders,  $X_2$ -engine displacement (cu. inches),  $X_3$ -horsepower,  $X_4$ -vehicle weight (lbs.),  $X_5$ -time to accelerate from 0 to 60 mph (sec.),  $X_6$ -model year (modulo 100), and origin of car (1=American, 2=European, 3=Japanese). We transform the variable origin of car to the pairwise variable  $(X_7, X_8)$ , where  $(X_7, X_8) = (1, 0), (0, 1), (0, 0)$  corresponding to American cars, European cars and Japanese cars. Before applying our dimension reduction method, we standardize all the covariates separately. The bandwidth is taken as three cases,  $h = n^{-1/3}$ ,  $h = n^{-1/2}$  and  $h = n^{-4/5}$ .

Denote the estimated central subspace by our method as  $\hat{\mathcal{S}}_{Y|X}$  in this dataset. To determine the number of dimension  $d$ , we use bootstrap to select the dimension based on [Ye and Weiss \(2003\)](#), and BIC type criterion proposed by [Zhu et al. \(2006\)](#). We re-sample  $n = 300$  dataset without replacement and the standardized covariates, and we obtained a number of 500 bootstrap estimates  $\hat{\mathcal{S}}_{Y|X}^b$ ,  $b = 1, \dots, 500$  based on our method. To measure the distance between the data estimator  $\hat{\mathcal{S}}_{Y|X}$  and bootstrap estimator  $\hat{\mathcal{S}}_{Y|X}^b$ , we adopt the vector correction coefficient  $q$  ([Hotelling, 1936](#)) and use  $\arccos(q)$  as a measure, see more details in [Ye and Weiss \(2003\)](#). The mean of bootstrap distance measure  $\arccos(\hat{q}^b)$  is reported in Table 3.

**Table 3** Bootstrap mean of  $\arccos(\hat{q}^b)$  between  $\hat{\mathcal{S}}_{Y|X}^b$  and  $\hat{\mathcal{S}}_{Y|X}$ .

	$d = 1$	$d = 2$	$d = 3$	$d = 4$	$d = 5$	$d = 6$	$d = 7$
$h = n^{-1/3}$	21.9695	14.1481	20.3936	14.9362	15.5745	20.5963	30.1306
$h = n^{-1/2}$	22.8674	14.2178	21.6424	14.9306	15.4797	20.3341	28.8257
$h = n^{-4/5}$	24.8663	14.1931	21.8925	14.7187	15.6498	21.5143	27.2937

From Table 3, the case of  $d = 2$  has the smallest value and then we suggest the dimension in this dataset is 2. The BIC type criterion proposed by [Zhu et al. \(2006\)](#) also suggests that the dimension is 2. Now we use our method to this dataset and use the bandwidth  $h = n^{-1/2}$  here, the two directions are  $\hat{\beta}_1 = (-0.11, 0.73, -0.59, 0.14, -0.15, -0.20, -0.09, -0.07)^\tau$  and  $\hat{\beta}_2 = (0.01, -0.55, -0.43, 0.68, -0.13, -0.15, 0.05, -0.08)^\tau$ . The first two estimated direction indicate that the three predictors  $X_2$ -engine displacement (cu. inches),  $X_3$ -horsepower and  $X_4$ -vehicle weight (lbs.) are of dominating effects. The plots of  $Y$  against  $\hat{\beta}_1^\tau X$  and  $\hat{\beta}_2^\tau X$  are shown in Figure 1. Figure 1 presents the scatter plot with the response  $Y$ -miles per gallon and the first two estimated directions. In the left panel plot, the response  $Y$  has a linear trend in the first direction for American and Japanese cars. In the right panel plot, the response  $Y$  has a slight linear trend in the second direction



**Figure 1** Cars data. The left panel, the scatter plots of  $Y$  against the first direction. The right panel, the scatter plots of  $Y$  against the second direction. The origins of cars are denoted by “+” for American cars, “o” for European cars and “x” for Japanese cars.

for European and Japanese cars and has no trend for American cars.

## 5 Discussion and Further research

This paper gives a topic of the estimation procedure for dimension reduction by using the conditional density function. We present the asymptotic results of the proposed estimators and investigate the numerical performance. We can study the proposed methods in this paper to consider the estimation in the divergent parameters (Wu and Li, 2014; Zhu et al., 2006) and variable selection problems (Chen et al., 2010). One can also use the proposed methods in this paper to consider the measurement errors data (Li and Yin, 2007), longitudinal data (Bi and Qu, 2015; Li and Yin, 2009), missing data (Ding and Wang, 2011; Guo et al., 2014) in a future work. The research for this topic is ongoing.

## Appendix A: Appendix section

In this section, we give the proofs of our main results. Without loss of generality, in the following we assume  $EX = 0$ .

