

## Statistica Sinica Preprint No: SS-2017-0532

<b>Title</b>	On supervised reduction and its dual
<b>Manuscript ID</b>	SS-2017-0532
<b>URL</b>	<a href="http://www.stat.sinica.edu.tw/statistica/">http://www.stat.sinica.edu.tw/statistica/</a>
<b>DOI</b>	10.5705/ss.202017.0532
<b>Complete List of Authors</b>	Peirong Xu and Tao Wang
<b>Corresponding Author</b>	Richard J. Samworth
<b>E-mail</b>	neowangtao@sjtu.edu.cn
Notice: Accepted version subject to English editing.	

## ON SUPERVISED REDUCTION AND ITS DUAL

Peirong Xu and Tao Wang

*Shanghai Normal University and Shanghai Jiao Tong University*

*Abstract:* Existing methods for dimension reduction in regression estimate a subspace in the primal predictor-based space, and then obtain the set of reduced predictors by projecting the original predictor vector onto this subspace. We propose a principled method for estimating a sufficient reduction in the dual sample-based space, on the basis of a supervised inverse regression model. Reduction is done without the need to estimate the subspace. Our method extends the duality between principal component analysis and principal coordinate analysis. We study the asymptotic behavior of the proposed method, and demonstrate that it is robust to model misspecification. We present simulation results to support the theoretical conclusion, and illustrate the application of the method in real data analysis.

*Key words and phrases:* Data visualization, inverse model-based reduction, multidimensional scaling, sufficient dimension reduction, supervised coordinate analysis.

## 1. Introduction

Dimension reduction is a long-standing and prominent problem in regression analysis (Cook, 2007). Classical methods for dimension reduction in regression transform the predictors and then fit a least squares model using the transformed variables. For example, the widely used principal component regression involves extracting the first few principal components of the predictors, and then using these components as the predictors in a linear model. One concern of principal component regression is that the directions in which the predictors show the most variation are not necessarily the directions that are associated with the response. Many methods have been proposed to deal with this issue, such as partial least squares and sliced inverse regression (Li, 1991). A common goal of these methods is to reduce the dimensionality of the predictors without losing any information about the response.

Suppose we have a response  $Y \in \mathbb{R}$  and a vector of predictors  $\mathbf{X} \in \mathbb{R}^p$ . Formally, the aim is to estimate a reduction  $\mathcal{R} : \mathbb{R}^p \rightarrow \mathbb{R}^d, d \leq p$ , such that

$$Y \perp\!\!\!\perp \mathbf{X} \mid \mathcal{R}(\mathbf{X}), \quad (1.1)$$

where  $\perp\!\!\!\perp$  indicates independence.  $\mathcal{R}(\mathbf{X})$  is called a sufficient reduction for the regression of  $Y$  onto  $\mathbf{X}$  (Cook, 1998). Sufficient dimension reduction has been an active research area since the introduction of sliced inverse

## 1 INTRODUCTION

---

regression and sliced average variance estimation (Cook and Weisberg, 1991).

In this paper we focus on linear reductions that are linear combinations of the predictors:  $\mathcal{R}(\mathbf{X}) = \boldsymbol{\eta}^\top \mathbf{X}$  for some  $p \times d$  matrix  $\boldsymbol{\eta}$ .

Depending on the stochastic nature of  $Y$  and  $\mathbf{X}$ , there are three paradigms for determining a sufficient reduction: forward reduction, inverse reduction, and joint reduction, and they are equivalent when  $Y$  and  $\mathbf{X}$  are jointly distributed (Cook, 2007). Without requiring a pre-specified model for  $Y \mid \mathbf{X}$ , inverse reduction is promising in regressions with many predictors. To estimate a reduction inversely, methods such as sliced inverse regression exploit properties of the conditional moments of  $\mathbf{X} \mid Y$ . These inverse moment-based methods impose constraints on the marginal distribution of  $\mathbf{X}$ . Alternatively, inverse model-based approaches directly specify a model for the inverse regression of  $\mathbf{X}$  onto  $Y$ . Much of existing work relies on normal models. See Adraghi and Cook (2009) for a recent review of inverse reduction methods.

Sufficient reduction permits us to restrict attention to a few new predictors  $\boldsymbol{\eta}^\top \mathbf{X}$ , upon which subsequent modeling and prediction can be built. Indeed, the original intent behind (1.1) is to provide a framework for dimension reduction to facilitate graphical analyses (Cook, 1998). Previous studies have largely focused on properties of estimators of the subspace spanned by

## 1 INTRODUCTION

---

the columns of  $\boldsymbol{\eta}$ . However, the inference object more relevant to subsequent data analyses is not the subspace but the reduction itself. Estimating sufficient reductions is relatively new. Cook et al. (2012) studied the asymptotic behavior of a class of methods for sufficient reduction in large abundant regressions, where most predictors contribute some information on the response. In the modern “small  $n$  and large  $p$ ” setting, Wang et al. (2018) recently proposed a method for estimating sparse reductions using a novel representation of sliced inverse regression. The former is inverse model-based, and the latter is inverse moment-based.

In this article, we propose a new approach for estimating a sufficient reduction. This is motivated by the well-known duality between principal component analysis and principal coordinate analysis (Gower, 1966), also known as classical multidimensional scaling (Hastie et al., 2009). Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top \in \mathbb{R}^{n \times p}$  be the data matrix of predictor values. Without loss of generality, assume that each column of  $\mathbf{X}$  has been centered to have mean zero. The singular value decomposition offers a way of expressing principal component analysis. Let  $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$  be the singular value decomposition of  $\mathbf{X}$ ; that is,  $\mathbf{U} = (\mathbf{U}_1, \dots, \mathbf{U}_p)$  is  $n \times p$  with orthonormal columns,  $\mathbf{V} = (\mathbf{V}_1, \dots, \mathbf{V}_p)$  is  $p \times p$  orthogonal, and  $\mathbf{D}$  is a  $p \times p$  diagonal matrix with diagonal entries  $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$ .  $\mathbf{V}_j$  is called the  $j$ th principal

## 1 INTRODUCTION

---

component direction, and  $d_j \mathbf{U}_j$  the  $j$ th principal component score vector. In the terminology of dimension reduction in regression,  $\mathbf{UD} = \mathbf{XV}$  are linear reductions used in principal component regression. Geometrically, each row of  $\mathbf{UD}$  represents the coordinates of the corresponding row of  $\mathbf{X}$  with respect to the orthonormal basis  $\mathbf{V}$ . In this sense, principal component analysis can be viewed as an ordination method. Indeed, it is equivalent to classical multidimensional scaling, and an alternative way of obtaining principal components is to perform an eigen-decomposition of the Gram matrix  $\mathbf{XX}^\top = \mathbf{UD}^2\mathbf{U}^\top$ .

Instead of estimating the directions, one can directly determine the projection coordinates of the predictor vector onto the subspace spanned by these directions. In the context of moment-based inverse reduction, Zhang et al. (2008) calculated the projection coordinates by applying classical multidimensional scaling to slice means, and then interpolated the projection of a new predictor vector using these coordinates. The method can be thought of as a dual version of sliced inverse regression. At the population level, however, it is not clear what quantity is being taking as the target. In this paper, we propose a principled method for estimating a sufficient reduction under the inverse model-based reduction scheme. Reduction is done without the need to estimate the subspace. For the first time we

## 2 A NAIVE INVERSE REGRESSION MODEL

---

study the asymptotics of predictor reduction in terms of prediction and under model misspecification.

