

On Sliced Inverse Regression With High-Dimensional Covariates

Lixing ZHU, Baiqi MIAO, and Heng PENG

Sliced inverse regression is a promising method for the estimation of the central dimension-reduction subspace (CDR space) in semiparametric regression models. It is particularly useful in tackling cases with high-dimensional covariates. In this article we study the asymptotic behavior of the estimate of the CDR space with high-dimensional covariates, that is, when the dimension of the covariates goes to infinity as the sample size goes to infinity. Strong and weak convergence are obtained. We also suggest an estimation procedure of the Bayes information criterion type to ascertain the dimension of the CDR space and derive the consistency. A simulation study is conducted.

KEY WORDS: Central dimension-reduction subspace; Convergence rate; Dimensionality determination; Sliced inverse regression.

1. INTRODUCTION

Large-scale data analysis has recently received much attention. Reducing the dimensions of data without loss of information is a natural manner in which to proceed with such analysis. In this area, Li (1991) considered a semiparametric regression model in which the response variable y depends on the covariates $\mathbf{x} = (x_1, \dots, x_p)^T$ through K linear combinations of the components, x_i 's. His model can be generalized to a conditional independence statement. Let β_1, \dots, β_K be K orthogonal $p \times 1$ vectors, the norms of which are $\|\beta_i\| = \sqrt{\beta_{i1}^2 + \dots + \beta_{ip}^2} = 1$ for $i = 1, \dots, K$, where β_{ij} , $j = 1, \dots, p$, are the components of β_i . When $(\beta_1^T \mathbf{x}, \dots, \beta_K^T \mathbf{x})$ is given, y is independent of \mathbf{x} , that is,

$$y \perp\!\!\!\perp \mathbf{x} | (\beta_1^T \mathbf{x}, \dots, \beta_K^T \mathbf{x}), \quad (1)$$

where " $\perp\!\!\!\perp$ " represents independence.

A dimension-reduction subspace (Cook 1994, 1998) is defined as the column space of any $p \times K$ ($K \leq p$) matrix $\mathbf{B} = (\beta_1, \dots, \beta_K)$ such that (1) holds. For uniqueness, we are interested in a subspace with minimal dimensions. Under mild conditions, the minimal subspace is often uniquely defined in practice and coincides with the intersection of all subspaces that satisfy (1) (Cook 1994, 1996). This intersection is called the *central dimension-reduction (CDR) space* and is written as $S_{y|\mathbf{x}}$. In this article we assume that $S_{y|\mathbf{x}}$ exists.

Sliced inverse regression (SIR) (Li 1991) can be used to estimate $S_{y|\mathbf{x}}$. This method is popular in the literature. Li (1991) derived the asymptotic results under the assumption that \mathbf{x} is normally distributed, Hsing and Carroll (1992) proved the asymptotic normality of the slicing estimator when each slice contains two points, Zhu and Ng (1995) obtained the asymptotic normality for general cases, and Zhu and Fang (1996) studied the asymptotic behavior of a kernel estimator. Other important estimation approaches have also been treated in the literature, including sliced average variance estimates (SAVE)

(Cook and Weisberg 1991; Cook 2000), parametric inverse regression (PIR) (Bura and Cook 2001a,b), and the proposed hybrid methods of Zhu, Ohtaki, and Li (2006) that are convex combinations of SIR and SAVE. Another relevant work is that of Cook and Li (2002).

All of the existing studies of asymptotics for model (1) are related to cases with a fixed dimension p and a sample size, say n , that is comparably high. However, it is often the case that the number of covariates is also large. Li (1991) first raised this issue and suggested the importance of studying it further. Recent research with covariates of a very high dimension also demonstrates this necessity. These kinds of datasets include gene expression data, DNA microarray data, and consumer financial history data. Golub et al. (1999) brought to life a whole new branch of data analysis under the name of microarray analysis. Data with high-dimensional covariates also appear in other areas. Zhu, Zhu, and Li (2006) analyzed a dataset of agricultural meteorological disasters. They explored the relationship between the yields of three crops and meteorological conditions. The dimension p of the covariate was also comparably high. Greenshtain and Ritov (2004) studied the asymptotics of lasso methods when p goes to infinity as n tends to infinity. Bickel and Levina (2004) investigated the asymptotic properties for Fisher's linear discriminant analysis (for relevant data analysis see Levina 2002). Conventional methods are limited for this kind of data analysis, and dimension-reduction methods may be helpful. Among others, Antoniadis, Lambert-Lacroix, and Leblanc (2003) applied a relevant method, MAVE, proposed by Xia, Tong, Li, and Zhu (2002), to analyze gene expression data for classification, and Bura and Pfeiffer (2003) used SIR and SAVE for DNA microarray data. Related work was done by Chiaromonte and Martinelli (2002).

Furthermore, to apply SIR to estimate $S_{y|\mathbf{x}}$, the linearity condition as designed by Li (1991) must be checked. Note that the linearity condition cannot be always met in cases in which the dimensions are small or moderate. Hall and Li (1993) proved that when $p \rightarrow \infty$ as $n \rightarrow \infty$, the linear combinations of the covariates are approximately normally distributed. This ensures that the linearity condition for SIR is asymptotically satisfied. Therefore, for cases in which the dimensions are high, SIR can work well. This stimulates us to study the asymptotic behavior of SIR with high-dimensional covariates.

Lixing Zhu is Professor of Statistics, Department of Mathematics, Hong Kong Baptist University, and Cheung Kong Chair Professor at Renmin University of China under the Cheung Kong Scholars Program of the Ministry of Education and Li Ka Shing Foundation, Hong Kong, People's Republic of China (E-mail: lzhu@hkbu.edu.hk). Baiqi Miao is Professor, University of Science and Technology of China. Heng Peng is Post-Doctoral Fellow, Princeton University, Princeton, NJ 08544. This research was supported by a grant from the University Grants Council of Hong Kong (HKU7181/02H). The authors are grateful to the editor, an associate editor, and the two referees for their constructive comments and suggestions, which led to the great improvement of earlier draft. They also thank the associate editor and one of the referees for their generous help in improving the presentation of the article.

The asymptotic research for the $p \rightarrow \infty$ case is very challenging, however. The literature contains some relevant studies on covariance matrices. It is impossible to uniformly formulate the asymptotic results for different ratios $\lambda = p/n$ when p diverges to infinity as n tends to infinity (see Bai and Saranadasa 1996; Bai 1999 for details). This demonstrates the significant difference in convergence between the case with a fixed p and that with a varying p .

The estimate for a matrix based on SIR has a much more complicated structure than the sample covariance matrix. In this article we explore the asymptotic structure of SIR matrix estimate, then investigate the strong and weak convergence when $p \rightarrow \infty$ as $n \rightarrow \infty$.

Another important issue is determinating the dimension of $\mathcal{S}_{y|x}$. For SIR, Li (1991) suggested a sequential chi-squared test procedure to determine the dimensions. Schott (1994) extended this procedure to the situation in which the vector of explanatory variables is sampled from an elliptically symmetric distribution. Velilla (1998) also considered a sequential test. Bura and Cook (2001a) suggested a general weighted chi-squared sequential test that does not require the assumption that covariates follow a normal distribution. It should be noted that the significance level at each step in a sequential test procedure does not determine the significance level of the entire procedure, and that the retained dimension depends on the choice of significance level. Taking this into account, Ferré (1998) proposed an approach to determine the dimension of $\mathcal{S}_{y|x}$. Although his method is somewhat similar to that of Li (1991), it does not need a sequential test, and it selects the dimension automatically.

We suggest a new procedure of the Bayes information criterion (BIC) (Schwarz 1978) type for determinating dimensions. Consistency of the estimator of the dimensions is obtained.

To the best of our knowledge, this is the first article to address the problem of high-dimensional covariates in the context of SIR, and it also presents a useful new methodology for dimension estimation. For example, the results can be used to deal with the problem of agricultural meteorological disasters, as studied by Zhu, Ohtaki, and Li (2006).

The article is organized as follows. In Section 2 we give a brief description of SIR. We provide the consistency results when $p \rightarrow \infty$ as $n \rightarrow \infty$ in Section 3. In Section 4 we present our BIC-type procedure for determining the dimension of $\mathcal{S}_{y|x}$ and the consistency result for this procedure. We report simulations in Section 5 and give technical proofs of the theorems in the Appendix.

2. A BRIEF DESCRIPTION OF SLICED INVERSE REGRESSION

Denote by $\mathbf{P}_{(\cdot)}$ the projection operator in the standard inner product (see Cook 1998). When \mathbf{x} is standardized, assume that the linearity condition is satisfied, that is,

$$E(\mathbf{x}|\mathbf{P}_{\mathcal{S}_{y|x}}\mathbf{x}) = \mathbf{P}_{\mathcal{S}_{y|x}}\mathbf{x}. \quad (2)$$

Then, under (1), the centered “inverse regression,” $E(\mathbf{x}|y) - E(\mathbf{x})$, is confined to $\mathcal{S}_{y|x}$. This means that

$$\mathcal{S}_{E(\mathbf{x}|y)} \subset \mathcal{S}_{y|x},$$

where $\mathcal{S}_{E(\mathbf{x}|y)}$ represents the space spanned by the centered “inverse regression” $E(\mathbf{x}|y) - E(\mathbf{x})$. As Chiaromonte, Cook, and Li (2002) discussed, if we assume that

$$y \perp\!\!\!\perp \mathbf{x} | \mathbf{P}_{\mathcal{S}_{E(\mathbf{x}|y)}}\mathbf{x}, \quad (3)$$

and if (2) holds, then $\mathcal{S}_{E(\mathbf{x}|y)} = \mathcal{S}_{y|x}$. In other words, an estimate of $\mathcal{S}_{E(\mathbf{x}|y)}$ would provide a good approximation to $\mathcal{S}_{y|x}$. If condition (3) is not satisfied, then, following the suggestion of Li (1991) and Cook and Weisberg (1991), we can explore higher conditional moments of \mathbf{x} given y or consider the so-called “central k th-moment subspace” (CKMS) (Yin and Cook 2002). We assume that $\mathcal{S}_{E(\mathbf{x}|y)} = \mathcal{S}_{y|x}$ and restrict our attention to SIR.

In the general case, $\Sigma_{\mathbf{x}}^{-1}\mathcal{S}_{E(\mathbf{x}|y)} = \mathcal{S}_{y|x}$. Note that $\{\Sigma_{\mathbf{x}}^{-1}\eta_1, \dots, \Sigma_{\mathbf{x}}^{-1}\eta_K\}$ is a basis of $\mathcal{S}_{y|x}$, where η_1, \dots, η_K are the eigenvectors associated with the K largest eigenvalues of the matrix $\text{cov}\{E(\mathbf{x}|y)\}$. Therefore, to estimate $\mathcal{S}_{y|x}$, we need only construct an estimate of $\text{cov}\{E(\mathbf{x}|y)\}$. In contrast, it is well known that

$$\Sigma_{\mathbf{x}} = \text{cov}\{E(\mathbf{x}|y)\} + E\{\text{cov}(\mathbf{x}|y)\} =: \text{cov}\{E(\mathbf{x}|y)\} + \Lambda_p.$$

Thus, by estimating Λ_p and $\Sigma_{\mathbf{x}}$, we can also obtain an estimate of $\text{cov}\{E(\mathbf{x}|y)\}$, and then the estimate of $\mathcal{S}_{y|x}$.

Alternatively, we can derive the estimate of $\mathcal{S}_{y|x}$ using the standardized variable $\mathbf{z}, \mathbf{z} = \Sigma_{\mathbf{x}}^{-1/2}\{\mathbf{x} - E(\mathbf{x})\}$. As is known, $\mathcal{S}_{y|z} = \Sigma_{\mathbf{x}}^{1/2}\mathcal{S}_{y|x}$ (see Li 1991; Cook 1998, chaps. 10 and 11). Hence $\{\Sigma_{\mathbf{x}}^{-1/2}\eta'_1, \dots, \Sigma_{\mathbf{x}}^{-1/2}\eta'_K\}$ is also a basis of $\mathcal{S}_{y|x}$, where η'_1, \dots, η'_K are the eigenvectors associated with the K largest eigenvalues of matrix $\text{cov}\{E(\mathbf{z}|y)\}$.

Li (1991) suggested two ways to estimate the basis of $\mathcal{S}_{y|x}$. Once we have an estimate $\hat{\Sigma}_{\mathbf{x}}$ of $\Sigma_{\mathbf{x}}$, we can either estimate $\text{cov}\{E(\mathbf{x}|y)\}$ directly or estimate Λ_p first and then estimate $\text{cov}\{E(\mathbf{x}|y)\}$ by $\Sigma_{\mathbf{x}} - \Lambda_p$. The same can be applied to \mathbf{z} . The following estimate is for Λ_p .