### A.1 Proof of Theorem 2.1

**Proof.** Invoking the assumption of that the conditional density function of  $Y|X$  being the same as that of  $Y|B^T X$ , i.e.,  $f_{Y|X}(y|x) = f_{Y|B^T X}(y|B^T X)$ ,

when the linear condition (2.5) of  $X$  holds, we have

$$\begin{aligned} \mathbf{E}[X f_{Y|X}(y|X)] &= \mathbf{E}[X f_{Y|B^\tau X}(y|B^\tau X)] \\ &= \mathbf{E}\{\mathbf{E}[X f_{Y|B^\tau X}(y|B^\tau X)|B^\tau X]\} = \mathbf{E}\{f_{Y|B^\tau X}(y|B^\tau X)\mathbf{E}[X|B^\tau X]\} \\ &= \mathbf{E}[f_{Y|B^\tau X}(y|B^\tau X)P_B^\tau X] = P_B^\tau \mathbf{E}[X f_{Y|B^\tau X}(y|B^\tau X)], \end{aligned} \quad (\text{A.1})$$

where  $P_B^\tau(\Sigma_X) = B(B^\tau \Sigma_X B)^{-1} B^\tau \Sigma_X$ . From the expression (A.1), we see that  $\Sigma_X^{-1} \mathbf{E}[X f_{Y|X}(y|X)] \subseteq \mathbf{S}_{Y|X}$ . Fan et al. (1996) shows that  $\mathbf{E}(K_h(Y - y)|X = x) \rightarrow f_{Y|X}(y|x)$  as  $h \rightarrow 0, n \rightarrow \infty$ . Using this results, we have

$$\begin{aligned} \mathbf{D}_h(y) &= \mathbf{E}[X K_h(Y - y)] = \mathbf{E}\{X \mathbf{E}[K_h(Y - y)|X]\} \\ &\rightarrow \mathbf{E}\{X f_{Y|X}(y|X)\} = \mathbf{E}[X f_{Y|B^\tau X}(y|B^\tau X)] \\ &= P_B^\tau \mathbf{E}[X f_{Y|B^\tau X}(y|B^\tau X)] = \mathbf{D}_B(y), \end{aligned}$$

i.e,  $\Sigma_X^{-1} \mathbf{D}_h(y) \rightarrow \Sigma_X^{-1} \mathbf{D}_B(y) \subseteq \mathbf{S}_{Y|X}$  as  $h \rightarrow 0, n \rightarrow \infty$ . If we take  $\tilde{Y}$  as an independent copy of  $Y$ , then we obtain that the kernel matrix  $\mathbf{V}_h = \Sigma_X^{-1} \mathbf{D}_h \Sigma_X^{-1} = \Sigma_X^{-1} \mathbf{E}[\mathbf{D}_h^{\otimes 2}(\tilde{Y})] \Sigma_X^{-1} \rightarrow \Sigma_X^{-1} \mathbf{E}[\mathbf{D}_B^{\otimes 2}(\tilde{Y})] \Sigma_X^{-1} = \mathbf{V}_B \subseteq \mathbf{S}_{Y|X}$ . So that  $\mathbf{V}_h$  will seek out the central subspace  $\mathbf{S}_{Y|X}$  as  $h \rightarrow 0, n \rightarrow \infty$ .  $\square$

## A.2 Proof of Theorems 2.2

**Proof. Step 2.1** In this step, we prove the consistency of the estimator  $\hat{\mathbf{D}}_h$ .

$$\begin{aligned} \hat{\mathbf{D}}_h &= \frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^n (X_i - \bar{X})(X_j - \bar{X})^\tau K_h(Y_i - Y_s) K_h(Y_j - Y_s) \\ &= \frac{6a_n}{n(n-1)(n-2)} \sum_{1 \leq i < j < s \leq n} u_h((X_i, Y_i), (X_j, Y_j), (X_s, Y_s)) + R_n \\ &= a_n \hat{\mathbf{D}}_h^0 + R_n \end{aligned} \quad (\text{A.2})$$

where  $a_n = \frac{(n-1)(n-2)}{n^2}$  and

$$\begin{aligned} &u_h((X_i, Y_i), (X_j, Y_j), (X_s, Y_s)) \\ &= \frac{1}{6} [(X_i X_j^\tau + X_j X_i^\tau) K_h(Y_i - Y_s) K_h(Y_j - Y_s) \\ &\quad + (X_i X_s^\tau + X_s X_i^\tau) K_h(Y_i - Y_j) K_h(Y_s - Y_j) \\ &\quad + (X_j X_s^\tau + X_s X_j^\tau) K_h(Y_j - Y_i) K_h(Y_s - Y_i)]. \end{aligned} \quad (\text{A.3})$$

From expressions (A.2) and (A.3), we see that  $\hat{\mathbf{D}}_h^0$  is a  $U$ -statistics with symmetric kernel  $u_h(\cdot)$ . Furthermore, as  $h \rightarrow 0$ ,  $n \rightarrow \infty$ ,

$$\begin{aligned}
& \mathbf{E}u_h((X_1, Y_1), (X_2, Y_2), (X_3, Y_3)) \\
&= \mathbf{E}\{\mathbf{E}u_h((X_1, Y_1), (X_2, Y_2), (X_3, Y_3)) | (X_3, Y_3)\} \\
&= \mathbf{E}[\mathbf{D}_h^{\otimes 2}(Y_3)] = \mathbf{E}[\mathbf{D}_h^{\otimes 2}(\tilde{Y})] \\
&\rightarrow \mathbf{E}_{\tilde{Y}}\{\mathbf{E}_X[Xf_{Y|B^\tau X}(\tilde{Y}|B^\tau X)]\mathbf{E}_X[X^\tau f_{Y|B^\tau X}(\tilde{Y}|B^\tau X)] | \tilde{Y}\} \\
&= P_B^\tau \mathbf{E}_{\tilde{Y}}\{\mathbf{E}_X[Xf_{Y|B^\tau X}(\tilde{Y}|B^\tau X)]\mathbf{E}_X[X^\tau f_{Y|B^\tau X}(\tilde{Y}|B^\tau X)] | \tilde{Y}\} P_B \\
&= \mathbf{E}[\mathbf{D}_B^{\otimes 2}(\tilde{Y})] = \mathbf{D}_B,
\end{aligned}$$