To express the projection coordinates explicitly, an inverse regression model is introduced in Section 2, without requiring normal errors. Since the coordinates are unconstrained, sufficient reduction is achieved by classical multidimensional scaling, or principal component analysis by duality. To perform supervised reduction, we model the coordinates in a parametric way in Section 3, and extend the method of Section 2 when the error structure is known. Reduction with an unknown error structure is considered in Section 4. Some theoretical conclusions are provided. Simulation results and a real data application are presented in Section 5. A concluding discussion is given in Section 6. Proofs can be found in the Supplementary Materials.

### 2. A naive inverse regression model

The subspace spanned by the columns of  $\boldsymbol{\eta}$  is called a dimension-reduction subspace. The parsimonious target of sufficient dimension reduction is the central subspace  $\mathcal{S}_{Y|X}$ , defined as the intersection of all dimension-reduction subspaces (Cook, 1998). Let  $\mathbb{Y}$  denote the sample space of  $Y$ , and let

$$\mathcal{S}_{E(\mathbf{X}|Y)} = \text{span}\{E(\mathbf{X} | Y = y) - E(\mathbf{X}), y \in \mathbb{Y}\}$$

denote the subspace spanned by the centered inverse regression curves. We

## 2 A NAIVE INVERSE REGRESSION MODEL

---

have the following proposition.

**Proposition 1.** *Assume (C1)  $\mathcal{S}_{E(\mathbf{X}|Y)} \subseteq \text{Var}(\mathbf{X})\mathcal{S}_{Y|\mathbf{X}}$  and (C2)  $\text{Var}(\mathbf{X} | Y)$  is positive definite and is non-random, then*

$$\text{Var}(\mathbf{X} | Y)\mathcal{S}_{Y|\mathbf{X}} = \text{Var}(\mathbf{X})\mathcal{S}_{Y|\mathbf{X}}.$$

Conditions (C1) and (C2) are generally regarded as mild in the sufficient dimension reduction literature. Condition (C1) holds if  $E(\mathbf{X} | \boldsymbol{\eta}^\top \mathbf{X})$  is a linear function of  $\boldsymbol{\eta}^\top \mathbf{X}$ , where  $\boldsymbol{\eta}$  is a basis matrix for  $\mathcal{S}_{Y|\mathbf{X}}$ . A slightly stronger condition is given by (C1')  $\mathcal{S}_{E(\mathbf{X}|Y)} = \text{Var}(\mathbf{X})\mathcal{S}_{Y|\mathbf{X}}$ . See Li and Wang (2007) for a good recent discussion of them.

Throughout this paper, conditions (C1') and (C2) are assumed to be true. Then, by Proposition 1,

$$\mathcal{S}_{E(\mathbf{X}|Y)} = \boldsymbol{\Delta}\mathcal{S}_{Y|\mathbf{X}},$$

where  $\boldsymbol{\Delta} = \text{Var}(\mathbf{X} | Y)$ . This implies that  $\mathcal{S}_{Y|\mathbf{X}} = \text{span}(\boldsymbol{\Delta}^{-1}\boldsymbol{\Gamma})$ , where  $\boldsymbol{\Gamma} \in \mathbb{R}^{p \times d}$  is a basis matrix for  $\mathcal{S}_{E(\mathbf{X}|Y)}$ . Let  $\mathbf{X}_y$  denote a random vector distributed as  $\mathbf{X} | (Y = y)$ . The above argument motivates the inverse regression model

$$\mathbf{X}_y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\mathbf{v}_y + \boldsymbol{\epsilon}, \quad (2.2)$$

where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^\top \in \mathbb{R}^p$ ,  $\mathbf{v}_y \in \mathbb{R}^d$  is an unknown vector-valued function of  $y$ ,  $\boldsymbol{\epsilon}$  is a  $p$ -dimensional random vector with mean vector zero and

## 2.1 Reduction via classical multidimensional scaling

covariance matrix  $\Delta$ , and  $\epsilon$  is independent of  $Y$ . Since  $\Gamma$  is not generally identifiable, we require that  $\Delta^{-1/2}\Gamma$  be a  $p \times d$  matrix with ortho-normal columns, that is,  $\Gamma^\top \Delta^{-1}\Gamma$  is the  $d \times d$  identity matrix. Let  $\mu_y = E(\mathbf{X}_y) = E(\mathbf{X} | Y = y)$ . Then  $\mathbf{v}_y = \Gamma^\top \Delta^{-1}(\mu_y - \mu)$ . We assume that  $\text{Var}(\mathbf{v}_Y)$  is positive definite.

### 2.1 Reduction via classical multidimensional scaling

Assume for the moment that  $\Delta$  is known. Without loss of generality, assume that  $\Delta = \mathbf{I}_p$ , the  $p \times p$  identity matrix. This implies that  $\Gamma$  is a semi-orthogonal matrix, and  $\mathcal{S}_{Y|X} = \text{span}(\Gamma)$ . Otherwise, multiply both sides of equation (2.2) by  $\Delta^{-1/2}$  and replace  $(\mathbf{X}_y, \mu, \Gamma, \epsilon, \mathcal{S}_{Y|X})$  by  $(\Delta^{-1/2}\mathbf{X}_y, \Delta^{-1/2}\mu, \Delta^{-1/2}\Gamma, \Delta^{-1/2}\epsilon, \Delta^{1/2}\mathcal{S}_{Y|X})$ .

Suppose that the data consist of  $n$  independent observations,  $\mathbf{x}_{y_1}, \dots, \mathbf{x}_{y_n}$ .

For two observations indexed by  $y$  and  $y'$ , define  $d_{yy'} = \|\mu_y - \mu_{y'}\|_2^2$ . We have

$$d_{yy'} = \|\Gamma\mathbf{v}_y - \Gamma\mathbf{v}_{y'}\|_2^2 = \mathbf{v}_y^\top \mathbf{v}_y - 2\mathbf{v}_y^\top \mathbf{v}_{y'} + \mathbf{v}_{y'}^\top \mathbf{v}_{y'}.$$

Let  $\mathbf{D} = (d_{yy'}) \in \mathbb{R}^{n \times n}$ ,  $\mathbf{w} = (\mathbf{v}_{y_1}^\top \mathbf{v}_{y_1}, \dots, \mathbf{v}_{y_n}^\top \mathbf{v}_{y_n})^\top \in \mathbb{R}^n$ , and  $\mathbf{V} = (\mathbf{v}_{y_1}, \dots, \mathbf{v}_{y_n})^\top \in \mathbb{R}^{n \times d}$ . In matrix notation

$$\mathbf{D} = \mathbf{w}\mathbf{1}_n^\top + \mathbf{1}_n\mathbf{w}^\top - 2\mathbf{V}\mathbf{V}^\top,$$

## 2.1 Reduction via classical multidimensional scaling

where  $\mathbf{1}_n$  is the  $n$ -vector of ones. Let  $\mathbf{P}_n = \mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}_n^\top$ . Then

$$\mathbf{P}_n \mathbf{D} \mathbf{P}_n = -2\mathbf{P}_n \mathbf{V} \mathbf{V}^\top \mathbf{P}_n = -2\mathbf{V} \mathbf{V}^\top,$$

and hence

$$\mathbf{V} \mathbf{V}^\top = -\frac{1}{2} \mathbf{P}_n \mathbf{D} \mathbf{P}_n.$$

Here, without loss of generality, we assume that the columns of  $\mathbf{V}$  are centered, that is,  $\sum_{i=1}^n \mathbf{v}_{y_i}$  is the  $d$ -vector of zeros.