Let $\{(y_i, \mathbf{x}_i), i = 1, \dots, n\}$ be an iid sample and, according to the value of y_i , denote the order statistics by $(y_{(i)}, \mathbf{x}_{(i)}), i = 1, \dots, n$, where $y_{(1)} \leq \dots \leq y_{(n)}$ and $\mathbf{x}_{(i)}$ are called the concomitant of $y_{(i)}$ (see Yang 1977). We introduce a double subscript (h, j) in which the first element refers to the slice number and the second number refers to the order number of an observation in the given slice, that is,

$$y_{(h,j)} = y_{(c(h-1)+j)}, \quad \mathbf{x}_{(h,j)} = \mathbf{x}_{(c(h-1)+j)},$$

where $c > 0$ is the number of $y_{(i)}$ in every slice. The estimate $\hat{\Lambda}_p$ of Λ_p has the form

$$\hat{\Lambda}_p = \frac{1}{H} \sum_{h=1}^H \left\{ \frac{1}{c-1} \sum_{j=1}^c \left(\mathbf{x}_{(h,j)} - \frac{1}{c} \sum_{\ell=1}^c \mathbf{x}_{(h,\ell)} \right) \times \left(\mathbf{x}_{(h,j)} - \frac{1}{c} \sum_{\ell=1}^c \mathbf{x}_{(h,\ell)} \right)^T \right\}, \quad (4)$$

where $H = [(n+c-1)/c]$ is the number of slices where $[a]$ is the largest integer part of a . In practice, the number of $y_{(i)}$ in the last slice may be less than c , but this has little effect on SIR, because the number of slices H is very large.

When the dimension p of \mathbf{x} is fixed, Hsing and Carroll (1992) proved the root- n consistency of the estimate $\hat{\Lambda}_p$ for the case in which each slice contains two of the ordered $y_{(i)}$, that is, $c = 2$.

Zhu and Ng (1995) established the root- n consistency for the case where c is an arbitrarily fixed number and the case where $c \rightarrow \infty$ as $n \rightarrow \infty$. We study the asymptotic properties of $\hat{\Lambda}_p$ and $\hat{\Sigma}_x - \Lambda_p$.

3. ASYMPTOTIC PROPERTIES OF SLICED INVERSE REGRESSION WITH LARGE p

Let (y_i, \mathbf{x}_i) , $i = 1, \dots, n$, be an independent and identically distributed (iid) sample. Denote the inverse regression curve and the associated residual by $\mathbf{m}(y) = E(\mathbf{x}|y)$ and

$$\varepsilon_i = \mathbf{x}_i - \mathbf{m}(y_i) \quad \text{and} \quad \varepsilon_{(i)} = \mathbf{x}_{(i)} - \mathbf{m}(y_{(i)}), \quad i = 1, \dots, n. \quad (5)$$

Here $\mathbf{m}(y)$ is a p -dimensional vector with components $m_i(y) = E(x_i|y)$, $i = 1, \dots, p$. Note that the $\varepsilon_{(i)}$'s are also concomitants of the $y_{(i)}$. It is clear that the ε_i 's are iid and that (see Yang 1977) they are independent with mean 0 when the order statistics, $y_{(i)}$'s, are given. Because $p \rightarrow \infty$ as $n \rightarrow \infty$, all p , \mathbf{x}_i , Σ_x , $y_{(i)}$, $\varepsilon_{(i)}$, and $\mathbf{m}(y_{(i)})$ depend on n and should be denoted by p_n , $x_i^{(n)}$, $\Sigma_x^{(n)}$, $y_{(i)}^{(n)}$, $\varepsilon_{(i)}^{(n)}$, and $\mathbf{m}_n(y_{(i)}^{(n)})$. For notational simplicity, we omit " n " from the subscripts and superscripts unless stated otherwise.

To study the asymptotic behavior, we need smoothness conditions on the inverse regression curve $\mathbf{m}(y)$. Let $\Pi_n(B)$ be the collection of all n -point partitions, $-B \leq y_{(1)}^* \leq \dots \leq y_{(n)}^* \leq B$, of the interval $[-B, B]$, where $B > 0$ and $n \geq 1$. Any vector-valued or real-valued function, $\mathbf{m}(y)$, is said to have a *total variation of order r* if, for any fixed $B > 0$,

$$\lim_{n \rightarrow \infty} \frac{1}{n^r \sqrt{p}} \sup_{\Pi_n(B)} \sum_{i=1}^{n-1} \|\mathbf{m}(y_{(i+1)}^*) - \mathbf{m}(y_{(i)}^*)\| = 0, \quad (6)$$

where $\|\cdot\|$ is the Euclidean norm—namely, $\|\mathbf{A}\|$ is the root of the squared sum of all elements of \mathbf{A} . A similar definition of total variation has been given by Hsing and Carroll (1992) and Zhu and Ng (1995), except for the value \sqrt{p} in the denominator. We have this value because in the Euclidean norm of $\mathbf{m}(y_{(i+1)}^*) - \mathbf{m}(y_{(i)}^*)$ there are p terms to be summed. We use \sqrt{p} to average these terms. There are many functions of total variation; a special case is the function that satisfies the Lipschitz condition. This can be easily seen because for some $C > 0$, $\|\mathbf{m}(y_{(i+1)}^*) - \mathbf{m}(y_{(i)}^*)\| \leq C\sqrt{p}\|y_{(i+1)}^* - y_{(i)}^*\|$ and $(1/\sqrt{p}) \sum_{i=1}^n \|\mathbf{m}(y_{(i+1)}^*) - \mathbf{m}(y_{(i)}^*)\| \leq C\|y_{(n)}^* - y_{(1)}^*\| \leq 2CB$. Thus, such a function is of a total variation of any order $r > 0$.

Furthermore, for the components $m_i(\cdot)$ of $\mathbf{m}(\cdot)$, if there is a nondecreasing real-valued function M and a real number B_0 such that for any two points, say y_1 and y_2 , both in $(-\infty, -B_0]$ or $[B_0, \infty)$,

$$|m_i(y_1) - m_i(y_2)| \leq |M(y_1) - M(y_2)|, \quad i = 1, \dots, p, \quad (7)$$

then we say that the function $\mathbf{m}(y)$ is *nonexpansive in the metric of M_{B_0} on both sides of B_0* .

The following theorem states the strong and weak convergence of $\hat{\Lambda}_p - \Lambda_p$.

Theorem 1. Assume that the following conditions are satisfied:

(A) $\sup_{i \leq p} E(|x_i|^l) < C_1$ for some constants $l \geq 4$ and C_1 , where $\mathbf{x} = (x_1, \dots, x_p)^T$.

(B) The inverse regression function $\mathbf{m}(y)$ has a total variation of order r , $0 < r < 1/4$. When p converges to infinity, y_i are dependent on n , and we assume the following condition.

(C) Let $F_n(\cdot)$ be the distribution from which the sample y_i 's are drawn. Let Y be the random variable following the distribution $F_n(\cdot)$. Suppose that there are two random variables, W_1 and W_2 , with distributions G_1 and G_2 such that $G_1(x) \leq F_n(x) \leq G_2(x)$; that is, $W_2 < Y < W_1$, where " $<$ " represents "stochastically less than." Similarly " $>$ " represents "stochastically greater than."

(D) For any B_0 , $\mathbf{m}(y)$ is nonexpansive in the metric of $M_{B_0}(y)$ on both sides of a positive number B_0 , such that $l \times r > 2$ for the l of condition (A),

$$E(|M(W_1)|^l) < \infty \quad \text{and} \quad E(|M(W_2)|^l) < \infty.$$

Then we have for any fixed $c \geq 2$, when $p = o(n^{\min\{1/2, 1-2/l-2r\}})$ as $n \rightarrow \infty$,

$$\|\hat{\Lambda}_p - \Lambda_p\| = O_p(pn^{-1/2})$$

and

$$\|(\hat{\Sigma}_x - \hat{\Lambda}_p) - (\Sigma_x - \Lambda_p)\| = O_p(pn^{-1/2}).$$

When $p = o(n^{1/4}/(\log n))$, the foregoing two terms almost surely have the convergence rate $o(pn^{-1/4}(\log n))$.

Remark 1. Condition (A) is almost necessary for obtaining the rate of weak convergence with root n . When p is fixed, the existence of a fourth moment is a necessary condition (see Petrov 1995). Condition (B) assumes the smoothness of the inverse regression $\mathbf{m}(y)$. It is slightly stronger than that of Zhu and Ng (1995) when p is fixed; therefore, it is quite mild. As for conditions (C) and (D), they are special for our problem. In the case where p is fixed and the distribution F_n of y_i is also independent of n , conditions (C) and (D) are reduced to those of Zhu and Ng (1995). Hence the conditions are needed to uniformly bind the distribution series $F_n(\cdot)$.

Remark 2. For high-dimension p , it is a natural concern whether $\hat{\Sigma}_x$ will be ill-conditioned, because the inverse matrix $\hat{\Sigma}_x^{-1}$ is needed when standardized variable \mathbf{z} is used to construct the estimate of the CDR space. It is noteworthy that for a fixed p , when n is large, $\hat{\Sigma}_x^{-1}$ cannot be ill-conditioned if Σ_x is of full rank. **When p is large, \mathbf{x} is normally distributed,** and Σ_x is of full rank, the smallest and largest eigenvalues of $\Sigma_x^{-1/2} \hat{\Sigma}_x \Sigma_x^{-1/2} - \mathbf{I}_p$ converge to 0 almost surely as $n \rightarrow \infty$ and $p/n \rightarrow 0$. That is, for large p , once n is large enough, $\hat{\Sigma}_x$ and $\hat{\Sigma}_x^{-1}$ are also of full rank. Thus the ill-conditioning problem can also be avoided (see Bai 1999). **In our case, $p = o(n^{\min\{1/2, 1-2/l-2r\}})$ satisfies this requirement.**

Remark 3. In the cases where all moments of \mathbf{x} exist and the total variation of order r is small, $p = o(n^{\min\{1/2, 1-2/l-2r\}})$ can be close to $O(n^{1/2})$. It is noteworthy that $p = O(n^{1/2})$ is the fastest possible rate if we do not assume extra conditions. This is because each element of $\hat{\Lambda}_p - \Lambda_p$ has an optimal rate $n^{-1/2}$ and $\|\hat{\Lambda}_p - \Lambda_p\|$ is the root of the squared sum of p^2 elements. The resulting rate of $\|\hat{\Lambda}_p - \Lambda_p\|$ must have a factor p . But some data, such as DNA microarray data, have even higher dimension than the sample size. In this case the asymptotic behavior of the estimate is very difficult to study unless $\hat{\Lambda}_p - \Lambda_p$ has a special structure. This is an important issue that merits further study.

Remark 4. The choice of the number, c , of data points in each slice is of practical importance. When p is fixed, Zhu and Ng (1995) derived that for a wide range of c , the convergence of the estimate can be achieved. When p varies with n , we have not obtained a similar result for the case with $c \rightarrow \infty$ and $c/n \rightarrow 0$ as $n \rightarrow \infty$ due to technical difficulties. However, we guess that if the convergence rate of c is much slower than that of p , then the consistency may still hold. We will consider this problem in a future study. When c is proportional to n (i.e., the number of slices is fixed), the CDR space can still be estimated, because the space that is determined by an approximation of $\mathbf{\Sigma}_x - \mathbf{\Lambda}_p$ is a subspace contained in the CDR space although the estimator is not consistent with the matrix $\mathbf{\Sigma}_x - \mathbf{\Lambda}_p$. In this case the asymptotics are much easier to study. The basic idea is as follows. In each slice there are n/H data points where H is a fixed number. Hence we can immediately use central limit theorems to a H -dimensional vector, the components of which are the sample covariance matrices based on the data in each slice. This is because each sample covariance matrix can be rewritten asymptotically as a sum of iid random variables; see the relevant work of Li and Zhu (2004).

Remark 5. By the known result of strong convergence about the sum of independent random variables (see Petrov 1995, thm. 9.20, p. 278), we can easily derive that almost surely $\|\mathbf{\Sigma}_x - \widehat{\mathbf{\Sigma}}_x\| = O(pn^{-1/2} \log n)$. The proof is given in the Appendix. Similar to Corollary 1, we can also derive the convergence of the eigenvalues of $\widehat{\mathbf{\Sigma}}_x$. Let $\{\gamma_i\}$ and $\{\widehat{\gamma}_i\}$ be the eigenvalues of $\mathbf{\Sigma}_x$ and $\widehat{\mathbf{\Sigma}}_x$. Under the condition that $\lim_{p \rightarrow \infty} \min_{1 \leq i \leq p} \{\gamma_i\} > b > 0$, we also have that when n is large, $\lim_{p \rightarrow \infty} \min_{1 \leq i \leq p} \{\widehat{\gamma}_i\} > b/2 > 0$. Note that

$$\widehat{\mathbf{\Sigma}}_x^{-1} - \mathbf{\Sigma}_x^{-1} = \mathbf{\Sigma}_x^{-1}(\mathbf{\Sigma}_x - \widehat{\mathbf{\Sigma}}_x)\widehat{\mathbf{\Sigma}}_x^{-1}.$$

By Lemma A.2 in the Appendix, the symmetry of both $\widehat{\mathbf{\Sigma}}_x$ and $\mathbf{\Sigma}_x^{-1}$ and Corollary 1, it is not difficult to show that

$$\begin{aligned} \|\mathbf{\Sigma}_x^{-1} - \widehat{\mathbf{\Sigma}}_x^{-1}\| &\leq \frac{1}{(\min_{1 \leq i \leq p} \{\gamma_i\} \min_{1 \leq i \leq p} \{\widehat{\gamma}_i\})} \|\mathbf{\Sigma}_x - \widehat{\mathbf{\Sigma}}_x\| \\ &= O(pn^{-1/2} \log n) \quad \text{a.s.} \end{aligned}$$

Together with Theorem 1 and Corollary 1, we can obtain the same convergence rates as those of Theorem 1 for the estimate of the basis, $\{\mathbf{\Sigma}_x^{-1}\boldsymbol{\eta}_1, \dots, \mathbf{\Sigma}_x^{-1}\boldsymbol{\eta}_K\}$. Hence, when we use the standardized variable $\mathbf{z} = \widehat{\mathbf{\Sigma}}_x^{-1/2}(\mathbf{x} - E(\mathbf{x}))$, the asymptotic results of Theorem 1 and Corollary 1 are also similar.