where  $\tilde{Y}$  is an independent copy of  $Y$ ,  $P_B = B(B^\tau \Sigma_X B)^{-1} B^\tau \Sigma_X$ , and  $\mathbf{E}_X(\cdot)$  stands for taking expectation about  $X$  and  $\mathbf{E}_{\tilde{Y}}(\cdot)$  stands for taking expectation about  $\tilde{Y}$ . Since  $\mathbf{E}u_h((X_1, Y_1), (X_2, Y_2), (X_3, Y_3)) \rightarrow \mathbf{D}_B$ , thus,  $\mathbf{E}\hat{\mathbf{D}}_h^0 \rightarrow \mathbf{D}_B$ . By the convergence of the  $U$ -statistics, we have  $\hat{\mathbf{D}}_h^0 \xrightarrow{P} \mathbf{D}_B$ . Moreover,  $a_n \rightarrow 1$ , then  $a_n \hat{\mathbf{D}}_h^0 \xrightarrow{P} \mathbf{D}_B$ .

Next, we prove  $R_n = o_P(1)$  in the following.

$$\begin{aligned}
R_n &= \frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^n [(\bar{X}^{\otimes 2} - X_i \bar{X}^\tau - \bar{X} X_j^\tau) K_h(Y_i - Y_s) K_h(Y_j - Y_s)] \\
&\quad + \frac{1}{n^3} \sum_{i=1}^n \sum_{j \neq i} [ (X_i - \bar{X})(X_j - \bar{X})^\tau K_h(0) K_h(Y_j - Y_i) ] \\
&\quad + \frac{1}{n^3} \sum_{j=1}^n \sum_{i \neq j} [ (X_i - \bar{X})(X_j - \bar{X})^\tau K_h(0) K_h(Y_i - Y_j) ] \\
&\quad + \frac{1}{n^3} \sum_{i=1}^n \sum_{s \neq i} [ (X_i - \bar{X})^{\otimes 2} K_h^2(Y_i - Y_s) ] \\
&\quad + \frac{1}{n^3} \sum_{i=1}^n [ (X_i - \bar{X})^{\otimes 2} K_h^2(0) ] \\
&= R_{1n} + R_{2n} + R_{3n} + R_{4n} + R_{5n}.
\end{aligned} \tag{A.4}$$

Note that  $\mathbf{E}X = 0$ , then  $\sqrt{n}\bar{X} = O_P(1)$ . Similar to proof of the consistency

of  $\hat{\mathbf{D}}_h^0$ , we have

$$\begin{aligned}
 R_{1n} &= \bar{X}^{\otimes 2} \frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^n K_h(Y_i - Y_s) K_h(Y_j - Y_s) \\
 &\quad + \left[ \frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^n X_i K_h(Y_i - Y_s) K_h(Y_j - Y_s) \right] \bar{X}^\tau \\
 &\quad + \bar{X} \left[ \frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^n X_j^\tau K_h(Y_i - Y_s) K_h(Y_j - Y_s) \right] \quad (\text{A.5}) \\
 &= O_P \left( \frac{1}{\sqrt{n}} \right).
 \end{aligned}$$

It follows that

$$\begin{aligned}
 R_{2n} &= \frac{K(0)}{nh} \left[ \frac{1}{n^2} \sum_{i=1}^n \sum_{j \neq i}^n (X_i X_j^\tau - X_i \bar{X}^\tau - \bar{X} X_j^\tau - \bar{X}^{\otimes 2}) K_h(Y_i - Y_j) \right] \\
 &= O_P \left( \frac{1}{nh} \right). \quad (\text{A.6})
 \end{aligned}$$

The analysis of  $R_{3n}$ ,  $R_{4n}$  are similar to  $R_{2n}$  and will be  $O_P \left( \frac{1}{nh} \right)$ . For  $R_{5n}$ ,

$$R_{5n} = \frac{K(0)^2}{(nh)^2} \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^{\otimes 2} = O_P \left( \frac{1}{n^2 h^2} \right). \quad (\text{A.7})$$

Based on (A.5)-(A.7), provided  $h \rightarrow 0$  and  $nh \rightarrow \infty$  as  $n \rightarrow \infty$ ,  $R_n = o_P(1)$ .