Since  $d_{yy'}$  is actually unknown, we replace it by  $\hat{d}_{yy'} = \|\mathbf{x}_y - \mathbf{x}_{y'}\|_2^2 - 2p$ .

It is easy to show that  $\hat{d}_{yy'}$  is an unbiased estimate of  $d_{yy'}$ . Let  $\hat{\mathbf{D}} = (\hat{d}_{yy'})$

and  $\mathbf{X} = (\mathbf{x}_{y_1}, \dots, \mathbf{x}_{y_n})^\top \in \mathbb{R}^{n \times p}$ . Then,

$$\mathbf{V} \mathbf{V}^\top \approx -\frac{1}{2} \mathbf{P}_n \hat{\mathbf{D}} \mathbf{P}_n = \mathbf{P}_n \mathbf{X} \mathbf{X}^\top \mathbf{P}_n.$$

Write the eigen-decomposition of  $\mathbf{P}_n \mathbf{X} \mathbf{X}^\top \mathbf{P}_n$  as

$$\mathbf{P}_n \mathbf{X} \mathbf{X}^\top \mathbf{P}_n = \sum_{i=1}^n \lambda_i \boldsymbol{\alpha}_i \boldsymbol{\alpha}_i^\top,$$

where  $\lambda_1 \geq \dots \geq \lambda_n \geq 0$  are the eigenvalues, and  $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n$  are the corresponding eigenvectors. By the Eckart–Young theorem, a solution for  $\mathbf{V}$  is given by

$$\tilde{\mathbf{V}} = (\lambda_1^{1/2} \boldsymbol{\alpha}_1, \dots, \lambda_d^{1/2} \boldsymbol{\alpha}_d).$$

Write  $\tilde{\mathbf{V}} = (\tilde{\mathbf{v}}_{y_1}, \dots, \tilde{\mathbf{v}}_{y_n})^\top$ . In the statistics literature, the reduction from  $\mathbf{x}_y$  to  $\tilde{\mathbf{v}}_y$  is known as the classical multidimensional scaling or principal

## 2.1 Reduction via classical multidimensional scaling

---

coordinate analysis. We can view  $\tilde{\mathbf{v}}_y$  as the vector of coordinates of  $\mathbf{x}_y$  with respect to the ortho-normal basis  $\mathbf{\Gamma}$ . From the viewpoint of dimension reduction in regression,  $\tilde{\mathbf{V}}$  then contains all the regression information on the response. In subsequent analyses, graphical displays and regression methods can be exploited to examine the relationship between the response and the vector of coordinates.

As such, it is important to predict the coordinates of a new observation,  $\mathbf{x}_{y^*}, y^* \in \mathbb{Y}$ . This can be done by the classical method of adding a point to vector diagrams (Gower, 1968; Zhang et al., 2008). For each  $i \in \{1, \dots, n\}$ , we define  $\tilde{s}_i = \|\tilde{\mathbf{v}}_{y_i}\|_2^2 - \|\mathbf{x}_{y^*} - \mathbf{x}_{y_i}\|_2^2$ . Let  $\tilde{\mathbf{s}} = (\tilde{s}_1, \dots, \tilde{s}_n)^\top \in \mathbb{R}^n$ . The predicted coordinates  $\tilde{\mathbf{v}}_{y^*}$  of  $\mathbf{x}_{y^*}$  is then

$$\tilde{\mathbf{v}}_{y^*} = \frac{1}{2}(\tilde{\mathbf{V}}^\top \tilde{\mathbf{V}})^{-1} \tilde{\mathbf{V}}^\top \tilde{\mathbf{s}}. \quad (2.3)$$

In the classical sufficient reduction, one is interested mainly in the matrix  $\mathbf{\Gamma}$  or the subspace  $\mathcal{S}_{Y|X}$  spanned by it. The above procedure operates in the space of coordinates of  $\mathbf{x}_y$  with respect to the ortho-normal basis  $\mathbf{\Gamma}$ . It achieves dimension reduction while at the same time avoiding the estimation of  $\mathbf{\Gamma}$ , and is thus appealing.

## 2.2 Subspace estimation

Once  $\mathbf{v}_y$  is found, it becomes natural to use least squares for estimating  $\mathbf{\Gamma}$  in model (2.2), if desired. Specifically, we estimate  $\mathbf{\Gamma}$  by minimizing the residual sum-of-squares

$$\|\mathbf{P}_n \mathbf{X} - \tilde{\mathbf{V}} \mathbf{\Gamma}^\top\|_F^2.$$

Here,  $\|\cdot\|_F$  denotes the Frobenius matrix norm. The minimizer can be shown to be

$$\tilde{\mathbf{\Gamma}} = \mathbf{X}^\top \tilde{\mathbf{V}} (\tilde{\mathbf{V}}^\top \tilde{\mathbf{V}})^{-1}. \quad (2.4)$$

Write  $\tilde{\mathbf{\Gamma}} = (\tilde{\mathbf{\Gamma}}_1, \dots, \tilde{\mathbf{\Gamma}}_d)$ . The estimate of  $\mathcal{S}_{Y|X}$  is then given by  $\text{span}(\tilde{\mathbf{\Gamma}})$ .

After some further manipulations,  $\tilde{\mathbf{\Gamma}}_j$  can be shown to equal the  $j$ th the principal component direction of  $\mathbf{P}_n \mathbf{X}$ , and thus the first  $d$  principal component score vectors of  $\mathbf{P}_n \mathbf{X}$  produces a sufficient reduction. Consequently, our method coincides with that of Cook (2007) under a PC regression model. The PC regression model is the same as the inverse regression model (2.2), except that  $\epsilon$  is assumed to be normally distributed, and the estimation under this model is by the method of maximum likelihood.

### 2.3 A toy example

Before we proceed, let us consider a simple simulation with  $p = 5$  and  $d = 2$ . Observations on  $(\mathbf{X}, Y)$  were generated from the inverse regression model (2.2) as follows. First,  $Y = y$  was sampled from a normal distribution with mean 0 and variance 4. Then,  $\mathbf{X}_y = \mathbf{x}_y$  was generated according to  $\mathbf{X}_y = \mathbf{\Gamma} \mathbf{v}_y + \boldsymbol{\epsilon}$ , where  $\mathbf{\Gamma} = (\mathbf{\Gamma}_1, \mathbf{\Gamma}_2)^\top$  with  $\mathbf{\Gamma}_1 = (1, 0, 0, 0, 0)^\top$  and  $\mathbf{\Gamma}_2 = (0, 1, 0, 0, 0)^\top$ ,  $\mathbf{v}_y = (y, y^2/3)^\top$ , and the error vector was sampled from a normal distribution with mean vector  $\mathbf{0}$  and covariance matrix  $\mathbf{\Delta} = \text{diag}(1, 1, 5, 5, 5)$ .