Consider the convergence of the eigenvectors associated with the largest K eigenvalues of $\widehat{\mathbf{\Sigma}}_x - \widehat{\mathbf{\Lambda}}_p$ and the projection spaces of these eigenvectors. Let $\lambda_{p1} \geq \dots \geq \lambda_{pp}$ denote the eigenvalues of $\mathbf{\Sigma}_x - \mathbf{\Lambda}_p$ and $\widehat{\lambda}_{p1} \geq \dots \geq \widehat{\lambda}_{pp}$ denote the eigenvalues of $\widehat{\mathbf{\Sigma}}_x - \widehat{\mathbf{\Lambda}}_p$. Furthermore, let $P_{\lambda_{pi}}$ and $P_{\widehat{\lambda}_{pi}}$ be the projection spaces associated with λ_{pi} and $\widehat{\lambda}_{pi}$. Note that if an eigenvalue is distinct from others, then the corresponding projection space is the space spanned by the associated eigenvector. Theorem 1 implies the following corollary.

Corollary 1. Assume that the conditions of Theorem 1 hold. For each $i = 1, \dots, K$, $|\lambda_{pi} - \widehat{\lambda}_{pi}|$ converges in probability to 0 at the same rates as in Theorem 1. In addition, assume that $\lambda_{p1} \geq$

$\dots \geq \lambda_{pK}$ are distinct in the sense that

$$\liminf_{n \rightarrow \infty} |\lambda_{p(i-1)} - \lambda_{pi}| > 0, \quad i = 2, \dots, K.$$

Then, for $i = 1, \dots, K$, $\|P_{\lambda_{pi}} - P_{\widehat{\lambda}_{pi}}\|$ also converge in probability to 0 at the same rate as in Theorem 1.

4. DETERMINATION OF THE DIMENSION OF $\mathcal{S}_{y|x}$

Differing from the methods in the literature, we suggest a procedure of the BIC type for determining the dimension. It can be used to handle the problems with high-dimensional covariates. The procedure is easy to implement and the consistency of the estimate is established.

Zhao, Krishnaiah, and Bai (1986a,b) studied a problem of detecting the number of signals from noise. By analyzing the eigenvalues of the covariance matrix of the sample, they proposed a BIC-type model selection procedure to determine the number of signals and proved the consistency of the estimate. There is some similarity between their problem and the determination of the dimension of $\mathcal{S}_{y|x}$. We now borrow their idea to construct a determination criterion.

Let $\Psi = \mathbf{\Sigma}_x - \mathbf{\Lambda}_p = \text{cov}\{E(\mathbf{x}|y)\}$ and $\widehat{\Psi} = \widehat{\mathbf{\Sigma}}_x - \widehat{\mathbf{\Lambda}}_p$. Because the smallest $p - K$ eigenvalues of Ψ are 0, we can consider the largest K eigenvalues the signals and the value of K the number of signals. To apply the method used by Zhao et al. (1986a,b), we artificially add an identity matrix \mathbf{I}_p , which acts as a covariance matrix of white noise to Ψ and $\widehat{\Psi}$.

Let $\Omega = \Psi + \mathbf{I}_p$ and $\widehat{\Omega} = \widehat{\Psi} + \mathbf{I}_p$. Let $\theta_1 \geq \theta_2 \geq \dots \geq \theta_p$ be the eigenvalues of Ω and $\widehat{\theta}_1 \geq \widehat{\theta}_2 \geq \dots \geq \widehat{\theta}_p$ be the eigenvalues of $\widehat{\Omega}$. It is clear that $\theta_i = \lambda_i + 1$, where λ_i are the eigenvalues of Ψ . Determining the dimension of $\mathcal{S}_{y|x}$ now becomes estimating K , the number of the eigenvalues of $\Omega > 1$.

Borrowing the idea of maximum likelihood estimation in the normal distribution case, we define

$$\log L(\boldsymbol{\theta}) = -\frac{n}{2} \log |\Omega| - \frac{n}{2} \text{tr} \Omega^{-1} \widehat{\Omega}, \quad (8)$$

because it is a function of $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)$. Let Θ_k be the set consisting of all values such that $\theta_1 \geq \theta_2 \geq \dots \geq \theta_k > 1$ and $\theta_{k+1} = \dots = \theta_p = 1$. In addition, let τ denote the number of $\widehat{\theta}_i$'s that are > 1 . According to Zhao et al. (1986a,b), we can have a specific form of $\sup_{\boldsymbol{\theta} \in \Theta_k} \log L(\boldsymbol{\theta})$,

$$\begin{aligned} \sup_{\boldsymbol{\theta} \in \Theta_k} \log L(\boldsymbol{\theta}) &= -\frac{n}{2} \sum_{i=1}^p \log \widehat{\theta}_i - \frac{np}{2} \\ &\quad + \frac{n}{2} \sum_{i=1+\min(\tau, k)}^p (\log \widehat{\theta}_i + 1 - \widehat{\theta}_i). \end{aligned} \quad (9)$$

Note that this supreme does not involve the unknowns relating to the matrix Ω and its eigenvalues, θ_i 's. For defining the estimator \widehat{K} of the true dimension K , (9) is equivalent to

$$\frac{n}{2} \sum_{i=1+\min(\tau, k)}^p (\log \widehat{\theta}_i + 1 - \widehat{\theta}_i).$$

and thus we can define a procedure based on it. A criterion of BIC type is defined as follows. Let

$$\log L_k = \frac{n}{2} \sum_{i=1+\min(\tau, k)}^p (\log \widehat{\theta}_i + 1 - \widehat{\theta}_i) \quad (10)$$

and

$$G(k) = \log L_k - \frac{C_n k(2p - k + 1)}{2}, \quad (11)$$

where the second term is the penalty term, C_n is a penalty constant, and $k(2p - k + 1)/2$ equals the number of free parameters of (8) needed to estimate when $\theta \in \Theta_k$. It can be justified as follows. When k is fixed, the SIR matrix Ψ is $p \times p$ being of rank k . Hence the last $p - k$ columns of this matrix are the linear combinations of the first k columns. The total number of the parameters in these k columns is kp . Furthermore, because of the symmetry of the matrix, in the left-upper $k \times k$ submatrix $k(k - 1)/2$ elements are the same as those above the diagonal. Hence the total number of the free parameters is $kp - k(k - 1)/2 = k(2p - k + 1)/2$.

The estimator of K is defined as the maximizer \hat{K} of $G(k)$ over $k \in \{0, \dots, p - 1\}$, that is,

$$G(\hat{K}) = \max_{0 \leq k \leq p-1} G(k). \quad (12)$$

The following result states the weak and strong consistency of \hat{K} .

Theorem 2. Assume that $p = O(n^s)$, K is a constant independent of n , $\|\Omega - \hat{\Omega}\| = O_p(n^{-t})$ [or $O(n^{-t})$ a.s.], $t > 0$, $2t > s$, and C_n satisfies the following:

- (a) $\lim_{n \rightarrow +\infty} C_n/n^{1-s} = 0$
- (b) $\lim_{n \rightarrow +\infty} C_n/n^{1-2t} = \infty$.

Then $\hat{K} - K = o_p(1)$ or $o(1)$ a.s.

Remark 6. From the proof of Theorem 1, we know that there are s and t satisfying the designed conditions of Theorem 2. If $p = O(n^s)$ and $C_n = O(n^a)$, then, from the results of Theorem 1, we derive that for any small $\eta > 0$,

$$\|\Omega - \hat{\Omega}\| = O_p(n^{2s-1/2}) \text{ or } O(n^{s+\eta-1/4}) \text{ a.s.}$$

Thus, from Theorem 2, the power a of C_n can be chosen within the range $1 - 2t = 4s < a < 1 - s$ for the weak convergence and $1/2 + 2s < a < 1 - s$ for the strong convergence. If p is fixed, then $s = 0$, and we can select a in the interval $[0, 1)$ and $[1/2, 1)$ for the weak and strong convergence. Therefore, the range is fairly wide.

Remark 7. In our procedure, we artificially add a covariance matrix, \mathbf{I}_p , to the matrix Ψ . There is an alternative to constructing Ω . Note that $\Sigma_x = \Psi + \Lambda_p$. Then, if Λ_p is nonsingular, we can replace the previous Ω and $\hat{\Omega}$ by

$$\Omega = \Sigma_x \Lambda_p^{-1} = \Psi \Lambda_p^{-1} + \mathbf{I}_p \quad \text{and}$$

$$\hat{\Omega} = \hat{\Sigma}_x \hat{\Lambda}_p^{-1} = \hat{\Psi} \hat{\Lambda}_p^{-1} + \mathbf{I}_p.$$

Remark 8. Note that the BIC-type method consistently estimates the dimension of the CDR space under the condition that the estimator of SIR matrix is consistent at some rate. The limit distribution is not necessary. **Therefore, the method is easy to use and is readily extended to other dimension-reduction methods, such as SAVE and MAVE. In contrast, the popular sequential test method can be easily used when the limit distribution of the estimator of the SIR matrix is available and its limiting variance is easy to estimate;** otherwise, we must use a Monte Carlo approximation to simulate its distribution. Permutation

(e.g., Cook and Weisberg 1991) or the bootstrap may be applicable if the consistency of these approximations can be verified.

5. SIMULATION STUDY

To investigate the performance of estimation and dimensionality determination by the BIC-type method, we carried out a set of simulations. In this section we report part of the results of these simulations. We used five models. To measure the distance between the CDR space and its estimator, we first obtained the eigenvectors $\{\hat{\beta}_1, \dots, \hat{\beta}_K\}$ associated with the first K largest eigenvalues of $\hat{\Sigma}_x - \hat{\Lambda}_p$, then multiplied these vectors from the left by $\hat{\Sigma}_x^{-1}$ to form an estimator \mathcal{B} of $\mathcal{S}_{y|x}$. As suggested by Li (1991), for any i , we use the squared multiple correlation coefficient $R^2(\hat{\beta}_i)$ to measure the distance between $\hat{\beta}_i^T \mathbf{x}$ and \mathcal{B} , where

$$R^2(\hat{\beta}_i) = \max_{\beta \in \mathcal{S}_{y|x}} \frac{(\hat{\beta}_i \Sigma_x \beta^T)^2}{\hat{\beta}_i \Sigma_x \hat{\beta}_i^T \cdot \beta \Sigma_x \beta^T}, \quad (13)$$

and use the average of $R^2(\hat{\beta}_i)$'s to measure the distance between \mathcal{B} and $\mathcal{S}_{y|x}$. We can also use the squared trace correlation, the average of the squared canonical correlation coefficients between $\hat{\beta}_1^T \mathbf{x}, \dots, \hat{\beta}_K^T \mathbf{x}$ and $\beta_1^T \mathbf{x}, \dots, \beta_K^T \mathbf{x}$, denoted by $R^2(\hat{\mathcal{B}})$, as our criterion (see also Hooper 1959). These two criteria are similar, but in the simulation we report all values of $R^2(\hat{\beta}_i)$'s, so that the performance of the estimation can be described more clearly.

Example 1. Consider the model

$$y = x_1 \cdot (x_2 + x_3 + 1) + \varepsilon, \quad (14)$$

where ε is normally distributed with mean 0 and variance σ^2 . In the simulation, $\sigma = .5$. The variables x_1, x_2 , and x_3 are the components of $\mathbf{x} = \{x_1, \dots, x_p\}$ ($p \geq 3$) and are independent of ε . \mathbf{x} is assumed to have a multivariate normal distribution $N(\mathbf{0}, \mathbf{I}_p)$, and the CDR space is spanned by two vectors, $\beta_1 = (1, 0, \dots, 0)$ and $\beta_2 = (0, 1, 1, \dots, 0)$.

We report the simulation results with $n = 400$ and 800 in Table 1. To examine the sensitivity of the estimator to the selection of the number of slices, we considered $c = 10, 20, 40$, and 80 . All of the reported results are the averages of 500 independent replicates, and the values in parentheses are the standard deviations.

Looking at Table 1, we see that when n is fixed and p gets larger, the values of $R^2(\beta_i)$'s get smaller. Clearly, dimensionality has an impact on the performance of the estimation. As expected, when n gets larger, the estimation improves. Comparing the case with $(p = 10, n = 400)$ with the case with $(p = 20, n = 800)$ shows that the magnitudes of R^2 are similar, as does comparing the case $(p = 20, n = 400)$ and the case $(p = 40, n = 800)$ also shows this. Thus these examples indicate that the ratio p/n plays an important role in the convergence of the estimator (see the discussion in Remark 2).