**Step 2.2** Note that  $\hat{\Sigma}_X = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^{\otimes 2} \xrightarrow{P} \Sigma_X$ , then  $\hat{\Sigma}_X - \Sigma_X = o_P(1)$ . We have proved that  $\hat{\mathbf{D}}_h - \mathbf{D}_B = o_P(1)$ , then

$$\begin{aligned}
 \hat{\mathbf{V}}_h - \mathbf{V}_B &= \hat{\Sigma}_X^{-1} \hat{\mathbf{D}}_h \hat{\Sigma}_X^{-1} - \Sigma_X^{-1} \mathbf{D}_B \Sigma_X^{-1} \quad (\text{A.8}) \\
 &= \hat{\Sigma}_X^{-1} (\Sigma_X - \hat{\Sigma}_X) \Sigma_X^{-1} \hat{\mathbf{D}}_h \hat{\Sigma}_X^{-1} + \Sigma_X^{-1} (\hat{\mathbf{D}}_h - \mathbf{D}_B) \hat{\Sigma}_X^{-1} \\
 &\quad + \Sigma_X^{-1} \mathbf{D}_B \hat{\Sigma}_X^{-1} (\Sigma_X - \hat{\Sigma}_X) \Sigma_X^{-1} = o_P(1),
 \end{aligned}$$

which indicates that  $\hat{\mathbf{V}}_h \xrightarrow{P} \mathbf{V}_B \subseteq \mathcal{S}_{Y|X}$ , and we have completed the proof of Theorem 2.2.  $\square$

### A.3 Proof of Theorem 2.3

**Proof.** In this section, we drive the asymptotic distribution of  $\hat{\mathbf{V}}_h$ . Based on (A.8),

$$\begin{aligned} & \sqrt{n}(\hat{\mathbf{V}}_h - \mathbf{V}_B) \\ &= \sqrt{n}(\hat{\Sigma}_X^{-1} \hat{\mathbf{D}}_h \hat{\Sigma}_X^{-1} - \Sigma_X^{-1} \mathbf{D}_B \Sigma_X^{-1}) \\ &= \hat{\Sigma}_X^{-1} \sqrt{n}(\Sigma_X - \hat{\Sigma}_X) \Sigma_X^{-1} \hat{\mathbf{D}}_h \hat{\Sigma}_X^{-1} + \Sigma_X^{-1} \sqrt{n}(\hat{\mathbf{D}}_h - \mathbf{D}_B) \hat{\Sigma}_X^{-1} \\ & \quad + \Sigma_X^{-1} \mathbf{D}_B \hat{\Sigma}_X^{-1} \sqrt{n}(\Sigma_X - \hat{\Sigma}_X) \Sigma_X^{-1}. \end{aligned}$$

First, we derive  $\sqrt{n}(\hat{\mathbf{D}}_h - \mathbf{D}_B)$  as a sum of i.i.d random variable and an asymptotic negligible part. According to (A.2), (A.4) and (A.6)-(A.7), if  $\sqrt{nh^2} \rightarrow 0$  and  $\sqrt{nh} \rightarrow \infty$ , we have

$$\begin{aligned} & \sqrt{n}(\hat{\mathbf{D}}_h - \mathbf{D}_B) \tag{A.9} \\ &= \sqrt{n}(a_n \hat{\mathbf{D}}_h^0 + R_{1n} - \mathbf{D}_B) + \sqrt{n}(R_{2n} + R_{3n} + R_{4n} + R_{5n}) \\ &= \sqrt{n}(a_n(\hat{\mathbf{D}}_h^0 - \mathbf{E}\hat{\mathbf{D}}_h^0) + R_{1n}) + \sqrt{n}(a_n \mathbf{E}\hat{\mathbf{D}}_h^0 - \mathbf{D}_B) + O_P\left(\frac{1}{\sqrt{nh}}\right) \\ &= \sqrt{n}(a_n(\hat{\mathbf{D}}_h^0 - \mathbf{E}\hat{\mathbf{D}}_h^0) + R_{1n}) + O(a_n \sqrt{nh^2} + \sqrt{n}(a_n - 1)) + O_P\left(\frac{1}{\sqrt{nh}}\right) \\ &= a_n \sqrt{n}(\hat{\mathbf{D}}_h^0 - \mathbf{E}\hat{\mathbf{D}}_h^0) + \sqrt{n}R_{1n} + o_P(1). \end{aligned}$$

Note that  $\hat{\mathbf{D}}_h^0$  is a  $U$ -statistic, and by Hajek projection (Serfling, 1980), we have  $\sqrt{n}(\hat{\mathbf{D}}_h^0 - \mathbf{E}\hat{\mathbf{D}}_h^0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{E}(\hat{\mathbf{D}}_h^0 | (X_i, Y_i)) - \mathbf{E}\hat{\mathbf{D}}_h^0] + o_P(1)$ , which implies

$$\begin{aligned} & \sqrt{n}(\hat{\mathbf{D}}_h^0 - \mathbf{E}\hat{\mathbf{D}}_h^0) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{E}(\hat{\mathbf{D}}_h^0 | (X_i, Y_i)) - \mathbf{E}\hat{\mathbf{D}}_h^0] + O_P(\sqrt{n}(a_n - 1)) + o_P(1) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{E}(\hat{\mathbf{D}}_h^0 | (X_i, Y_i)) - \mathbf{E}\hat{\mathbf{D}}_h^0] + o_P(1). \tag{A.10} \end{aligned}$$

We calculate the second part  $\sqrt{n}R_{1n}$  in (A.5). Based on (A.5),

$$\begin{aligned} \sqrt{n}R_{1n} &= \left[ \frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^n X_i K_h(Y_i - Y_s) K_h(Y_j - Y_s) \right] \sqrt{n} \bar{X}^\tau \\ & \quad + \sqrt{n} \bar{X} \left[ \frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^n X_j^\tau K_h(Y_i - Y_s) K_h(Y_j - Y_s) \right] + O_P\left(\frac{1}{\sqrt{n}}\right). \end{aligned}$$