In the upper plot of Figure 1, the 2-dimensional coordinates of 200 samples from classical multidimensional scaling (CMDS, or principal component analysis by duality) are displayed, with each sample colored according to the response value. There appears to be little discernible relationship between the response and the coordinates (i.e., principal component scores). This lack of pattern is to be expected: aside from the subscript  $y$ , nothing on the right-hand side of (2.2) is observable. Consequently, dimension reduction under this model is based solely on the predictors, and hence is unsupervised. The lower plot shows the results of applying the supervised reduction method in Section 4. We see that the response increases as we move from the left to the right, and that the middle and extreme of response

### 3 A SUPERVISED INVERSE REGRESSION MODEL

---

values are somewhat separated by the second coordinate. In other words, some proportion of variability in the response can be explained using the new coordinates.

As was the case in this toy example, in many applications the response is expected to play an important role in supervising our reduction. Indeed, this is the main motivation for most modern dimension-reduction methods, including those developed in the framework of sufficient dimension reduction. We elaborate on this in the next section.

#### 3. A supervised inverse regression model

To facilitate supervised reduction, we can model the coordinate vectors as

$$\mathbf{v}_y = \boldsymbol{\beta} \mathbf{f}_y,$$

where  $\boldsymbol{\beta} \in \mathbb{R}^{d \times r}$  has rank  $d \leq \min(r, p)$ , and  $\mathbf{f}_y \in \mathbb{R}^r$  is a known vector-valued function of  $y$ . Usually,  $\mathbf{f}_y$  is required to contain a reasonably flexible set of basis functions, like slice indicator functions or B-spline basis functions.

This parameterization has been widely used in model-based reduction; see, for example, Cook and Forzani (2008), Cook et al. (2012), and Wang and Zhu (2013). Replacing  $\mathbf{v}_y$  in model (2.2) by  $\boldsymbol{\beta} \mathbf{f}_y$ , we are led to the following

### 3 A SUPERVISED INVERSE REGRESSION MODEL

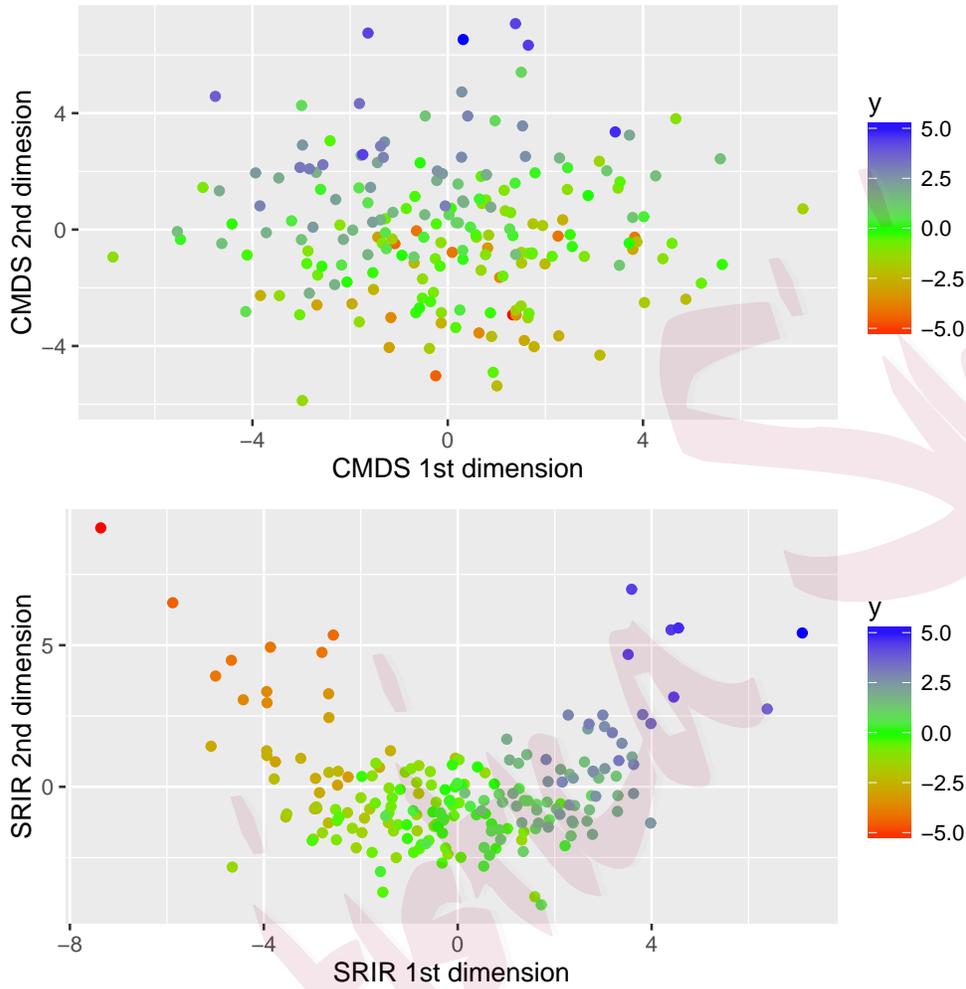


Figure 1: 2-dimensional plots of 200 samples from the inverse regression model (2.2). The simulation setup is described in the text. Top: The axes represent the first and second CMDS coordinates. Bottom: The axes represent the first and second coordinates produced by the supervised reduction method in Section 4.

### 3 A SUPERVISED INVERSE REGRESSION MODEL

model

$$\mathbf{X}_y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\beta}\mathbf{f}_y + \boldsymbol{\epsilon}. \quad (3.5)$$

We will refer to this model as a supervised inverse regression model. Without loss of generality, we assume that the sample mean vector of  $\mathbf{f}_y$  is zero.