For selecting the number of slices, which is equivalent to selecting the number c of data points in each slice, we found that generally the range of the slice numbers for selection is fairly wide. As Zhu and Ng (1995) proved, the root- n consistency

Table 1. Estimation of the CDR Space With Model (14)

p	$c = 10$		$c = 20$		$c = 40$		$c = 80$	
	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$
$n = 400$								
10	.90 (.07)	.71 (.22)	.91 (.05)	.78 (.16)	.93 (.04)	.80 (.10)	.92 (.05)	.76 (.14)
20	.78 (.12)	.51 (.21)	.84 (.06)	.61 (.18)	.83 (.07)	.66 (.13)	.82 (.07)	.59 (.15)
30	.69 (.14)	.35 (.21)	.78 (.08)	.52 (.14)	.75 (.09)	.53 (.16)	.74 (.09)	.49 (.14)
40	.60 (.13)	.29 (.17)	.71 (.10)	.42 (.16)	.72 (.09)	.43 (.16)	.70 (.9)	.39 (.13)
50	.53 (.17)	.26 (.18)	.61 (.11)	.32 (.16)	.63 (.10)	.34 (.15)	.63 (.10)	.33 (.12)
$n = 800$								
10	.96 (.02)	.87 (.09)	.96 (.02)	.91 (.04)	.96 (.02)	.92 (.05)	.97 (.02)	.92 (.04)
20	.90 (.04)	.72 (.14)	.92 (.04)	.81 (.10)	.92 (.03)	.83 (.06)	.93 (.03)	.83 (.06)
30	.85 (.05)	.62 (.16)	.88 (.04)	.73 (.11)	.89 (.04)	.74 (.08)	.88 (.04)	.74 (.08)
40	.80 (.06)	.51 (.18)	.83 (.04)	.64 (.12)	.85 (.05)	.69 (.09)	.84 (.5)	.67 (.10)
50	.75 (.09)	.41 (.18)	.80 (.07)	.57 (.13)	.81 (.05)	.62 (.12)	.80 (.05)	.60 (.10)

can be maintained in a rather wide range of slice numbers, although the asymptotic variance is affected. The simulation results given herein support this theoretical conclusion. **Another finding of interest is that the choice of c is not seriously affected by the dimension p . From $p = 10$ –50, the values of R^2 with $c = 20, 40, 80$ are very similar.**

To determine the dimension, we must consider the selection of C_n when we use the BIC proposed in this article. For a different problem, the choice of C_n has been discussed by Bai, Krishnaiah, and Zhao (1989). In our simulations, we considered $C_n = c^{-1}W_n$ where c is the number of data points in each slice. We used the term c^{-1} because Zhu and Ng (1995) have proven that the asymptotic variance of $\hat{\Lambda}_p$ depends on $(c-1)^{-1}$, which then plays the role of a variance σ^2 that is used in the BIC for variable selection (see Eubank 1999). Selecting C_n is then equivalent to selecting W_n . We tried several values. From Theorem 2, W_n can be selected in a fairly wide range. Note that Schwarz (1978) used $\log n$, and from Theorem 2, a rate of n^α for some α satisfies the conditions. Therefore, we started with a rate of $\log n$ and a rate of $n^{1/3}$, and tried several combinations of these two values: $W_n = .1 \log n$, $.5 \log(n)$, $.1n^{1/3}$, $.5n^{1/3}$, $(.5 \log(n) + .1n^{1/3})$, and $(.5 \log(n) + .1n^{1/3})/2$. Table 2 gives the results of $W_n = (.5 \log(n) + .1n^{1/3})/2$, which performs better overall than the other W_n 's.

A competitor in the literature is the now-standard sequential test (ST) method. ST includes the chi-squared test (CST) sug-

gested by Li (1991) when the distribution of X is normal, as well as the weighted chi-squared test (WCST) proposed by Bura and Cook (2001a) when X is not normally distributed. These two methods were compared. It is noteworthy that when $p \rightarrow \infty$, no asymptotic results are available for sequential tests. But in practice both n and p are fixed, and we can still use this method. The proportions of the correct decisions based on the CST reported in Table 2 were obtained in the following way. Let the nominal level be .05. Set the hypothesis $H_0: d = 1$ versus $H_1: d > 1$. We recorded the simulation experiments that reject H_0 . Then set the hypothesis $H_0: d = 2$ versus $H_1: d > 2$. Among the recorded experiments in the previous test, we recorded the decisions that reject the null hypothesis. This procedure proceeds until $H_0: d = K$ versus $H_1: d > K$ for $K \leq p - 1$, where K is the true dimension of the CDR space. The results reported in Table 2 are the proportions of correct decisions at step K of the 500 simulation experiments.

Because X is normally distributed in this example, we use the CST. Interestingly, from Table 2, we can see that a large c (or, equivalently, a small number of slices) favors the CST, as does high dimensions. The results reported in Table 2 based on the BIC do not show this trend. When we choose the respective c that favors its corresponding methodology, the BIC outperforms the CST for $p \leq 30$, whereas the ST works better when $p = 40, 50$.

Table 2. Proportion of Correct Decisions About the Dimension of the CDR Space With Model (14)

p	BIC with $W_n = (.5 \log(n) + .1n^{1/3})/2$								Chi-squared test							
	$n = 400$				$n = 800$				$n = 400$				$n = 800$			
	$c = 10$	$c = 20$	$c = 40$	$c = 80$	$c = 10$	$c = 20$	$c = 40$	$c = 80$	$c = 10$	$c = 20$	$c = 40$	$c = 80$	$c = 10$	$c = 20$	$c = 40$	$c = 80$
10	.78	.92	.93	.94	.94	.99	.99	.98	.39	.71	.82	.82	.70	.94	.94	.95
20	.73	.82	.81	.87	.85	.95	.95	.90	.24	.48	.54	.67	.46	.80	.88	.93
30	.69	.62	.54	.65	.81	.82	.90	.70	.17	.34	.40	.51	.31	.65	.75	.89
40	.43	.26	.20	.32	.63	.76	.79	.51	.11	.26	.34	.39	.27	.54	.70	.86
50	.12	.06	.12	.13	.42	.54	.72	.19	.10	.21	.28	.35	.18	.50	.60	.75

Table 3. Estimation of the CDR Space With Model (15)

<i>n</i>	<i>c</i> = 5		<i>c</i> = 10		<i>c</i> = 20		<i>c</i> = 40	
	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$
Distribution (A)								
100	1.00 (0)	.55 (.28)	1.00 (0)	.61 (.27)	.99 (.01)	.57 (.28)	.99 (.01)	.50 (.29)
200	1.00 (0)	.65 (.28)	1.00 (0)	.75 (.23)	1.00 (0)	.79 (.20)	1.00 (0)	.74 (.22)
300	1.00 (0)	.737 (.25)	1.00 (0)	.82 (.18)	1.00 (0)	.88 (.11)	1.00 (0)	.86 (.16)
400	1.00 (0)	.79 (.21)	1.00 (0)	.87 (.15)	1.00 (0)	.90 (.11)	1.00 (0)	.92 (.08)
Distribution (B)								
100	.98 (.03)	.56 (.30)	.98 (.02)	.64 (.31)	.99 (.01)	.67 (.28)	.98 (.02)	.66 (.30)
200	.99 (.01)	.59 (.32)	.99 (.01)	.65 (.31)	.99 (.01)	.78 (.30)	.99 (.01)	.71 (.22)
300	.99 (.01)	.64 (.32)	1.00 (0)	.73 (.27)	1.00 (0)	.79 (.24)	1.00 (0)	.82 (.19)
400	1.00 (.01)	.67 (.31)	1.00 (0)	.76 (.27)	1.00 (0)	.82 (.23)	1.00 (0)	.87 (.17)

Example 2. In this example we consider a model used by Velilla (1998),

$$y = (4 + x_1)(x_2 + x_3 + 2) + \varepsilon, \quad (15)$$

where ε has a normal distribution with mean 0 and variance σ^2 with $\sigma = .5$ in the simulations. The variables x_1, x_2 , and x_3 are the components of $\mathbf{x} = \{x_1, \dots, x_5\}$ and are independent of ε . This model is similar to (14) but has different distributions of \mathbf{x} , as follows:

(A) \mathbf{x} has a normal distribution $N_5(\mathbf{0}, \Sigma)$, where $\Sigma = \text{diag}(2, 2, 2, 4, 2)$.

(B) \mathbf{x} is obtained by

$$\mathbf{x} = \mathbf{C}\mathbf{v} + \mathbf{B}(\mathbf{B}'\mathbf{B})^{-1}\mathbf{w},$$

where

$$\mathbf{C} = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ -1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & -1 \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}.$$

Here $\mathbf{v} = (v_1, v_2, v_3)'$ consists of three iid uniform random variables on $(-4, 4)$, and $\mathbf{w} = (w_1, w_2)'$ has two independent components, of which w_1 is distributed with a mixture $.5N(0, 4) + .5N(0, 16)$ of two normal distributions and w_2 has the $U(-4, 4)$ distribution, where $U(-4, 4)$ is the uniform distribution on $(-4, 4)$. Velilla (1998) proved that \mathbf{x} in case (B) also satisfies the linearity condition designed by Li (1991), although it is not elliptically distributed.

Because $p = 5$ is relatively small, we can consider small sample sizes. In the simulation, the sample size was $n = 100, 200, 300, 400$. Tables 3 and 4 report the results with model (15). Case (A) is used to ascertain whether difference between the variances of the components of \mathbf{x} affects the estimation efficiency. Looking at Table 3, we find that in a normal case, the effect is not significant. For the selection of c , the pairs $(n = 100, c = 10)$, $(n = 200, 300, c = 20)$, and $(n = 400, c = 40)$ give large values of R^2 . It seems that, overall, $H = 10$ may work well. For case (B), SIR still works, although its performance is not as good as that of case (A) with normal covariates. Again, the estimation is not very sensitive to the choice of the slice number; $c = 20$ works well, and $H = 10$ may still be considered as a working value for practical use.

For the determination of dimensionality, Table 4 reports the results of the BIC and the CST for (A) and WCST for (B) with $c = 5, 10, 20$. The normal distribution of \mathbf{x} clearly favors the two methods, but for the very nonnormal case (B), estimation for K is more difficult. A comparison of these two methods indicates that the BIC is clearly better than the CST/WCST regardless of whether or not the distribution of \mathbf{x} is normal.

Example 3. In this example we consider a model in which the CDR space is spanned by $\beta_1 = (1, 0, \dots, 0)^T$ in mean and $\beta_2 = (0, 1, 1, \dots, 0)^T$ in variance as

$$y = (4 + x_1) + (x_2 + x_3 + 2) \cdot \varepsilon. \quad (16)$$

Here ε has a normal distribution with mean 0 and variance σ^2 with $\sigma = .5$ in the simulation. The variables x_1, \dots, x_3 are

Table 4. Proportion of Correct Decisions About the Dimension of the CDR Space With Model (15)

<i>n</i>	<i>BIC with $W_n = (.5 \log(n) + .1n^{1/3})/2$</i>						<i>Chi-squared test</i>			<i>Weighted chi-squared test</i>		
	<i>Distribution (A)</i>			<i>Distribution (B)</i>			<i>Distribution (A)</i>			<i>Distribution (B)</i>		
	<i>c</i> = 5	<i>c</i> = 10	<i>c</i> = 20	<i>c</i> = 5	<i>c</i> = 10	<i>c</i> = 20	<i>c</i> = 5	<i>c</i> = 10	<i>c</i> = 20	<i>c</i> = 5	<i>c</i> = 10	<i>c</i> = 20
100	.53	.57	.60	.33	.39	.34	.17	.31	.36	.24	.32	.37
200	.68	.81	.88	.44	.41	.53	.33	.57	.75	.28	.36	.50
300	.79	.92	.95	.46	.51	.68	.44	.73	.90	.30	.50	.58
400	.88	.95	.98	.53	.63	.75	.57	.81	.93	.38	.55	.65

Table 5. Estimation of the CDR Space With Model (16)

<i>n</i>	<i>c</i> = 5		<i>c</i> = 10		<i>c</i> = 20		<i>c</i> = 40	
	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$
Distribution (A)								
100	.95 (.06)	.51 (.29)	.97 (.03)	.53 (.29)	.97 (.03)	.50 (.30)	.95 (.04)	.33 (.25)
200	.98 (.02)	.52 (.31)	.98 (.03)	.61 (.30)	.99 (.01)	.64 (.28)	.98 (.01)	.64 (.27)
300	.99 (.01)	.61 (.28)	.99 (.01)	.70 (.28)	.99 (.01)	.76 (.23)	.99 (.01)	.72 (.30)
400	.99 (.01)	.63 (.30)	.99 (.01)	.76 (.23)	.99 (.01)	.81 (.21)	.99 (.01)	.83 (.17)
Distribution (B)								
100	.99 (.01)	.30 (.29)	.99 (.01)	.32 (.29)	.99 (.01)	.40 (.31)	.97 (.03)	.30 (.26)
200	.99 (.01)	.30 (.28)	.99 (.01)	.36 (.30)	.99 (.01)	.41 (.32)	.99 (.01)	.31 (.33)
300	1.00 (0)	.30 (.28)	.99 (0)	.36 (.31)	.99 (0)	.43 (.32)	.99 (0)	.41 (.27)
400	1.00 (0)	.34 (.29)	1.00 (0)	.36 (.31)	1.00 (0)	.46 (.32)	1.00 (0)	.48 (.33)

the components of $\mathbf{x} = \{x_1, \dots, x_p\}$ ($p = 5$) and are independent of ε , and \mathbf{x} has the same distributions (A) and (B) as in model (15). For this model, the situation is similar to that of Example 2; the methods work better in the normal case (A). As to the selection of c , we come to a similar conclusion that SIR is not sensitive to the number of slices and $c = 20$ works well with the sample size that we used. In addition, $H = 10$ may be a good choice for practical use.