The analysis of  $\frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^n X_i K_h(Y_i - Y_s) K_h(Y_j - Y_s)$  is similar to  $\hat{D}_h$  and is expressed as a sum of  $U$ -statistic and a negligible part. Note that

$$\begin{aligned} & \mathbf{E}[X_1 K_h(Y_1 - Y_3) K_h(Y_2 - Y_3)] \\ &= P_B^T \mathbf{E}_{\tilde{Y}} \{ \mathbf{E}_X [X f_{Y|B^T X}(\tilde{Y} | B^T X)] g(\tilde{Y}) \} + O(h^2) \\ &= \mathbf{E} \{ \mathbf{D}_B(\tilde{Y}) g(\tilde{Y}) \} + O(h^2) \end{aligned}$$

where  $g(\cdot)$  is the density function of  $Y$ . Then, we have

$$\begin{aligned} \sqrt{n} R_{1n} &= \mathbf{E} \{ \mathbf{D}_B(\tilde{Y}) g(\tilde{Y}) \} \sqrt{n} \bar{X}^T + \sqrt{n} \bar{X} \mathbf{E} \{ \mathbf{D}_B^T(\tilde{Y}) g(\tilde{Y}) \} \\ &\quad + o_P(1). \end{aligned} \quad (\text{A.11})$$

Together with (A.10) and (A.11), we have

$$\begin{aligned} & \sqrt{n}(\hat{D}_h - \mathbf{D}_B) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{E}(\hat{D}_h^0 | (X_i, Y_i)) - \mathbf{E} \hat{D}_h^0 + \mathbf{E} \{ \mathbf{D}_B(\tilde{Y}) g(\tilde{Y}) \} X_i^T \\ &\quad + X_i \mathbf{E} \{ \mathbf{D}_B^T(\tilde{Y}) g(\tilde{Y}) \}] + o_P(1). \end{aligned} \quad (\text{A.12})$$

Second, the  $\sqrt{n}(\hat{\Sigma}_X - \Sigma_X) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i^{\otimes 2} - \Sigma_X) + o_P(\frac{1}{\sqrt{n}})$ . Note that  $\hat{D}_h \xrightarrow{P} \mathbf{D}_B$  and  $\hat{\Sigma}_X \xrightarrow{P} \Sigma_X$ , then

$$\begin{aligned} \sqrt{n}(\hat{V}_h - \mathbf{V}_B) &= \sqrt{n}(\hat{\Sigma}_X^{-1} \hat{D}_h \hat{\Sigma}_X^{-1} - \Sigma_X^{-1} \mathbf{D}_B \Sigma_X^{-1}) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \Sigma_X^{-1} \mathbf{T}_h(X_i, Y_i) \Sigma_X^{-1} + o_P(1), \end{aligned}$$

where

$$\begin{aligned} & \mathbf{T}_h(X_i, Y_i) \\ &= (X_i^{\otimes 2} - \Sigma_X) \Sigma_X^{-1} \mathbf{D}_B + \mathbf{D}_B \Sigma_X^{-1} (X_i^{\otimes 2} - \Sigma_X) \\ &\quad + \mathbf{E}(\hat{D}_h^0 | (X_i, Y_i)) - \mathbf{E} \hat{D}_h^0 + \mathbf{E} \{ \mathbf{D}_B(\tilde{Y}) g(\tilde{Y}) \} X_i^T + X_i \mathbf{E} \{ \mathbf{D}_B^T(\tilde{Y}) g(\tilde{Y}) \}. \end{aligned} \quad (\text{A.13})$$

Note that  $\{\mathbf{T}_h(X_i, Y_i), i = 1, \dots, n\}$  are i.i.d random variables. For any  $a \in \mathbb{R}^{\frac{p(p+1)}{2}}$  and  $a \neq 0$ , if  $a^T \text{Cov}(\text{vech}(\Sigma_X^{-1} \mathbf{T}_h(X_1, Y_1) \Sigma_X^{-1})) a \rightarrow \sigma_a^2 > 0$  as  $h \rightarrow 0$ ,  $n \rightarrow \infty$ , the CLT theorem entails that  $\sqrt{n} a^T (\text{Vech}(\Sigma_X^{-1} \mathbf{T}_h(X_i, Y_i) \Sigma_X^{-1}) - \text{Vech}(\Sigma_X^{-1} \mathbf{E} \mathbf{T}_h(X_i, Y_i) \Sigma_X^{-1})) \xrightarrow{L} N(0, \sigma_a^2)$ , and we write it as

$$\sqrt{n}(\hat{V}_h - \mathbf{V}_B) \xrightarrow{L} \mathcal{H}. \quad (\text{A.14})$$

We have completed the proof of Theorem 2.3.  $\square$

#### A.4 Proof of Theorem 2.4

**Proof.** In this section, we get the asymptotic distribution of the nonzero eigenvalues of  $\hat{\mathbf{V}}_h$ . Denote  $\lambda_1 > \lambda_2 > \dots > \lambda_l$  is the distinct eigenvalues of positive semi-definite matrix  $\mathbf{V}_B$  with the multiplicity of  $\lambda_i$  being  $m_i$ ,  $i = 1, \dots, l$  and  $m_1 + m_2 + \dots + m_l = p$ . There exists orthogonal matrices  $Q$  such that

$$Q^T \mathbf{V}_B Q = \begin{bmatrix} \lambda_1 I_{m_1} & 0 & \dots & 0 \\ 0 & \lambda_2 I_{m_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_l I_{m_l} \end{bmatrix}$$