The process of dimension reduction based on classical multidimensional scaling is essentially the same as before. Note that, under (3.5),

$$d_{yy'} = \|\boldsymbol{\Gamma}\boldsymbol{\beta}\mathbf{f}_y - \boldsymbol{\Gamma}\boldsymbol{\beta}\mathbf{f}_{y'}\|_2^2 = \mathbf{f}_y^\top \boldsymbol{\beta}^\top \boldsymbol{\beta} \mathbf{f}_y - 2\mathbf{f}_y^\top \boldsymbol{\beta}^\top \boldsymbol{\beta} \mathbf{f}_{y'} + \mathbf{f}_{y'}^\top \boldsymbol{\beta}^\top \boldsymbol{\beta} \mathbf{f}_{y'}.$$

Let  $\boldsymbol{\pi} = (\mathbf{f}_{y_1}^\top \boldsymbol{\beta}^\top \boldsymbol{\beta} \mathbf{f}_{y_1}, \dots, \mathbf{f}_{y_n}^\top \boldsymbol{\beta}^\top \boldsymbol{\beta} \mathbf{f}_{y_n})^\top \in \mathbb{R}^n$  and  $\mathbf{F} = (\mathbf{f}_{y_1}, \dots, \mathbf{f}_{y_n})^\top \in \mathbb{R}^{n \times r}$ . In matrix form

$$\mathbf{D} = \boldsymbol{\pi} \mathbf{1}_n^\top + \mathbf{1}_n \boldsymbol{\pi}^\top - 2\mathbf{F}\boldsymbol{\beta}^\top \boldsymbol{\beta} \mathbf{F}^\top.$$

Since  $\mathbf{P}_n \mathbf{F} = \mathbf{F}$ , a simple calculation shows that

$$\boldsymbol{\beta}^\top \boldsymbol{\beta} = -\frac{1}{2}(\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{F}^\top \mathbf{D} \mathbf{F}(\mathbf{F}^\top \mathbf{F})^{-1}.$$

Substituting  $\mathbf{D}$  by  $\hat{\mathbf{D}}$ , we arrive at

$$\begin{aligned} (\mathbf{F}^\top \mathbf{F})^{1/2} \boldsymbol{\beta}^\top \boldsymbol{\beta} (\mathbf{F}^\top \mathbf{F})^{1/2} &\approx -\frac{1}{2}(\mathbf{F}^\top \mathbf{F})^{-1/2} \mathbf{F}^\top \hat{\mathbf{D}} \mathbf{F} (\mathbf{F}^\top \mathbf{F})^{-1/2} \\ &= (\mathbf{F}^\top \mathbf{F})^{-1/2} \mathbf{F}^\top \mathbf{X} \mathbf{X}^\top \mathbf{F} (\mathbf{F}^\top \mathbf{F})^{-1/2}. \end{aligned}$$

Write the eigen-decomposition of the term in the last line as

$$(\mathbf{F}^\top \mathbf{F})^{-1/2} \mathbf{F}^\top \mathbf{X} \mathbf{X}^\top \mathbf{F} (\mathbf{F}^\top \mathbf{F})^{-1/2} = \sum_{j=1}^r \rho_j \boldsymbol{\phi}_j \boldsymbol{\phi}_j^\top,$$

## 4 REDUCTION WHEN $\Delta$ IS UNKNOWN

where  $\rho_1 \geq \dots \geq \rho_r \geq 0$  are the eigenvalues, and  $\phi_1, \dots, \phi_r$  are the corresponding eigenvectors. A solution for  $\beta$  is then given by

$$\tilde{\beta} = (\rho_1^{1/2} \phi_1, \dots, \rho_d^{1/2} \phi_d)^\top (\mathbf{F}^\top \mathbf{F})^{-1/2}.$$

Furthermore,

$$\tilde{\mathbf{v}}_y = \tilde{\beta} \mathbf{f}_y,$$

and the vector of coordinates  $\tilde{\mathbf{v}}_{y^*}$  associated with a new observation  $\mathbf{x}_{y^*}$  is, again, computed from (2.3).

### 4. Reduction when $\Delta$ is unknown

#### 4.1 The proposed method

In practice,  $\Delta$  is seldom known in advance, and has to be estimated from the data. Throughout the paper, we estimate  $\Delta$  by the residual sample covariance matrix from the multivariate linear regression of  $\mathbf{X}_y$  on  $\mathbf{f}_y$ :

$$\hat{\Delta} = \frac{1}{n} \mathbf{X}^\top (\mathbf{I}_n - \mathbf{P}_{\mathbf{F}}) \mathbf{X},$$

where  $\mathbf{P}_{\mathbf{F}} = \mathbf{F}(\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{F}^\top$  is the hat matrix. Asymptotic properties of  $\hat{\Delta}$  can be found in Lemmas 1 and 2 in the Supplementary Materials. Theorem 3.1 of Cook and Forzani (2008) shows that this estimator and the maximum likelihood estimator under normality of errors are different but related.

#### 4.1 The proposed method

We fix  $\Delta$  at  $\hat{\Delta}$ , and base the analysis on the standardized data  $\mathbf{X}\hat{\Delta}^{-1/2}$ .

For simplicity, let us focus on the supervised inverse regression model (3.5).

Replacing  $\mathbf{X}$  by  $\mathbf{X}\hat{\Delta}^{-1/2}$ , we compute

$$(\mathbf{F}^\top \mathbf{F})^{-1/2} \mathbf{F}^\top \mathbf{X} \hat{\Delta}^{-1} \mathbf{X}^\top \mathbf{F} (\mathbf{F}^\top \mathbf{F})^{-1/2}$$

and its eigen-decomposition

$$\sum_{j=1}^r \hat{\rho}_j \hat{\phi}_j \hat{\phi}_j^\top.$$

We estimate  $\beta$  and  $\mathbf{v}_y$  by

$$\hat{\beta} = (\hat{\rho}_1^{1/2} \hat{\phi}_1, \dots, \hat{\rho}_d^{1/2} \hat{\phi}_d)^\top (\mathbf{F}^\top \mathbf{F})^{-1/2}$$

and

$$\hat{\mathbf{v}}_y = \hat{\beta} \mathbf{f}_y.$$

Define  $\hat{s}_i = \|\hat{\mathbf{v}}_{y_i}\|_2^2 - \|\hat{\Delta}^{-1/2}(\mathbf{x}_{y^*} - \mathbf{x}_{y_i})\|_2^2$ . Let  $\hat{\mathbf{s}} = (\hat{s}_1, \dots, \hat{s}_n)^\top$  and

$\hat{\mathbf{V}} = (\hat{\mathbf{v}}_{y_1}, \dots, \hat{\mathbf{v}}_{y_n})^\top$ . Then, the vector of coordinates of a new observation  $\mathbf{x}_{y^*}$  is predicted by

$$\hat{\mathbf{v}}_{y^*} = \frac{1}{2} (\hat{\mathbf{V}}^\top \hat{\mathbf{V}})^{-1} \hat{\mathbf{V}}^\top \hat{\mathbf{s}}. \quad (4.6)$$

We call our method supervised reduction via inverse regression (SRIR).

As mentioned earlier, the advantage of working with coordinate vectors is that reduction can be done without the need to estimate  $\Gamma$  or  $\mathcal{S}_{Y|X}$ .