For determining dimensionality, we considered $c = 5, 10, 20$ when the sample size was ≤ 300 . Looking at Table 6, we clearly see that the dimension is difficult to estimate, especially in the nonnormal case (B). This is because, in contrast to the model of Example 2, one dimension is included in the variance. Hence this model does not favor the SIR. Therefore, we considered a larger sample size to explore how large a sample size is needed to obtain a good estimator; $n = 400, 800, 1,600, 3,200$ were tested. For sequential test method, similar to those of Example 2, in the normal case (A) we use the CST and in the nonnormal case (B) we use the WCST. In a normal case, the estimation will be satisfactory when $n \geq 800$; however, for the nonnormal case (B), $n = 3,200$ is still not large enough. This observation demonstrates that we may need to use other means to obtain a better structure of data such as that suggested by Cook and Nachtsheim (1994). Between the BIC and the CST/WCST,

the former generally performs better, although in case (B) the WCST performs slightly better when the sample size is large.

A comparison of models (15) and (16) provides evidence that SIR has difficulty estimating the CDR space when basis is in the variance part.

Example 4. Finally, we consider a model with a three-dimensional CDR space,

$$y = x_1(x_2 + x_3 + 2) + (x_4 + x_5 + 2)^3 + \varepsilon, \quad (17)$$

where ε is normally distributed with mean 0 and variance σ^2 . In the simulation, $\sigma = .5$. The variables x_1, x_2, x_3, x_4 , and x_5 are the components of $\mathbf{x} = \{x_1, \dots, x_p\}$ ($p \geq 5$) and are independent of ε . $\mathbf{x} = (x_1, x_2, \dots, x_p)$ is normally distributed with $N_p(\mathbf{0}, \mathbf{I}_p)$. The CDR space is spanned by $\beta_1 = (1, 0, \dots, 0)^T$, $\beta_2 = (0, 1, 1, \dots, 0)^T$, and $\beta_3 = (0, 0, 0, 1, 1, \dots, 0)^T$.

Based on our experience, we believe that a small sample size cannot provide a good estimation of the CDR space of this model. Here we report the results with $n = 400$ and $n = 3,200$ to demonstrate the difficulty of estimating the CDR space and the large sample size needed to obtain an acceptable estimation.

Table 7 shows that the $R^2(\hat{\beta}_3)$'s are small, even when $p = 10$. When n is increased to 3,200, the $R^2(\hat{\beta}_3)$'s become acceptably large when $p \leq 20$. Increasing the dimension p significantly reduces the performance of the estimation. Clearly, for larger p ,

Table 6. Proportion of Correct Decisions About the Dimension of the CDR Space With Model (16)

BIC with $W_n = (.5 \log(n) + .1n^{1/3})/2$							Chi-squared test			Weighted chi-squared test			
n	Distribution (A)			Distribution (B)			n	Distribution (A)			Distribution (B)		
	$c = 5$	$c = 10$	$c = 20$	$c = 5$	$c = 10$	$c = 20$		$c = 5$	$c = 10$	$c = 20$	$c = 5$	$c = 10$	$c = 20$
100	.28	.28	.23	.20	.16	.13	100	.08	.15	.16	.02	.06	.09
200	.29	.43	.48	.13	.13	.13	200	.12	.24	.38	.04	.08	.08
300	.34	.49	.61	.09	.11	.13	300	.21	.34	.51	.04	.05	.10
n	Distribution (A)			Distribution (B)			n	Distribution (A)			Distribution (B)		
	$c = 20$	$c = 40$	$c = 80$	$c = 20$	$c = 40$	$c = 80$		$c = 20$	$c = 40$	$c = 80$	$c = 20$	$c = 40$	$c = 80$
400	.72	.79	.69	.11	.17	.19	400	.60	.71	.63	.09	.16	.19
800	.92	.98	.99	.11	.20	.28	800	.87	.91	.93	.14	.24	.31
1,600	.99	1.00	.99	.12	.27	.47	1,600	.95	.95	.94	.20	.33	.52
3,200	.99	1.00	1.00	.15	.38	.76	3,200	.96	.93	.96	.32	.52	.77

Table 7. Estimation of the CDR Space With Model (17)

p	$c = 10$			$c = 20$			$c = 40$			$c = 80$		
	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_3)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_3)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_3)$	$R^2(\hat{\beta}_1)$	$R^2(\hat{\beta}_2)$	$R^2(\hat{\beta}_3)$
$n = 400$												
10	1.00 (.001)	.70 (.19)	.30 (.22)	1.00 (.001)	.77 (.15)	.33 (.23)	1.00 (.001)	.80 (.13)	.33 (.23)	1.00 (.002)	.74 (.14)	.24 (.18)
20	.99 (.002)	.46 (.22)	.20 (.16)	1.00 (.002)	.56 (.19)	.21 (.15)	.99 (.002)	.60 (.17)	.21 (.16)	.99 (.003)	.53 (.06)	.14 (.12)
30	.99 (.003)	.30 (.18)	.16 (.13)	.99 (.003)	.40 (.18)	.16 (.13)	.99 (.003)	.47 (.17)	.16 (.13)	.99 (.004)	.41 (.14)	.12 (.11)
40	.99 (.004)	.23 (.14)	.13 (.11)	.99 (.003)	.31 (.17)	.13 (.11)	.99 (.004)	.37 (.15)	.13 (.10)	.98 (.005)	.31 (.14)	.10 (.09)
$n = 3,200$												
10	1.00 (.0003)	.94 (.03)	.47 (.26)	1.00 (.0003)	.95 (.03)	.60 (.25)	1.00 (.0003)	.95 (.02)	.70 (.20)	1.00 (.0003)	.95 (.03)	.69 (.19)
20	1.00 (.0004)	.87 (.05)	.26 (.19)	1.00 (.0004)	.89 (.04)	.37 (.21)	1.00 (.0004)	.90 (.04)	.43 (.21)	1.00 (.0005)	.89 (.04)	.47 (.19)
30	1.00 (.0005)	.79 (.07)	.17 (.14)	1.00 (.0005)	.83 (.05)	.22 (.17)	1.00 (.0005)	.85 (.04)	.31 (.18)	1.00 (.0006)	.84 (.05)	.36 (.16)
40	1.00 (.0006)	.72 (.09)	.12 (.11)	1.00 (.0005)	.78 (.06)	.18 (.14)	1.00 (.0006)	.79 (.05)	.24 (.15)	1.00 (.0007)	.79 (.05)	.29 (.14)

we should use larger datasets to obtain a reasonable estimation. The selection of c seems similar to that of the previous examples. Within a large range of c (from 20 to 80), the performance of SIR is similar.

The determination of the dimension is also difficult. To provide insights into how the dimension of the CDR space can be estimated, in Table 8 we report the proportions of the estimated dimensions of the CDR space in the 500 replications. Neither the BIC nor the CST works for $n = 400$ (see Table 8); a large sample size is needed. We also report the simulation results with $n = 3,200$. Table 8 clearly shows that when the sample size is equal to 400, the BIC tends to estimate K as 2 and the CST seems to estimate K even smaller. This becomes clearer when the dimension p of X is large. When $n = 3,200$, both methods can obtain a high proportion of correct decisions, and the BIC outperforms the CST, especially when p is large. Another interesting observation is that, in contrast to the performance of estimation with model (14), the CST gets worse when p gets larger.

In summary, estimation of the SIR matrix is not very sensitive to the choice of c or, equivalently, to the number of slices. The choice of c is not greatly affected by the dimension p , and thus we can choose the number of slices regardless of the dimension of the covariates. Based on our simulations, the number of slices chosen can fall between 10 and 20. Compared with the sequential test method for determining the dimension of the CDR space, the advantages of the BIC are that it does not rely on the limit distribution of the estimator of the SIR matrix and that it

can be readily applied to other dimension-reduction methods, such as the SAVE and the MAVE. However, the need to select C_n is a disadvantage. Whether there is a data-driven selection for C_n deserves further study.

APPENDIX: PROOFS OF THEOREMS

A.1 Notation

For the sake of convenience, we have assumed that n/c is an integer. Here we define some notation. Recalling (4) and (5), we define $\mathbf{Y}_n = (y_{(1)}, \dots, y_{(n)})^T$, $\mathbf{X}_n = (\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(n)})_{p \times n}$ as a $p \times n$ matrix, $\mathbf{X}_h = (\mathbf{x}_{(h,1)}, \dots, \mathbf{x}_{(h,c)})_{p \times c}$, $h = 1, 2, \dots, H$, as $p \times c$ matrices, and for any integer t , $\mathbf{e}_t = (1, \dots, 1)^T$ as a $t \times 1$ column vector and $\mathcal{E}_t = (\mathbf{e}_t, \dots, \mathbf{e}_t)_{t \times p}^T$ as a $t \times p$ matrix. We further define more matrices. Set \mathbf{I}_c to be the $c \times c$ identity matrix and

$$\mathbf{P}_c = \mathbf{I}_c - \frac{1}{c} \mathbf{e}_c \mathbf{e}_c^T, \quad \mathbf{Q}_c = \mathbf{P}_c - \left(1 - \frac{1}{c}\right) \mathbf{I}_c = \frac{1}{c} (\mathbf{I}_c - \mathbf{e}_c \mathbf{e}_c^T);$$

$$\mathbf{M}_n = (\mathbf{m}(y_{(1)}), \dots, \mathbf{m}(y_{(n)}))_{p \times n},$$

$$\mathbf{M}_h = (\mathbf{m}(y_{(h,1)}), \dots, \mathbf{m}(y_{(h,c)}))_{p \times c};$$

$$\boldsymbol{\varepsilon}_n = (\boldsymbol{\varepsilon}_{(1)}, \dots, \boldsymbol{\varepsilon}_{(n)})_{p \times n}, \quad \boldsymbol{\varepsilon}_{h,c} = (\boldsymbol{\varepsilon}_{(h,1)}, \dots, \boldsymbol{\varepsilon}_{(h,c)})_{p \times c},$$

and

$$\tilde{\mathbf{M}}_n = [\{\mathbf{m}(y_{(1)}) - \mathbf{m}(y_{(2)})\}, \dots, \{\mathbf{m}(y_{(n-1)}) - \mathbf{m}(y_{(n)})\}, \mathbf{m}(y_{(n)})]_{p \times n},$$

Table 8. Proportions of the Estimated Dimensions of the CDR Space With Model (17), $H = 10$

p	BIC with $W_n = (.5 \log(n) + .1n^{1/3})/2$								Chi-squared test							
	$n = 400$				$n = 3,200$				$n = 400$				$n = 3,200$			
	10	20	30	40	10	20	30	40	10	20	30	40	10	20	30	40
$K = 1$.05	.07	.05	.04	0	0	0	.04	.14	.35	.44	.58	0	0	0	0
$K = 2$.86	.84	.85	.85	.02	.05	.07	0	.75	.59	.51	.40	.02	.15	.30	.39
$K = 3$.09	.09	.10	.11	.98	.95	.91	.89	.11	.06	.05	.02	.94	.83	.67	.59
$K = 4$	0	0	0	0	0	0	.02	.07	0	0	0	0	.04	.02	.03	.02

NOTE: The true dimension is three.

where the $\varepsilon_{(i)}$'s are as defined in (5) and $\varepsilon_{(h,j)} = \varepsilon_{(c(h-1)+j)}$. Furthermore, for any t ,

$$\mathbf{T}_t = \begin{pmatrix} 1 & -1 & & \\ & 1 & \ddots & \\ & & \ddots & -1 \\ & & & 1 \end{pmatrix}_{t \times t} \quad \text{and}$$

$$\mathbf{S}_t = \mathbf{T}_t^{-1} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ & 1 & \cdots & \vdots \\ & & \ddots & 1 \\ & & & 1 \end{pmatrix}_{t \times t}.$$

From the definitions of $\tilde{\mathbf{M}}_n$ and \mathbf{T}_n , it is clear that $\tilde{\mathbf{M}}_n^T = \mathbf{T}_n \mathbf{M}_n^T$.