From the results (A.14), we have  $\sqrt{n}(Q^T \hat{\mathbf{V}}_h Q - Q^T \mathbf{V}_B Q) \xrightarrow{L} Q^T \mathcal{H} Q$ , and we partition  $Q^T \hat{\mathbf{V}}_h Q$  in the similar way of  $Q^T \mathbf{V}_B Q$  as

$$Q^T \hat{\mathbf{V}}_h Q = \begin{bmatrix} V_{n,11} & V_{n,12} & \dots & V_{n,1l} \\ V_{n,21} & V_{n,22} & \dots & V_{n,2l} \\ \vdots & \vdots & \ddots & \vdots \\ V_{n,l1} & V_{n,l2} & \dots & V_{n,ll} \end{bmatrix}$$

We define that  $\lambda(A)$  stands for the vector of ordered eigenvalues of  $A$ , i.e., denoted as  $\lambda(A) = (\lambda_1(A), \lambda_2(A), \dots, \lambda_k(A))^T$ , and let  $1_{m_i} \in \mathbb{R}^{m_i}$  stand for the vector of ones for  $i = 1, \dots, l$ . Applying the result of Theorem 3.1. in Eaton and Tyler (1991), we have

$$\sqrt{n}(\lambda(Q^T \hat{\mathbf{V}}_h Q) - \lambda(Q^T \mathbf{V}_B Q)) = \sqrt{n} \begin{bmatrix} \lambda(V_{n,11}) & - & \lambda_1 1_{m_1} \\ \lambda(V_{n,22}) & - & \lambda_2 1_{m_2} \\ \vdots & & \vdots \\ \lambda(V_{n,ll}) & - & \lambda_l 1_{m_l} \end{bmatrix} + o_P(1)$$

Since  $\sqrt{n}(Q^T \hat{\mathbf{V}}_h Q - Q^T \mathbf{V}_B Q) \xrightarrow{L} Q^T \mathcal{H} Q$ , we have

$$\sqrt{n} \begin{bmatrix} V_{n,11} & - & \lambda_1 I_{m_1} \\ V_{n,22} & - & \lambda_2 I_{m_2} \\ \vdots & & \vdots \\ V_{n,ll} & - & \lambda_l I_{m_l} \end{bmatrix} \xrightarrow{L} \tilde{\mathcal{H}} = \begin{bmatrix} \tilde{\mathcal{H}}_{11} \\ \tilde{\mathcal{H}}_{22} \\ \vdots \\ \tilde{\mathcal{H}}_{ll} \end{bmatrix},$$

where  $\tilde{\mathcal{H}} = (\tilde{\mathcal{H}}_{i,j})$  is the partitioning of  $Q^T \mathcal{H} Q$  in blocks of order  $m_i \times m_j$ . Moreover, the eigenvalue vector  $\lambda(A)$  is a continuous function about matrix

$A$ , and  $Q$  is an orthogonal matrix, then  $\lambda(Q^\tau \hat{\mathbf{V}}_h Q) = \lambda(\hat{\mathbf{V}}_h)$ ,  $\lambda(Q^\tau \mathbf{V}_B Q) = \lambda(\mathbf{V}_B)$ . We have

$$\sqrt{n}(\lambda(\hat{\mathbf{V}}_h) - \lambda(\mathbf{V}_B)) \xrightarrow{L} \begin{bmatrix} \lambda(\tilde{\mathcal{H}}_{11}) \\ \lambda(\tilde{\mathcal{H}}_{22}) \\ \vdots \\ \lambda(\tilde{\mathcal{H}}_{ll}) \end{bmatrix}.$$

We have completed the proof of Theorem 2.4.  $\square$

#### A.5 Proof of Theorem 2.5

**Proof.** For any vector  $e$  is orthogonal to  $\mathbf{S}_{Y|X}$ , satisfying  $e^\tau e = 1$ . Note that  $\hat{\mathbf{V}}_h \hat{b}_j = \hat{\lambda}_j \hat{b}_j$ ,  $j = 1, \dots, k$ , then

$$\begin{aligned} \sqrt{n} e^\tau \hat{b}_j &= \frac{\sqrt{n}}{\hat{\lambda}_j} e^\tau \hat{\mathbf{V}}_h \hat{b}_j \\ &= \frac{1}{\lambda_j} e^\tau \sqrt{n} \{ \hat{\mathbf{V}}_h - \mathbf{V}_B \} b_j + \frac{\lambda_j - \hat{\lambda}_j}{\lambda_j \hat{\lambda}_j} e^\tau \sqrt{n} \{ \hat{\mathbf{V}}_h - \mathbf{V}_B \} b_j \\ &\quad + \frac{1}{\lambda_j} e^\tau \sqrt{n} \{ \hat{\mathbf{V}}_h - \mathbf{V}_B \} \{ \hat{b}_j - b_j \} + \frac{1}{\hat{\lambda}_j} \sqrt{n} e^\tau \mathbf{V}_B \hat{b}_j \\ &= I_{n,1} + I_{n,2} + I_{n,3} + I_{n,4}. \end{aligned}$$