## 4.2 Theoretical properties

Nevertheless, there are situations in which the inferential target is  $\mathcal{S}_{Y|X}$ , as is the case in the traditional inquiry of sufficient dimension reduction. To conduct dimension reduction in the original predictor space, we have to determine both  $\Delta$  and  $\Gamma$ , and it is generally infeasible to find a closed-form expression for these estimators; usually an alternating procedure is needed. Fortunately, the estimate  $\hat{\Delta}$  has nothing to do with  $\Gamma$ , suggesting a one-step estimate for  $\Gamma$ . Specifically, we estimate  $\Gamma$  by minimizing the residual sum-of-squares

$$\begin{aligned} \text{RRS}(\Gamma) &= \|\mathbf{P}_n \mathbf{X} \hat{\Delta}^{-1/2} - \hat{\mathbf{V}} \Gamma^\top \hat{\Delta}^{-1/2}\|_F^2 \\ &= \text{trace}\{(\mathbf{P}_n \mathbf{X} - \hat{\mathbf{V}} \Gamma^\top) \hat{\Delta}^{-1} (\mathbf{P}_n \mathbf{X} - \hat{\mathbf{V}} \Gamma^\top)^\top\}. \end{aligned}$$

The solution is

$$\hat{\Gamma} = \mathbf{X}^\top \hat{\mathbf{V}} (\hat{\mathbf{V}}^\top \hat{\mathbf{V}})^{-1} = \mathbf{X}^\top \mathbf{F} \hat{\beta}^\top (\hat{\beta} \mathbf{F}^\top \mathbf{F} \hat{\beta}^\top)^{-1}. \quad (4.7)$$

Note that  $\hat{\Gamma}$  depends on  $\hat{\Delta}$  (and hence the observed responses  $y_i$ ) through  $\hat{\beta}$ . Finally, we estimate  $\mathcal{S}_{Y|X}$  by  $\text{span}(\hat{\Delta}^{-1} \hat{\Gamma})$ .

## 4.2 Theoretical properties

The limiting behavior of  $\hat{\boldsymbol{v}}_{y^*}$  is considered in the following theorem.

**Theorem 1.** *Assume that  $\boldsymbol{v}_Y = \beta \mathbf{f}_Y$  has finite sixth moments, and that  $\epsilon$*

## 4.2 Theoretical properties

has finite fourth moments. Then, for some  $d \times d$  rotation matrix  $\mathbf{R}$ ,

$$\hat{\mathbf{v}}_{y^*} = \mathbf{R}(\mathbf{v}_{y^*} + \mathbf{\Gamma}^\top \mathbf{\Delta}^{-1} \boldsymbol{\epsilon}_{y^*}) + O_P\left(\frac{1}{\sqrt{n}}\right).$$

For two  $d$ -dimensional random vectors  $\mathbf{V}_1$  and  $\mathbf{V}_2$ , let  $\boldsymbol{\Sigma}_1$ ,  $\boldsymbol{\Sigma}_2$ , and  $\boldsymbol{\Sigma}_{12}$  denote the covariance matrix of  $\mathbf{V}_1$ , the covariance matrix of  $\mathbf{V}_2$ , and the covariance matrix between  $\mathbf{V}_1$  and  $\mathbf{V}_2$ , respectively. To assess the prediction accuracy, we use the multiple correlation coefficient, which is defined as the positive square root of

$$\rho^2(\mathbf{V}_1, \mathbf{V}_2) = \frac{1}{d} \text{trace}(\boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_{12}^\top \boldsymbol{\Sigma}_1^{-1}).$$

This measure takes the maximum value of 1 if  $\mathbf{V}_1$  and  $\mathbf{V}_2$  are linearly related, and takes the minimum 0 if the components of the two vectors are uncorrelated; see Hall and Mathiason (1990) and Li and Dong (2009). We have the following corollary.

**Corollary 1.** *Assume the conditions of Theorem 1. Then,*

$$\rho^2(\hat{\mathbf{v}}_{Y^*}, \mathbf{v}_{Y^*}) = \frac{1}{d} \text{trace}[\text{Var}(\mathbf{v}_{Y^*})\{\text{Var}(\mathbf{v}_{Y^*}) + \mathbf{I}_d\}^{-1}] + O_P\left(\frac{1}{\sqrt{n}}\right),$$

where the covariances in  $\rho^2(\hat{\mathbf{v}}_{Y^*}, \mathbf{v}_{Y^*})$  are taken with respect to the joint distribution of  $Y^*$  and  $\boldsymbol{\epsilon}_{Y^*}$ .

We now consider the situation in which  $\mathbf{f}_y$  is misspecified. Denote by

$$\{\text{Var}(\mathbf{f}_Y)\}^{-1} \text{Cov}(\mathbf{f}_Y, \mathbf{v}_Y) = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^\top$$

## 4.2 Theoretical properties

the singular value decomposition of  $\{\text{Var}(\mathbf{f}_Y)\}^{-1}\text{Cov}(\mathbf{f}_Y, \mathbf{v}_Y)$ ; that is,  $\mathbf{U} = (\mathbf{U}_1, \dots, \mathbf{U}_d)$  is  $r \times d$  with orthonormal columns,  $\mathbf{V} = (\mathbf{V}_1, \dots, \mathbf{V}_d)$  is  $d \times d$  orthogonal, and  $\mathbf{\Lambda}$  is a  $d \times d$  diagonal matrix with diagonal entries  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ . Let  $\mathbf{\Phi} = (\lambda_1 \mathbf{U}_1, \dots, \lambda_d \mathbf{U}_d)^\top$ .

**Theorem 2.** *Assume that  $\mathbf{f}_Y$  has finite sixth moments, and that  $\mathbf{v}_Y$  and  $\boldsymbol{\epsilon}$  both have finite fourth moments. If  $\text{Cov}(\mathbf{f}_Y, \mathbf{v}_Y)$  has full column rank, that is,  $\lambda_d > 0$ , then*

$$\hat{\mathbf{v}}_{y^*} = \mathbf{R}(\mathbf{c} + \mathbf{A}\mathbf{v}_{y^*} + \mathbf{A}\mathbf{\Gamma}^\top \mathbf{\Omega}^{-1} \boldsymbol{\epsilon}_{y^*}) + O_P\left(\frac{1}{\sqrt{n}}\right),$$

for some  $d \times d$  rotation matrix  $\mathbf{R}$ , where

$$\mathbf{c} = \frac{1}{2} \{ \mathbf{\Phi} \text{Var}(\mathbf{f}_Y) \mathbf{\Phi}^\top \}^{-1} \{ \mathbf{E}(\mathbf{\Phi} \mathbf{f}_Y \mathbf{f}_Y^\top \mathbf{\Phi}^\top \mathbf{\Phi} \mathbf{f}_Y) - \mathbf{E}(\mathbf{\Phi} \mathbf{f}_Y \mathbf{v}_Y^\top \mathbf{v}_Y) \},$$

$$\mathbf{A} = \{ \mathbf{\Phi} \text{Var}(\mathbf{f}_Y) \mathbf{\Phi}^\top \}^{-1} \mathbf{\Phi} \text{Cov}(\mathbf{f}_Y, \mathbf{v}_Y),$$

and

$$\mathbf{\Omega} = \text{Var}(\mathbf{X}) - \mathbf{\Gamma} \text{Cov}(\mathbf{v}_Y, \mathbf{f}_Y) \{ \text{Var}(\mathbf{f}_Y) \}^{-1} \text{Cov}(\mathbf{f}_Y, \mathbf{v}_Y) \mathbf{\Gamma}^\top.$$

This result indicates that, up to an affine transformation, that is, a linear transformation followed by a translation, the conclusion of Theorem 1 remains valid, given that  $\mathbf{f}_Y$  and  $\mathbf{v}_Y$  are sufficiently correlated. From a dimension reduction point of view, we can treat  $\mathbf{v}_{y^*}$  and  $\mathbf{c} + \mathbf{A}\mathbf{v}_{y^*}$  as the same reduction.

The following theorem gives us consistency of subspace estimation.