Now we can rewrite $\hat{\mathbf{A}}_p$ in terms of the foregoing notation in a matrix form. Because $\mathbf{X}_n = \mathbf{M}_n + \varepsilon_n$,

$$\begin{aligned} \hat{\mathbf{A}}_p &= \frac{1}{H(c-1)} \sum_{h=1}^H \mathbf{X}_h \mathbf{P}_c \mathbf{X}_h^T \\ &= \frac{1}{H(c-1)} \mathbf{X}_n (\mathbf{I}_H \otimes \mathbf{P}_c) \mathbf{X}_n^T \\ &= \frac{1}{H(c-1)} \varepsilon_n (\mathbf{I}_H \otimes \mathbf{P}_c) \varepsilon_n^T \\ &\quad + \frac{1}{H(c-1)} \{ \varepsilon_n (\mathbf{I}_H \otimes \mathbf{P}_c) \mathbf{M}_n^T + \mathbf{M}_n (\mathbf{I}_H \otimes \mathbf{P}_c) \varepsilon_n^T \} \\ &\quad + \frac{1}{H(c-1)} \mathbf{M}_n (\mathbf{I}_H \otimes \mathbf{P}_c) \mathbf{M}_n^T \\ &=: \mathbf{J}_1 + \mathbf{J}_2 + \mathbf{J}_3, \end{aligned} \quad (\text{A.1})$$

where “ \otimes ” represents the Kronecker product. Note that $\mathbf{P}_c + \mathbf{Q}_c = (1 - 1/c)\mathbf{T}_c$. Then $\mathbf{I}_H \otimes \mathbf{P}_c + \mathbf{I}_H \otimes \mathbf{Q}_c = (1 - 1/c)\mathbf{I}_H$. \mathbf{J}_1 is then decomposed as

$$\begin{aligned} \mathbf{J}_1 &= \frac{1}{H(c-1)} \left\{ \left(1 - \frac{1}{c}\right) \varepsilon_n \varepsilon_n^T \right\} + \frac{1}{H(c-1)} \{ \varepsilon_n (\mathbf{I}_H \otimes \mathbf{Q}_c) \varepsilon_n^T \} \\ &=: \mathbf{K}_1 + \mathbf{K}_2. \end{aligned} \quad (\text{A.2})$$

It is easy to see that

$$\|\mathbf{A}_p - \hat{\mathbf{A}}_p\| \leq \|\mathbf{A}_p - \mathbf{K}_1\| + \|\mathbf{K}_2\| + \|\mathbf{J}_2\| + \|\mathbf{J}_3\|. \quad (\text{A.3})$$

We need only derive the convergence rate of the foregoing four terms on the right side of (A.3). For this, we introduce the following lemmas.

A.2 Lemmas

Throughout this and the following sections, b is a constant that may take different values with each appearance.

Lemma A.1. Under the conditions of Theorem 1,

$$\lim_{n \rightarrow \infty} \frac{1}{n^r \sqrt{p}} \sum_{i=2}^n \{ \|\mathbf{m}(y_{(i)}) - \mathbf{m}(y_{(i-1)})\| \} = 0 \quad \text{a.s.}$$

Proof. If all $y_i \in [-B, B]$ for some $B > 0$, then the result follows from condition (B) of Theorem 1. Now we need only consider the case where the support of $y_i^{(n)}$ is unbounded. It suffices to show that the sums

$$S_{n1} = \frac{1}{n^r \sqrt{p}} \sum_{i=[n\hat{\theta}]+1}^{[n(1-\hat{\theta})]} \{ \cdots \},$$

$$S_{n2} = \frac{1}{n^r \sqrt{p}} \sum_{i=2}^{[n\hat{\theta}]} \{ \cdots \},$$

and

$$S_{n3} = \frac{1}{n^r \sqrt{p}} \sum_{i=[n(1-\hat{\theta})]+1}^n \{ \cdots \}$$

converge to 0 almost surely as $n \rightarrow \infty$, where $[a]$ is the largest integer part of a .

We deal with S_{n1} first. Recall the definition of $G_2(x)$ and $G_1(x)$ in condition (C) of Theorem 1. For any fixed δ , we define two events: For any $\beta > 0$,

$$A_n = \{y_{([n\delta])} > G_2^{-1}(\beta)\} \quad \text{and}$$

$$B_n = \{y_{([n(1-\delta)])} < G_1^{-1}(1-\beta)\}.$$

For $\delta < 1/2$, select $0 < \beta < \delta(1/2)^{1/\delta}$ and a B_0 such that $B_0 > G_1^{-1}(1-\beta)$ and $-B_0 < G_2^{-1}(\beta)$. Then for any small η , for n sufficiently large, we can derive that, combining (6) with condition (D) of Theorem 1,

$$\begin{aligned} P\{S_{n1} > \eta\} &= P\left\{ \frac{1}{n^r \sqrt{p}} \sum_{i=[n\delta]+1}^{[n(1-\delta)]} \|\mathbf{m}(y_{(i)}) - \mathbf{m}(y_{(i-1)})\| > \eta \right\} \\ &\leq P\left\{ \frac{1}{n^r \sqrt{p}} \sum_{i=[n\delta]+1}^{[n(1-\delta)]} \|\mathbf{m}(y_{(i)}) - \mathbf{m}(y_{(i-1)})\| > \eta, A_n B_n \right\} \\ &\quad + P\{(A_n B_n)^c\} \\ &\leq P(A_n^c) + P(B_n^c) + o(1), \end{aligned}$$

where $A_n B_n$ represents the intersection of two events A_n and B_n .

Now we deal with $P(A_n^c)$ and $P(B_n^c)$. Suppose that $w_{(i)}$ is the i th-order statistic of iid random variables $w_i, i = 1, \dots, n$, that has the distribution $G_2(x)$ of W_2 because by condition (C), $Y \succ W_2$, and then $y_{([n\delta])} \succ w_{([n\delta])}$. Thus,

$$P(A_n^c) = P\{y_{([n\delta])} \leq G_2^{-1}(\beta)\} \leq P\{w_{([n\delta])} \leq G_2^{-1}(\beta)\}.$$

From the properties of order statistics, we know that

$$\begin{aligned} P\{w_{([n\delta])} \leq G_2^{-1}(\beta)\} \\ \leq \frac{n!}{([n\delta]-1)!(n-[n\delta])!} \int_0^\beta t^{[n\delta]-1} (1-t)^{n-[n\delta]} dt. \end{aligned}$$

By Stirling's formula, $n!$ is asymptotically equal to $(\frac{n}{e})^n \sqrt{2\pi n}$; then we have

$$P(A_n^c) = O\left\{ \left(\frac{\beta}{\delta}\right)^{[n\delta]} \frac{1}{(1-\delta)^{n-[n\delta]}} \right\}.$$

Recall that $\frac{\beta}{\delta} < (\frac{1}{2})^{1/\delta}$ and $\delta < \frac{1}{2}$. Then there is a ρ_1 with $0 < \rho_1 < 1$,

$$P(A_n^c) = O(\rho_1^n). \quad (\text{A.4})$$

Similarly, we can obtain that for some ρ_2 with $0 < \rho_2 < 1$,

$$P(B_n^c) = O(\rho_2^n). \quad (\text{A.5})$$

From (A.4) and (A.5), the strong convergence of S_{n1} follows the Borel–Cantelli lemma.

Because the arguments for proving the convergence of S_{n2} and S_{n3} are similar, we present only the proof for S_{n3} here. Invoking condition (B) of Theorem 1 and the monotonic property of $\mathbf{M}(\cdot)$ assumed in (7), we have

$$\begin{aligned} S_{n3} &= \frac{1}{n^r \sqrt{p}} \sum_{i=[n(1-\delta)+1]}^n \|\mathbf{m}(y_{(i)}) - \mathbf{m}(y_{(i-1)})\| \\ &\leq n^{-r} \sum_{i=[n(1-\delta)+1]}^n |\mathbf{M}(y_{(i)}) - \mathbf{M}(y_{(i-1)})| \end{aligned}$$

$$= n^{-r} |\mathbf{M}(y_{[n(1-\delta)]}) - \mathbf{M}(y_{(n)})|$$

$$\leq n^{-r} |\mathbf{M}(y_{(n)})|.$$

Without loss of generality, we assume that $\mathbf{M}(y_{(n)})$ is positive. Then we need only prove that $n^{-r} \mathbf{M}(y_{(n)}) \rightarrow 0$ a.s. $n \rightarrow \infty$. Note that for any $\eta > 0$,

$$P\{n^{-r} \mathbf{M}(y_{(n)}) > \eta\} = 1 - [P\{n^{-r} \mathbf{M}(y) \leq \eta\}]^n$$

$$= 1 - [1 - P\{n^{-r} \mathbf{M}(y) > \eta\}]^n.$$

By condition (C), there is a W_1 such that $Y < W_1$, and from condition (D), it is easily derived that for a constant $b > 0$,

$$P\{n^{-r} \mathbf{M}(y) > \eta\} \leq P\{n^{-r} \mathbf{M}(W_1) > \eta\} \leq \frac{b}{(\eta n^r)^l}.$$

Using the inequality $(1 - \alpha)^n \geq 1 - n\alpha$, $0 \leq \alpha \leq 1$, we obtain that

$$P\{n^{-r} \mathbf{M}(y_{(n)}) > \eta\} = 1 - [1 - P\{n^{-r} \mathbf{M}(y) > \eta\}]^n \leq \frac{b}{(\eta n^r)^l}.$$

As $l \times r > 2$, the Borel–Cantelli lemma yields $n^{-r} \mathbf{M}(y_{(n)}) \rightarrow 0$ a.s. $n \rightarrow \infty$, and then the convergence of S_{n3} is proved. The proof of the lemma is finished.

Lemma A.2. If \mathbf{M}_1 is an $n \times p$ real-valued matrix, \mathbf{M}_2 is a $p \times n$ real-valued matrix, and \mathbf{M}_3 and \mathbf{M}_4 are two $n \times n$ real-valued matrices, then

$$\|\mathbf{M}_1 \mathbf{M}_2\| \leq \lambda_{\max}^{1/2}(\mathbf{M}_2 \mathbf{M}_2^T) \|\mathbf{M}_1\|,$$

$$\lambda_{\max}(\mathbf{M}_1 \mathbf{M}_2) = \lambda_{\max}(\mathbf{M}_2 \mathbf{M}_1),$$

and

$$\lambda_{\max}(\mathbf{M}_3 \mathbf{M}_4) \leq \lambda_{\max}(\mathbf{M}_3) \lambda_{\max}(\mathbf{M}_4),$$

where $\lambda_{\max}(\cdot)$ represents the maximum eigenvalue (see Bai 1999).

Lemma A.3. Under the conditions of Theorem 1, we have

$$\max_j \max_i |\varepsilon_{(i)}(j)| = \max_j \max_i |\varepsilon_i(j)|$$

$$= O_p(n^{1/l} p) \text{ or } O(n^{1/l} p (\log n)^{2/l}) \quad \text{a.s.,}$$

where $j = 1, \dots, p$ and $i = 1, \dots, n$, $\varepsilon_{(i)}(j)$ and $\varepsilon_i(j)$ are the j th components of random vector $\varepsilon_{(i)}$ and ε_i defined in (5).

Proof. For any $j = 1, \dots, p$, let $|\varepsilon_{nj}| = \max_{1 \leq i \leq n} |\varepsilon_{(i)}(j)|$. It is clear that $|\varepsilon_{nj}| = \max_i |\varepsilon_i(j)|$ and $\max_j \max_i |\varepsilon_{(i)}(j)| = \max_i |\varepsilon_{nj}|$. We investigate only the almost-sure convergence, because the convergence in probability is much easier to obtain. Choosing $\eta = bn^{1/l} p (\log n)^{1/2}$, together with the iid property of variables ε_i , we have that

$$P\left(\max_j |\varepsilon_{nj}| \geq \eta\right)$$

$$\leq \max_j P\{|\varepsilon_{nj}| \geq bn^{1/l} p (\log n)^{2/l}\}$$

$$\leq \max_j [1 - P\{|\varepsilon_{nj}| \leq bn^{1/l} p (\log n)^{2/l}\}]$$

$$= \max_j (1 - [1 - P\{|\varepsilon_1(j)| \geq bn^{1/l} p (\log n)^{2/l}\}]^n).$$

Furthermore, invoking condition (A), we can easily derive the finiteness of the l th moment of $|\varepsilon_i(j)|$, and

$$P\{|\varepsilon_1(j)| \geq bn^{1/l} p (\log n)^{2/l}\} \leq \frac{b}{n(\log n)^2}.$$

Hence

$$\max_j (1 - [1 - P\{|\varepsilon_1(j)| \geq bn^{1/l} p (\log n)^{2/l}\}]^n)$$

$$\leq 1 - \left\{1 - \frac{b}{n(\log n)^2}\right\}^n \leq 1 - \exp\left\{-\frac{b}{(\log n)^2}\right\}$$

$$\leq \frac{b}{(\log n)^2}.$$

The Borel–Cantelli lemma implies strong convergence of the subsequence $\{\max_{1 \leq j \leq p} |\varepsilon_{nmj}| : m \geq 1\}$ of $\{\max_{1 \leq j \leq p} |\varepsilon_{nj}| : n \geq 1\}$. Using the subsequence approach, we need only show that for any n with $e^m \leq n \leq e^{m+1}$,

$$P\left\{\max_{e^m \leq n \leq e^{m+1}} \left(\max_j |\varepsilon_{nj}| - \max_j |\varepsilon_{e^m j}|\right) \geq \eta\right\}$$

$$= P\left\{\max_j |\varepsilon_{e^{m+1} j}| - \max_j |\varepsilon_{e^m j}| \geq \eta\right\}$$

$$\leq \max_j P\left\{\max_{e^m \leq n \leq e^{m+1}} |\varepsilon_{nj}| \geq be^{m/l} m^{2/l}\right\}$$

$$\leq \frac{b}{m^2}.$$

The last inequality comes from the same argument as the foregoing with the upper bound $b/(\log n)^2$. This means that the difference between ε_{nj} and its subsequence also converges almost surely to 0 at the rate $n^{-1/l} p (\log n)^{2/l}$.