Using the proof of Theorem 2.3, we have

$$\sqrt{n} \{ \hat{\mathbf{V}}_h - \mathbf{V}_B \} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \Sigma_X^{-1} \mathbf{T}_h(X_i, Y_i) \Sigma_X^{-1} + o_P(1),$$

where  $\mathbf{T}_h(X_i, Y_i)$  is defined in (A.13), then

$$I_{n,1} = \sqrt{n} e^\tau \{ \hat{\mathbf{V}}_h - \mathbf{V}_B \} b_j = \frac{1}{\sqrt{n}} \sum_{i=1}^n e^\tau \Sigma_X^{-1} \mathbf{T}_h(X_i, Y_i) \Sigma_X^{-1} b_j + o_P(1).$$

Here,  $\{e^\tau \Sigma_X^{-1} \mathbf{T}_h(X_i, Y_i) \Sigma_X^{-1} b_j\}_{i=1}^n$  are i.i.d random variables, and if the limit of variance  $e^\tau \text{Cov}(\Sigma_X^{-1} \mathbf{T}_h(X_i, Y_i) \Sigma_X^{-1} b_j) e \rightarrow e^\tau W_j e > 0$  as  $h \rightarrow 0$ ,  $n \rightarrow \infty$ . The CLT theorem entails that

$$I_{n,1} \xrightarrow{L} N(0, e^\tau W_j e). \quad (\text{A.15})$$

Next, we show  $I_{n,2}, I_{n,3}, I_{n,4}$  are  $o_P(1)$ . Theorem 2.4 entails that  $\hat{\lambda}_j = \lambda_j + O_P(\frac{1}{\sqrt{n}})$ , and using  $I_{n,1} = O_P(1)$ , we have

$$I_{n,2} = \frac{\lambda_j - \hat{\lambda}_j}{\lambda_j \hat{\lambda}_j} I_{n,1} = \frac{O_P(\frac{1}{\sqrt{n}})}{\lambda_j(\lambda_j + O_P(\frac{1}{\sqrt{n}}))} O_P(1) = o_P(1). \quad (\text{A.16})$$

By the perturbation theory and  $\hat{\mathbf{V}}_h = \mathbf{V}_B + o_P(1)$ , we obtain  $\hat{b}_j = b_j + o_P(1)$ , and then

$$I_{n,3} = \frac{1}{\lambda_j} e^\tau \sqrt{n} \{ \hat{\mathbf{V}}_h - \mathbf{V}_B \} \{ \hat{b}_j - b_j \} = O_P(1) o_P(1) = o_P(1). \quad (\text{A.17})$$

As  $e$  is orthogonal to  $\mathbf{S}_{Y|X}$ , and  $\mathbf{V}_B \subseteq \mathbf{S}_{Y|X}$ , then  $e^\tau \mathbf{V}_B = 0$ , and also

$$I_{n,4} = \frac{1}{\hat{\lambda}_j} \sqrt{n} e^\tau \mathbf{V}_B \hat{b}_j = o_P(1). \quad (\text{A.18})$$

Together with the result (A.15)-(A.18), we have

$$\sqrt{n} e^\tau \hat{b}_j = I_1 + I_2 + I_3 + I_4 = I_1 + o_P(1) \xrightarrow{L} N(0, e^\tau W_j e).$$

We have completed the proof of Theorem 2.5.  $\square$

## Acknowledgements

The authors thank the editor, the associate editor and a referee for their constructive suggestions that helped us to improve the early manuscript. The research work of Zhang Jun is supported by the National Natural Sciences Foundation of China (Grant No. 11401391).

## References

- Chen, C. H. and Li, K. C. (1998) Can SIR be as popular as multiple linear regression? *Statist. Sinica*, **8**, 289–316.
- Cook, R. D. (1994) On the interpretation of regression plots. *J. Amer. Statist. Assoc.* **89**, 177–189.
- Cook, R. D. (1998) *Regression Graphics. Ideas for studying regressions through graphics*. New York: Wiley.
- Cook, R. D. and Weisberg, S. (1991) Discussion of “Sliced inverse regression for dimension reduction”. by K.-C.Li. *J. Amer. Statist. Assoc.*, **86**, 328–332.
- Eaton, M. and Tyler, D. (1991) On Wielandt’s Inequality and Its Applications to the Asymptotic Distribution of the Eigenvalues of a Random Symmetric Matrix. *Ann. Statist.* **19**, 260–271.
- Fan, J., Yao, Q. and Tong, H. (1996) Estimation of conditional densities and sensitivity measures in nonlinear dynamical systems. *Biometrika*. **83**, 189–196.
- Hotelling H. (1936) Relations between two sets of variates. *Biometrika*. **28**, 321–377.
- Li, B., Zha, H. and Chiaromonte, F. (2005) Contour regression: A general approach to dimension reduction. *Ann. Statist.* **33**, 1580–1616.
- Li, B., Wang S. L. (2007) On directional regression for dimension reduction. *J. Amer. Statist. Assoc.* **102**, 997–1008.
- Li, K. C. (1991) Sliced inverse regression for dimension reduction (with discussion). *J. Amer. Statist. Assoc.* **86**, 316–342.