**Theorem 3.** *Assume the conditions of Theorem 1 or Theorem 2. Then,  $\text{span}(\hat{\Delta}^{-1}\hat{\Gamma})$  is a  $\sqrt{n}$  consistent estimate of  $\mathcal{S}_{Y|X}$ .*

Let  $\hat{\Sigma} = \mathbf{X}^\top \mathbf{P}_F \mathbf{X}$ . Define  $\mathcal{S}_d(\hat{\Delta}, \hat{\Sigma})$  to be the span of  $\hat{\Delta}^{-1/2}$  times the first  $d$  eigenvectors of  $\hat{\Delta}^{-1/2} \hat{\Sigma} \hat{\Delta}^{-1/2}$ . One connection between our one-step subspace estimate and the maximum likelihood estimate is captured in the following theorem.

**Theorem 4.** *Assume that  $\mathbf{f}_y$  is correctly specified. Then,  $\text{span}(\hat{\Delta}^{-1}\hat{\Gamma}) = \mathcal{S}_d(\hat{\Delta}, \hat{\Sigma})$ . Consequently, if  $\epsilon$  is normally distributed, then  $\text{span}(\hat{\Delta}^{-1}\hat{\Gamma})$  is the maximum likelihood estimate of  $\mathcal{S}_{Y|X}$ .*

### 4.3 Choice of $d$

In practice, the structural dimension  $d$  is unknown, and the choice of  $d$  is essential to the proposed method. In the literature, there are two useful techniques for determining  $d$ : one is based on a sequential test (Li, 1991), and the other is by using an information criterion (Zhu et al., 2006). Let  $d_0$  denote the true dimension. To estimate  $d_0$ , we propose to use the Bayesian information criterion (Zhu et al., 2012). With

$$\text{BIC}_d = \frac{\sum_{j=1}^d \hat{\rho}_j^2}{\sum_{k=1}^r \hat{\rho}_k^2} - \frac{\log(n)}{n} \times \frac{d(d+1)}{2},$$

the estimated dimension is

$$\hat{d} = \arg \max_{1 \leq d \leq r} \text{BIC}_d. \quad (4.8)$$

**Theorem 5.** *Assume the conditions of Theorem 1 or Theorem 2. Then,  $\hat{d}$  converges to  $d_0$  in probability.*

## 5. Numerical studies

In this section, we first conduct Monte Carlo simulation studies to assess the finite sample performance of the proposed method. We then apply our method to the analysis of a real data set.

### 5.1 Simulations

Throughout the simulation study, we considered the structural dimension  $d = 2$ , the sample size  $n = 200$ , and the number of predictors  $p \in \{10, 20\}$ .

We set  $\mathbf{\Delta} = (\theta^{|i-j|})$ , with  $\theta$  taking 0 or 0.5. Let  $\mathbf{\Gamma}_{01} = (1, 1, -1, -1, 0, \dots, 0)^\top / 2$ ,  $\mathbf{\Gamma}_{02} = (1, 0, 1, 0, 1, 0, \dots, 0)^\top / \sqrt{3}$ , and  $\mathbf{\Gamma}_0 = (\mathbf{\Gamma}_{01}, \mathbf{\Gamma}_{02})$ . Set  $\mathbf{\Gamma} = \mathbf{\Gamma}_0(\mathbf{\Gamma}_0^\top \mathbf{\Delta}^{-1} \mathbf{\Gamma}_0)^{-1/2}$ .

We used the cubic polynomial basis  $(y, y^2, y^3)$  to fit the model, and then assessed the prediction accuracy on an independent test sample,  $\{(\mathbf{x}_{y_i^*}, y_i^*)\}$ , of size 100. Let  $\hat{\mathbf{v}}_{y_i^*}$  be the predicted vector of coordinates of  $\mathbf{x}_{y_i^*}$ . To measure the closeness between  $\hat{\mathbf{v}}_{y_i^*}$  and  $\mathbf{v}_{y_i^*}$ , we used the sample version

of the multiple correlation coefficient (MCC). For each configuration, the number of repetitions was 200.

**Example 1.** To gain first insights into the operating characteristics of the proposed method, consider the model

$$\mathbf{X}_y = \mathbf{\Gamma} \mathbf{v}_y + \boldsymbol{\epsilon},$$

where  $y$  is a draw from a normal distribution  $N(0, \sigma^2)$ ,  $\mathbf{v}_y = (y, y^2)^\top$ , and  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \boldsymbol{\Delta})$ . By Corollary 1,

$$\rho^2(\hat{\mathbf{v}}_{Y^*}, \mathbf{v}_{Y^*}) = g^2(\sigma) + O_P\left(\frac{1}{\sqrt{n}}\right).$$

Here,  $g(\sigma) = \sqrt{\sigma^2/(2\sigma^2 + 2) + \sigma^4/(2\sigma^4 + 1)}$  is an increasing function of  $\sigma$ . Six values of  $\sigma$  were explored: 0.5, 0.8, 1, 1.5, 2, and 3. Figure 2 depicts  $g(\sigma)$  and its sample estimate under different configurations. We see that there is an excellent agreement between theoretical prediction and empirical behavior.

Below we examined the behavior of our method in more detail. In addition to the prediction accuracy, we also assessed the performance in terms of subspace recovery. Specifically, we used the vector correlation coefficient (VCC) and the trace correlation coefficient (TCC) to measure the closeness between the true and estimated subspaces (Ye and Weiss, 2003). Let  $\hat{\mathbf{B}}$  and  $\mathbf{B}$  be basis matrices for the estimated and true subspaces,