For the convergence in probability, we do not need to use the subsequence approach. The rate is clearly $O(n^{-1/l} p)$. The proof is completed.

A.3 Proof of Theorem 1

The Proof for Strong Convergence. As described earlier, to prove the strong convergence of $\hat{\mathbf{A}}_p$ we need only prove convergence of the right side of (A.3). The proof is divided into four steps:

a. $|\mathbf{J}_2| = o(n^{-1+1/l+r} p^{3/2} (\log n)^{2/l})$ a.s. as $n \rightarrow \infty$.

From the definition of \mathbf{J}_2 in (A.1), we need only consider the term

$$\left\| \frac{1}{H(c-1)} \varepsilon_n (\mathbf{I}_H \otimes \mathbf{P}_c) \tilde{\mathbf{M}}_n^T \right\|.$$

Recalling the definitions of \mathbf{P}_c and $\mathbf{e}_c \mathbf{e}_c^T$ in Section A.1, we have $\mathbf{P}_c \mathbf{e}_c \mathbf{e}_c^T = 0$. Invoking the equations $\tilde{\mathbf{M}}_n^T = \mathbf{T}_n \tilde{\mathbf{M}}_n^T$ and $\mathbf{S}_n = \mathbf{T}_n^{-1}$, we can derive that, together with the properties of the Kronecker product and recalling that $\mathbf{S}_n = \mathbf{T}_n^{-1}$,

$$\|\varepsilon_n (\mathbf{I}_H \otimes \mathbf{P}_c) \tilde{\mathbf{M}}_n^T\| = \|\varepsilon_n (\mathbf{I}_H \otimes \mathbf{P}_c) \mathbf{S}_n \tilde{\mathbf{M}}_n^T\| = \|\varepsilon_n (\mathbf{I}_H \otimes (\mathbf{P}_c \mathbf{S}_c)) \tilde{\mathbf{M}}_n^T\|.$$

From the definitions of ε_n and \mathcal{E}_n in Section A.1, we have

$$\|\varepsilon_n (\mathbf{I}_H \otimes (\mathbf{P}_c \mathbf{S}_c)) \tilde{\mathbf{M}}_n^T\| \leq \max_{1 \leq i \leq n, 1 \leq j \leq p} |\varepsilon_{(i)}(j)| \|\mathcal{E}_n (\mathbf{I}_H \otimes (\mathbf{P}_c \mathbf{S}_c)) \tilde{\mathbf{M}}_n^T\|.$$

Also note that in matrix $\mathbf{I}_H \otimes (\mathbf{P}_c \mathbf{S}_c)$, all elements of the last column are 0 and the elements such as $a_{k_1 k_2} = 0$ if $|k_1 - k_2| > c$ and $a_{k_1 k_2} \leq 1$ otherwise. Then, by the first inequality of Lemma A.2, we have

$$\|\mathcal{E}_n (\mathbf{I}_H \otimes (\mathbf{P}_c \mathbf{S}_c)) \tilde{\mathbf{M}}_n^T\| \leq 2c \sum_{i=2}^n \|\mathbf{m}(y_{(i)}) - \mathbf{m}(y_{(i-1)})\|.$$

Applying Lemma A.3 to $\varepsilon_{(i)}(j)$ and Lemma A.1 to $\tilde{\mathbf{M}}_n^T$, we have

$$\max_{1 \leq j \leq p} \max_{1 \leq i \leq n} |\varepsilon_{(i)}(j)| = O(n^{1/l} p (\log n)^{2/l}) \quad \text{a.s.}$$

and

$$\sum_{i=2}^n \|\mathbf{m}(y_{(i)}) - \mathbf{m}(y_{(i-1)})\| = O(n^r \sqrt{p}) \quad \text{a.s.}$$

Hence we almost surely have

$$\frac{1}{H(c-1)} \|\varepsilon_n(\mathbf{I}_H \otimes (\mathbf{P}_c \mathbf{S}_c)) \tilde{\mathbf{M}}_n^T\| = O(n^{-1+1/l+r} p^{3/2} (\log n)^{1/l}).$$

b. $\|\mathbf{J}_3\| = O(n^{2r-1}p)$ a.s. as $n \rightarrow \infty$.

By Lemma A.2 and the property of the Kronecker product, we can obtain that, from the definition of \mathbf{M}_n^T and $\tilde{\mathbf{M}}_n^T$,

$$\|\mathbf{M}_n(\mathbf{I}_H \otimes \mathbf{P}_c) \mathbf{M}_n^T\| = \|\tilde{\mathbf{M}}_n(\mathbf{I}_H \otimes \mathbf{S}_c^T \mathbf{P}_c \mathbf{S}_c) \tilde{\mathbf{M}}_n^T\|.$$

Note that the last row and last column of the matrix $\mathbf{I}_H \otimes \mathbf{S}_c^T \mathbf{P}_c \mathbf{S}_c$ are both 0. Let $\tilde{\mathbf{M}}_{n1}$ be a matrix with elements equal to those of $\tilde{\mathbf{M}}_n$ except that the elements of the last column are 0. Then the matrix $\tilde{\mathbf{M}}_n(\mathbf{I}_H \otimes \mathbf{S}_c^T \mathbf{P}_c \mathbf{S}_c) \tilde{\mathbf{M}}_n^T$ can be rewritten as $\tilde{\mathbf{M}}_{n1}(\mathbf{I}_H \otimes \mathbf{S}_c^T \mathbf{P}_c \mathbf{S}_c) \tilde{\mathbf{M}}_{n1}^T$. It is also clear that $\lambda_{\max}^{1/2}(\mathbf{S}_c^T \mathbf{P}_c \mathbf{S}_c) \leq b$ for some $b > 0$. Thus, invoking Lemma A.1,

$$\begin{aligned} \|\mathbf{J}_3\| &= \frac{1}{H(c-1)} \|\tilde{\mathbf{M}}_n(\mathbf{I}_H \otimes \mathbf{P}_c) \tilde{\mathbf{M}}_{n1}^T\| \\ &\leq \frac{bc}{n(c-1)} \|\tilde{\mathbf{M}}_{n1}^T\|^2 = O(n^{2r-1}p) \quad \text{a.s.} \end{aligned}$$

The proof is finished.

c. $\|\mathbf{K}_2\| = O(n^{-1/4}p(\log n)^{1/2})$ a.s. as $n \rightarrow \infty$.

Note that $\|\mathbf{K}_2\|^4$ can be written as

$$\|\mathbf{K}_2\|^4 = \frac{1}{H^4(c-1)^4 c^4} \left\| \sum_{h=1}^H \varepsilon_{h,c} \mathbf{Q}'_c \varepsilon_{h,c}^T \right\|^4,$$

where $\varepsilon_{h,c}$ is defined in Section A.1 and

$$\mathbf{Q}'_c = c \mathbf{Q}_c = \begin{pmatrix} 0 & -1 & \cdots & -1 \\ -1 & \ddots & -1 & -1 \\ \vdots & & & \vdots \\ -1 & \cdots & -1 & 0 \end{pmatrix}_{c \times c}.$$

Furthermore,

$$\begin{aligned} E \left\| \sum_{h=1}^H \varepsilon_{h,c} \mathbf{Q}'_c \varepsilon_{h,c}^T \right\|^4 &= E \left[\text{tr} \left(\sum_{h_1=1}^H \varepsilon_{h_1,c} \mathbf{Q}'_c \varepsilon_{h_1,c}^T \sum_{h_2=1}^H \varepsilon_{h_2,c} \mathbf{Q}'_c \varepsilon_{h_2,c}^T \right) \right]^2 \\ &= \sum_{h_1=1}^H \sum_{h_2=1}^H \sum_{h_3=1}^H \sum_{h_4=1}^H E [\text{tr}(\varepsilon_{h_1,c} \mathbf{Q}'_c \varepsilon_{h_1,c}^T \varepsilon_{h_2,c} \mathbf{Q}'_c \varepsilon_{h_2,c}^T) \\ &\quad \times \text{tr}(\varepsilon_{h_3,c} \mathbf{Q}'_c \varepsilon_{h_3,c}^T \varepsilon_{h_4,c} \mathbf{Q}'_c \varepsilon_{h_4,c}^T)]. \end{aligned}$$

Together with the special structure of \mathbf{Q}'_c and the conditional independence of $\varepsilon_{(i)}$ given the $y_{(i)}$'s, we know that $E\{\text{tr}(\varepsilon_{h_i,c} \mathbf{Q}'_c \varepsilon_{h_i,c}^T)\} = 0$ for $i = 1, 2, 3, 4$. Then it is easy to derive that if $h_1 \neq h_j$, where $j = 2, 3, 4$,

$$E \text{tr}(\varepsilon_{h_1,c} \mathbf{Q}'_c \varepsilon_{h_1,c}^T \varepsilon_{h_2,c} \mathbf{Q}'_c \varepsilon_{h_2,c}^T) \text{tr}(\varepsilon_{h_3,c} \mathbf{Q}'_c \varepsilon_{h_3,c}^T \varepsilon_{h_4,c} \mathbf{Q}'_c \varepsilon_{h_4,c}^T) = 0.$$

For the other cases, invoking the conditions of Theorem 1 and the conditional independence of the $\varepsilon_{(i)}$'s given the $y_{(i)}$'s, we have that for some $b > 0$,

$$\begin{aligned} |E[\text{tr}(\varepsilon_{h_1,c} \mathbf{Q}'_c \varepsilon_{h_1,c}^T \varepsilon_{h_2,c} \mathbf{Q}'_c \varepsilon_{h_2,c}^T) \text{tr}(\varepsilon_{h_3,c} \mathbf{Q}'_c \varepsilon_{h_3,c}^T \varepsilon_{h_4,c} \mathbf{Q}'_c \varepsilon_{h_4,c}^T)]| \\ \leq p^4 c^8 E[|\varepsilon_{(j)}|^4] \leq b p^4 c^8. \end{aligned}$$

Thus

$$E \left(\left\| \sum_{h=1}^H \varepsilon_{h,c} \mathbf{Q}'_c \varepsilon_{h,c}^T \right\|^4 \right) \leq 3bH^2 p^4 c^8 = 3bn^2 p^4 c^6,$$

and then, for any $\eta > 0$,

$$P(\|\mathbf{K}_2\| > \eta) \leq \frac{E(\|\mathbf{K}_2\|^4)}{\eta^4} \leq \frac{3bn^2 p^4 c^6}{n^4 (c-1)^4 \eta^4} = O(n^{-2} p^4 / \eta^4).$$

Choosing $\eta = n^{-1/4} p (\log n)^{1/2}$, the Borel–Cantelli lemma implies the conclusion.

d. $\|\mathbf{K}_1 - \mathbf{\Lambda}_p\| = O(n^{-1/2} p (\log n))$ a.s. as $n \rightarrow \infty$.

Note that $\mathbf{\Lambda}_p$ is the expectation of \mathbf{K}_1 . This means that $\mathbf{K}_1 - \mathbf{\Lambda}_p$ is a centered sample mean. Each of its elements is a sum of the independent identically distributed variables with mean 0 and finite second moment. Because the fourth moment of X is finite, invoking the result of theorem IX 20 of Petrov (1995, p. 278), every element of $\|\mathbf{K}_1 - \mathbf{\Lambda}_p\|$ has order $O(n^{-1/2} \log n)$ a.s. Together with the definition of the Euclidean norm, the resulting convergence rate is $pn^{-1/2} \log n$. It can be immediately derived that $\|\mathbf{K}_1 - \mathbf{\Lambda}_p\| = O(n^{-1/2} p \log n)$ a.s. as $n \rightarrow \infty$.

Together with steps a–d, and $r < 1/4$, we complete the proof of the strong convergence of $\|\hat{\mathbf{\Lambda}}_p - \mathbf{\Lambda}_p\|$.

As for the convergence of the term $\|(\hat{\Sigma}_x - \hat{\Lambda}_p) - (\Sigma_x - \Lambda_p)\|$, we need only note that $\|\hat{\Sigma}_x - \Sigma_x\|$ can have the same convergence rate as $\|\hat{\Lambda}_p - \Lambda_p\|$. The proof is similar to step d, as follows. Consider every element of $\Sigma_x - \hat{\Sigma}_x$. Without loss of generality, we consider the left-upper element on the diagonal. We know that the left-upper element, $\sigma_{1,1}$, of Σ_x can be rewritten as the difference of the second moment and the square of the first moment of X , and similarly for $\hat{\Sigma}_x$. As the l th moment of X is finite, applying the result of Petrov (1995) to the difference relating to the second-moment and first-moment terms, we can easily obtain the convergence rate $O(n^{-1/2} (\log n))$. Together with the definition of the Euclidean norm again, the resulting convergence rate is $O(pn^{-1/2} (\log n))$.