- Li, K. C. (1992) On principal Hessian directions for data visualization and dimension reduction: another application of Stein's lemma. *J. Amer. Statist. Assoc.* **87**, 1025–1039.
- Li, K. C. and Duan, N. (1989) Regression analysis under link violation. *Ann. Statist.* **17**, 1009–1052.
- Serfling, R. J. (1980) *Approximation theorems of mathematical statistics*. New York: Wiley.
- Stein, C. (1981) Estimation the mean of a multivariate normal distribution. *Ann. Statist.* **9**, 1135–1151.
- Ye, Z. and Weiss, R. E. (2003). Using the bootstrap to select one of a new class of dimension reduction methods. *J. Amer. Statist. Assoc.* **98**, 968–979.
- Samarov, A. M. (1993) Exploring regression structure using nonparametric functional estimation. *J. Amer. Statist. Assoc.* **88**, 836–847.
- Zhu, L. X. Miao, B. Q. and Peng, H. (2006) Sliced inverse regression with large dimensional covariates. *J. Amer. Statist. Assoc.* **101**, 630–643.
- Wang, T. and Zhu, L. X. (2013) Sparse sufficient dimension reduction using optimal scoring. *Comput. Statist. Data Anal.* **57**, 223–232.
- Xia, Y. (2007) A constructive approach to the estimation of dimension reduction directions. *Ann. Statist.* **35**, 2654–2690.
- Xia, Y., Tong, H., Li, W. K. and Zhu, L. X. (2002). An adaptive estimation of dimension reduction (with discussion). *J. R. Statist. Soc. B.* **64**, 363–410.
- Zhu, X., Chen, F., Guo, X. and Zhu, L. X. (2016) Heteroscedasticity testing for regression models: a dimension reduction-based model adaptive approach. *Comput. Statist. Data Anal.* **103**, 263–283.
- Zhu, X., Chen, F., Guo, X. and Zhu, L. X. (2016) An adaptive-to-model test for partially parametric single-index models. *Stat. Comput.*
- Lansangan, Joseph Ryan G. and Barrios, Erniel B. (2017) Simultaneous dimension reduction and variable selection in modeling high dimensional data. *Comput. Statist. Data Anal.* **112**, 242–256.
- Luo, Wei and Zhu, Yeying and Ghosh, Debashis. (2017) On estimating regression-based causal effects using sufficient dimension reduction. *Biometrika* **104**, 51–65.
- Yoshida, Takuma. (2017) Nonlinear surface regression with dimension reduction method. *ASTA Adv. Stat. Anal.* **101**, 29–50.
- Deng, Jianqiu and Wang, Qihua. (2017) Dimension reduction estimation for probability density with data missing at random when covariables are present. *J. Statist. Plann. Inference.* **181**, 11–29.
- Sheng, Wenhui and Yin, Xiangrong. (2016) Sufficient dimension reduction via distance covariance. *J. Comput. Graph. Statist.* **25**, 91–104.
- Zhou, Jingke and Zhu, Lixing. (2016) Principal minimax support vector machine for sufficient dimension reduction with contaminated data. *Comput. Statist. Data Anal.* **94**, 33–48.
- Cui, Xia and Härdle, Wolfgang Karl and Zhu, Lixing. (2009) The EFM approach for single-index models. *Ann. Statist.* **39**, 1658–1688.
- Peng, Heng and Huang, Tao. (2011) Penalized least squares for single index models. *J. Statist. Plann. Inference.* **141**, 1362–1379.
- Chen, Xin and Zou, Changliang and Cook, Dennis R. (2010) Coordinate-independent sparse sufficient dimension reduction and variable selection. *Ann. Statist.* **38**, 3696–3723.
- Li, Gaorong and Peng, Heng and Dong, Kai and Tong, Tiejun. (2014) Simultaneous confidence bands and hypothesis testing for single-index models *Statist. Sinica*, **24**, 937–955.

- Liang, H. and Liu, X. and Li, R. and Tsai, C. L. (2010) Estimation and testing for partially linear single-index models. *Ann. Statist.* **38**, 3811–3836.
- Li, Bing and Yin, Xiangrong. (2007) On surrogate dimension reduction for measurement error regression: an invariance law. *Ann. Statist.* **35**, 2143–2172.
- Wu, Yichao and Li, Lexin. (2011) Asymptotic properties of sufficient dimension reduction with a diverging number of predictors. *Statist. Sinica*, **21**, 707–730.
- Ding, Xiaobo and Wang, Qihua. (2011) Fusion-refinement procedure for dimension reduction with missing response at random. *J. Amer. Statist. Assoc.* **106**, 1193–1207.
- Guo, Xu and Wang, Tao and Xu, Wangli and Zhu, Lixing. (2014) Dimension reduction with missing response at random. *Comput. Statist. Data Anal.* **69**, 228–242.
- Bi, Xuan and Qu, Annie. (2015) Sufficient dimension reduction for longitudinal data. *Statist. Sinica*, **251**, 787–807.
- Li, Lexin and Yin, Xiangrong. (2009) Longitudinal data analysis using sufficient dimension reduction method. *Comput. Statist. Data Anal.* **53**, 4106–4115.

J. Zhang  
College of Mathematics and Statistics  
Institute of Statistical Sciences  
Shenzhen-Hong Kong Joint Research  
Center for Applied Statistical Sciences  
Shenzhen  
China

B.He  
School of Mathematics and Computational Science  
Sun Yat-Sen University  
Guangzhou  
China  
E-mail: [mrhephd@163.com](mailto:mrhephd@163.com)

T.Lu  
Department of Mathematics and Statistics  
University of Nevada  
Reno, NV  
USA

S.Wen  
College of Mathematics and Statistics  
Shenzhen University  
Shenzhen  
China