The Proof for Weak Convergence. We now turn to the proof of the weak convergence. From the lemmas, we have that in step a we use the rate for the convergence in probability to obtain $\|\mathbf{J}_2\| = O(p(n^{-1+1/l+r} p^{3/2}))$, in step b we use the exactly the same argument to obtain $\|\mathbf{J}_3\| = O(p(pn^{-1+2r}))$, in step c we can choose $\eta = n^{1/2} p$ with no use of the Borel–Cantelli lemma to achieve the rate $\|\mathbf{K}_2\| = O(p(pn^{-1/2}))$, and in step d a standard argument for weak convergence is adopted to arrive at $\|\mathbf{K}_1 - \mathbf{\Lambda}_p\| = (pn^{-1/2})$ (see Petrov 1995). When $p = o(n^{\min\{1/2, 1-2/l-2r\}})$, the desired convergence rate can be achieved. The weak convergence of $\|(\hat{\Sigma}_x - \hat{\Lambda}_p) - (\Sigma_x - \Lambda_p)\|$ can be proved similarly. The proof of Theorem 1 is finished.

A.4 Proof of Corollary 1

Recall the definition of λ_{pi} before Corollary 1. Note that $\lambda_{p1} \geq \lambda_{p2} \geq \cdots \geq \lambda_{pK} > \lambda_{p(K+1)} = \cdots = \lambda_{pp} = 0$. It is clear that

$$\sum_{k=1}^p |\lambda_{pk} - \hat{\lambda}_{pk}|^2 \leq \|\mathbf{\Lambda}_p - \hat{\mathbf{\Lambda}}_p\|^2.$$

The convergence of $\hat{\lambda}_{pk}$ can be derived from Theorem 1 immediately.

For the convergence of the estimate $\hat{P}_{\lambda_{pi}}$ of the projection space $P_{\lambda_{pi}}$, because the nonzero eigenvalues are distinct, the projection spaces are just those of eigenvectors. For brevity, let $\mathbf{M}_p = \Sigma_x - \Lambda_p$ and $\hat{\mathbf{M}}_{pn} = \hat{\Sigma}_x - \hat{\Lambda}_p$. Therefore, invoking Theorem 1, Corollary 1, and the representation of $\hat{P}_{\lambda_{pi}}$ of Zhu and Ng (1995) for fixed p , the conclusion is reached immediately. We omit the details of the proof here.

A.5 Proof of Theorem 2

We deal with the strong convergence first. The proof for the weak convergence is similar. Throughout the proof, the argument is with

probability 1 unless stated otherwise. Let K be the true value of the dimension of $\Psi = \Sigma_X - \Lambda_p$. Note that

$$G(K) - G(k) = \log L_K - \log L_k - \frac{C_n(K-k)(2p-k-K+1)}{2}.$$

From Corollary 1, we have that for large n ,

$$\hat{\theta}_i > 1, \quad i = 1, \dots, K, \quad \text{and} \quad \min(\tau, K) = K,$$

where τ is the number of $\hat{\theta}_i$ with $\hat{\theta}_i > 1$. If $k < K$, then $\min(\tau, k) = k$. Thus, for large n ,

$$\log L_K - \log L_k = -\frac{1}{2}n \sum_{i=k+1}^K (\log \hat{\theta}_i + 1 - \hat{\theta}_i) = \frac{1}{2}n W_n(K, k),$$

where

$$W_n(K, k) = - \sum_{i=k+1}^K (\log \hat{\theta}_i + 1 - \hat{\theta}_i).$$

We have that

$$\lim_{n \rightarrow \infty} W_n(K, k) = W(K, k) \equiv - \sum_{i=k+1}^K (\log \theta_i + 1 - \theta_i) > 0.$$

Hence we have that for large n ,

$$\log L_K - \log L_k > \frac{1}{4}n W(K, k).$$

Note that $\lim_{n \rightarrow \infty} C_n/n^{1-s} = 0$ and $p = O(n^s)$. Then

$$\frac{C_n(K-k)(2p-k-K+1)}{n} \rightarrow 0. \quad (\text{A.6})$$

Therefore, recalling the formula of $G(K) - G(k)$, for large n , we have that

$$G(K) - G(k) > 0. \quad (\text{A.7})$$

If $k > K$, then we have, together with (10),

$$|\log L_K - \log L_k| \leq n \sum_{i=K+1}^p |\log \hat{\theta}_i + 1 - \hat{\theta}_i|.$$

Invoking the Taylor expansion, we derive that

$$\begin{aligned} |\log L_K - \log L_k| &\leq n \sum_{i=K+1}^p \frac{1}{2}(\hat{\theta}_i - 1)^2(1 + o(1)) \\ &\leq n \|\Omega - \hat{\Omega}\|^2 = O(n^{1-2t}) \quad \text{a.s.} \end{aligned}$$

Because $\lim_{n \rightarrow \infty} C_n/n^{1-2t} = \infty$, we can see that for large n ,

$$G(K) - G(k) = O(n^{1-2t}) + \frac{C_n(k-K)(2p-k-K+1)}{2} > 0. \quad (\text{A.8})$$

From (A.7) and (A.8), it follows that for large n ,

$$\hat{K} = K.$$

Thus the strong consistency of the estimate is proved.

[Received January 2004. Revised July 2005.]

REFERENCES

- Antoniadis, A., Lambert-Lacroix, S., and Leblanc, F. (2003), "Effective Dimension Reduction Methods for Tumor Classification Using Gene Expression Data," *Bioinformatics*, 19, 563–570.
- Bai, Z. D. (1999), "Methodologies in Spectral Analysis of Large-Dimensional Random Matrices, a Review," *Statistica Sinica*, 9, 611–677.
- Bai, Z. D., Krishnaiah, R. D., and Zhao, L. C. (1989), "On Rates of Convergence of Efficient Detection Criteria in Signal Processing With White Noise," *IEEE Transactions on Information Theory*, 35, 380–388.
- Bai, Z. D., and Saranadasa, H. (1996), "Effect of High Dimension by an Example of a Two-Sample Problem," *Statistica Sinica*, 6, 311–329.
- Bickel, P. J., and Levina, E. (2004), "Some Theory for Fisher's Discriminant Function, 'Naive Bayes,' and Some Alternatives When There Are Many More Variables Than Observations," *Bernoulli*, 10, 989–1010.
- Bura, E., and Cook, R. D. (2001a), "Extending Sliced Inverse Regression: The Weighted Chi-Squared Test," *Journal of the American Statistical Association*, 96, 996–1003.
- (2001b), "Estimating the Structural Dimension of Regressions via Parametric Inverse Regression," *Journal of the Royal Statistical Society, Ser. B*, 63, 393–410.
- Bura, E., and Pfeiffer, R. M. (2003), "Graphical Methods for Class Prediction Using Dimension Reduction Techniques on DNA Microarray Data," *Bioinformatics*, 19, 1252–1258.
- Carroll, R. J., and Li, K. C. (1995), "Binary Regressors in Dimension Reduction Models: A New Look at Treatment Comparisons," *Statistica Sinica*, 5, 667–688.
- Chiaromonte, F., and Cook, R. D. (1997), "On Foundations of Regression Graphics," Technical Report 616, University of Minnesota, School of Statistics.
- Chiaromonte, F., Cook, R. D., and Li, B. (2002), "Sufficient Dimensions Reduction in Regressions With Categorical Predictors," *The Annals of Statistics*, 30, 475–497.
- Chiaromonte, F., and Martinelli, J. (2002), "Dimension Reduction Strategies for Analyzing Global Gene Expression Data With a Response," *Mathematical Biosciences*, 176, 123–144.
- Cook, R. D. (1994), "On the Interpretation of Regression Plots," *Journal of the American Statistical Association*, 89, 177–189.
- (1996), "Graphics for Regressions With a Binary Response," *Journal of the American Statistical Association*, 91, 983–992.
- (1998), *Regression Graphics: Ideas for Studying Regressions Through Graphics*, New York: Wiley.
- (2000), "SAVE: A Method for Dimension Reduction and Graphics in Regression," *Communications in Statistics, Part A—Theory and Methods*, 29, 2109–2121.
- Cook, R. D., and Lee, H. (1999), "Dimension Reduction in Binary Response Regression," *Journal of the American Statistical Association*, 94, 1187–1200.
- Cook, R. D., and Li, B. (2002), "Dimension Reduction for Conditional Mean in Regression," *The Annals of Statistics*, 30, 455–474.
- Cook, R. D., and Nachtsheim, C. J. (1994), "Reweight to Achieve Elliptically Contoured Covariates in Regression," *Journal of the American Statistical Association*, 89, 592–599.
- Cook, R. D., and Weisberg, S. (1991), Comment on "Sliced Inverse Regression for Dimension Reduction," by K. C. Li, *Journal of the American Statistical Association*, 86, 328–332.
- Eubank, R. L. (1999), *Nonparametric Regression and Spline Smoothing* (2nd ed.), New York: Marcel Dekker.
- Ferré, L. (1998), "Determining the Dimension in Sliced Inverse Regression and Related Methods," *Journal of the American Statistical Association*, 93, 132–140.
- Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesorov, J. P., Coller, H., Loh, M. L., Downing, J. R., Caligiuri, M. A., Bloomfield, C. D., and Lander, E. S. (1999), "Molecular Classification and Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring," *Science*, 286, 531–537.
- Greenshtain, E., and Ritov, Y. (2004), "Persistence in High-Dimensional Linear Predictor Selection and the Virtue of Overparametrization," *Bernoulli*, 10, 971–988.
- Hall, P., and Li, C. K. (1993), "On Almost Linearity of Low-Dimensional Projections From High-Dimensional Data," *The Annals of Statistics*, 21, 867–889.
- Hooper, J. (1959), "Simultaneous Equations and Canonical Correlation Theory," *Econometrica*, 27, 245–256.
- Hsing, T., and Carroll, R. J. (1992), "An Asymptotic Theory for Sliced Inverse Regression," *The Annals of Statistics*, 20, 1040–1061.
- Levina, E. (2002), "Statistical Issues in Texture Analysis," unpublished doctoral thesis, University of California Berkeley, California.
- Li, K. C. (1991), "Sliced Inverse Regression for Dimension Reduction" (with discussions), *Journal of the American Statistical Association*, 86, 316–342.

- (1992), “On Principal Hessian Directions for Data Visualization and Dimension Reduction,” *Journal of the American Statistical Association*, 87, 1025–1040.
- Li, Y. X., and Zhu, L. X. (2004), “When Are Sliced Average Variance Estimates Convergent?” technical report, University of Hong Kong, Dept. of Statistics and Actuarial Science.
- Petrov, V. V. (1995), *Limit Theorems of Probability Theory: Sequences of Independent Random Variables*, New York: Oxford University Press.
- Pollard, D. (1984), *Convergence of Stochastic Processes*, New York: Springer-Verlag.
- Schott, J. R. (1994), “Determining the Dimensionality in Sliced Inverse Regression,” *Journal of the American Statistical Association*, 89, 141–148.
- Schwarz, G. (1978), “Estimating the Dimension of a Model,” *The Annals of Mathematical Statistics*, 30, 461–464.
- Velilla, S. (1998), “Assessing the Number of Linear Components in a General Regression Problem,” *Journal of the American Statistical Association*, 93, 1088–1098.
- Xia, Y. C., Tong, H., Li, W. K., and Zhu, L. X. (2002), “An Adaptive Estimation of Dimension Reduction Space,” *Journal of the Royal Statistical Society, Ser. B*, 64, 363–410.
- Yang, S. S. (1977), “General Distribution Theory of the Concomitants of Order Statistics,” *The Annals of Statistics*, 5, 996–1002.
- Yin, X. R., and Cook, R. D. (2002), “Dimension Reduction for the Conditional k th Moment in Regression,” *Journal of the Royal Statistical Society, Ser. B*, 64, 159–175.
- Zhao, L. C., Krishnaiah, P. R., and Bai, Z. D. (1986a), “On Detection of the Number of Signals in Presence of White Noise,” *Journal of Multivariate Analysis* 20, 1–25.
- (1986b), “On Detection of the Number of Signals When the Noise Covariance Matrix Is Arbitrary,” *Journal of Multivariate Analysis* 20, 26–49.
- Zhu, L. X., and Fang, K. T. (1996), “Asymptotics for Kernel Estimate of Sliced Inverse Regression,” *The Annals of Statistics*, 24, 1055–1068.
- Zhu, L. X., and Ng, K. W. (1995), “Asymptotics of Sliced Inverse Regression,” *Statistica Sinica*, 5, 727–736.
- Zhu, L. X., Ohtaki, M., and Li, Y. X. (2006), “Extensions of Sliced Inverse Regression-Based Algorithms,” technical report, University of Hong Kong, Dept. of Statistics and Actuarial Science.
- Zhu, L. X., Zhu, L. P., and Li, X. (2006), “Transformed Partial Least Squares for Multivariate Data,” *Statistica Sinica*, accepted.