## 5.1 Simulations

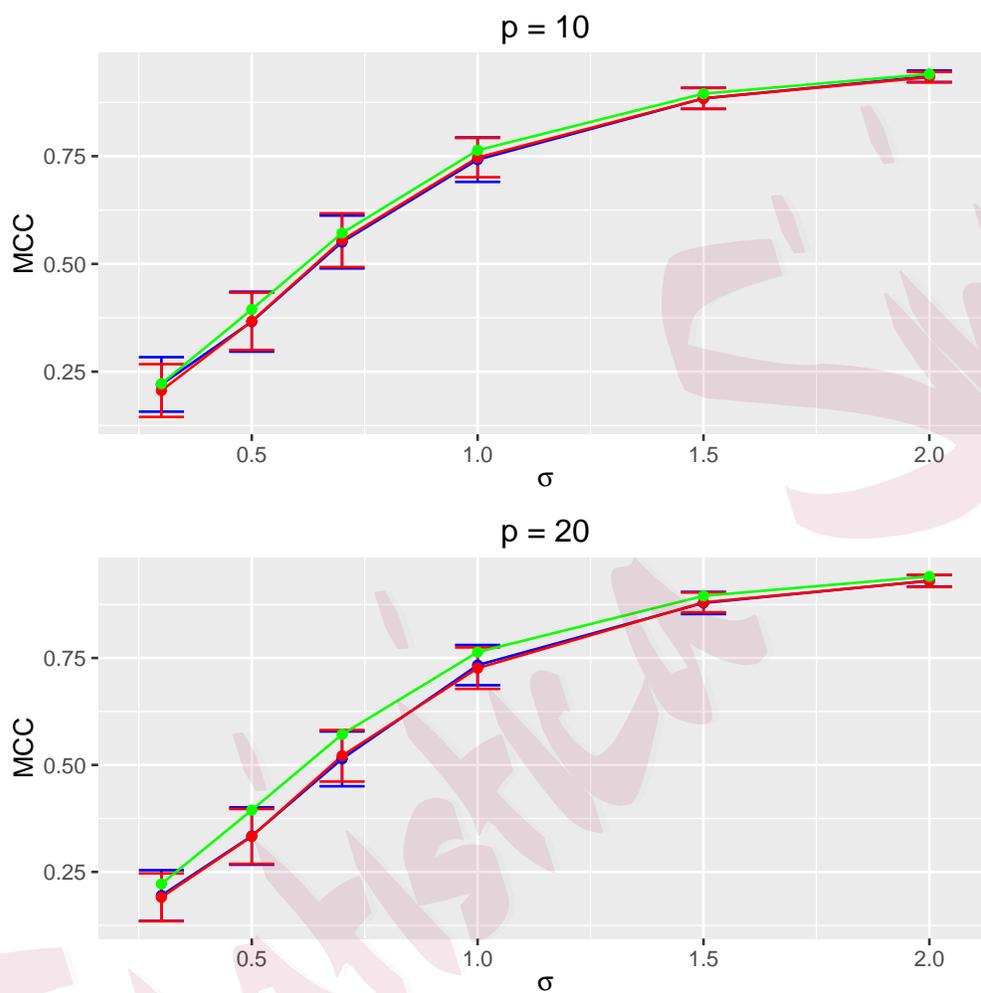


Figure 2: The estimated MCC curves, based on 200 data replications, for  $\theta = 0$  (red) and  $\theta = 0.5$  (blue). The error bars indicate one standard deviation. The theoretical MCC curve  $g(\sigma)$  is shown in green.

respectively. Denote by  $\rho_1 \geq \dots \geq \rho_d$  the ordered eigenvalues of  $\hat{\mathbf{B}}^\top \mathbf{B} \mathbf{B}^\top \hat{\mathbf{B}}$ .

VCC is defined as the square positive root of  $q^2(\hat{\mathbf{B}}, \mathbf{B}) = \prod_{i=1}^d \rho_i$ , and TCC

is defined as the square positive root of  $r^2(\hat{\mathbf{B}}, \mathbf{B}) = d^{-1} \sum_{i=1}^d \rho_i$ .

**Example 2.** Consider the model

$$\mathbf{X}_y = \mathbf{\Gamma} \boldsymbol{\beta} \mathbf{f}_y + \boldsymbol{\epsilon},$$

where  $\mathbf{f}_y = (y, |y|, y^2)^\top$ , and  $\boldsymbol{\epsilon}$  has mean vector  $\mathbf{0}$  and covariance matrix  $\boldsymbol{\Delta}$ .

Two types of non-Gaussian errors were explored, with covariance structures

the same as in the previous example. In the former,  $\boldsymbol{\epsilon}$  is from multivariate

$t$ -distribution with 5 degrees of freedom, and in the latter each component

of  $\boldsymbol{\epsilon}$  is uniformly distributed on  $[-\sqrt{3}, \sqrt{3}]$ . For the coefficient matrix  $\boldsymbol{\beta}$ ,

we set

$$\boldsymbol{\beta} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \text{ or } \begin{pmatrix} 1 & -0.5 & 0 \\ 0 & 0.5 & 1 \end{pmatrix}.$$

This corresponds to the setting where the cubic polynomial basis is correctly

specified or misspecified. Finally, we generated  $Y$  from the standard normal

distribution. The averaged values of MCC, VCC, and TCC, and their

standard deviations, based on 200 data replications, are reported in Tables

1-4. From Tables 1 and 3 we see that the prediction accuracy for non-

Gaussian errors is comparable to that under the Gaussian assumption (Figure

2,  $\sigma = 1$ ). Furthermore, our method performs well in terms of subspace

## 5.1 Simulations

estimation. Generally, the performance improves as the number of predictors decreases. From Tables 2 and 4, we see that our method is robust to misspecification of the basis functions, as expected from Theorems 2 and 3.

Table 1: Means and standard deviations (in parentheses) of MCC, VCC, and TCC, over 200 data applications.  $\epsilon$  is from multivariate  $t$ -distribution with 5 degrees of freedom and  $\mathbf{f}_y$  is correctly specified.

		SRIR			PC		
		MCC	VCC	TCC	MCC	VCC	TCC
$p = 10$	$\theta = 0$	0.737 (0.056)	0.900 (0.036)	0.949 (0.018)	0.725 (0.066)	0.845 (0.170)	0.929 (0.061)
	$\theta = 0.5$	0.740 (0.048)	0.916 (0.037)	0.958 (0.018)	0.456 (0.085)	0.120 (0.083)	0.406 (0.096)
$p = 20$	$\theta = 0$	0.716 (0.055)	0.802 (0.047)	0.897 (0.025)	0.667 (0.081)	0.581 (0.279)	0.816 (0.103)
	$\theta = 0.5$	0.728 (0.055)	0.827 (0.048)	0.911 (0.025)	0.374 (0.085)	0.050 (0.049)	0.304 (0.077)

Table 2: Means and standard deviations (in parentheses) of MCC, VCC, and TCC, over 200 data applications.  $\epsilon$  is from multivariate  $t$ -distribution with 5 degrees of freedom and  $\mathbf{f}_y$  is misspecified.

		SRIR			PC		
		MCC	VCC	TCC	MCC	VCC	TCC
$p = 10$	$\theta = 0$	0.768 (0.046)	0.899 (0.036)	0.949 (0.018)	0.758 (0.057)	0.843 (0.187)	0.930 (0.067)
	$\theta = 0.5$	0.767 (0.050)	0.917 (0.039)	0.958 (0.019)	0.502 (0.069)	0.130 (0.098)	0.431 (0.082)
$p = 20$	$\theta = 0$	0.752 (0.052)	0.815 (0.041)	0.903 (0.022)	0.711 (0.070)	0.607 (0.257)	0.834 (0.089)
	$\theta = 0.5$	0.752 (0.048)	0.826 (0.049)	0.911 (0.025)	0.457 (0.071)	0.058 (0.049)	0.353 (0.065)

## 5.1 Simulations

Table 3: Means and standard deviations (in parentheses) of MCC, VCC, and TCC, over 200 data applications. Each component of  $\epsilon$  is uniformly distributed on  $[-\sqrt{3}, \sqrt{3}]$  and  $\mathbf{f}_y$  is correctly specified.

		SRIR			PC		
		MCC	VCC	TCC	MCC	VCC	TCC
$p = 10$	$\theta = 0$	0.745 (0.044)	0.931 (0.024)	0.965 (0.012)	0.747 (0.043)	0.946 (0.028)	0.973 (0.014)
	$\theta = 0.5$	0.742 (0.048)	0.942 (0.027)	0.971 (0.013)	0.436 (0.075)	0.097 (0.073)	0.377 (0.088)
$p = 20$	$\theta = 0$	0.727 (0.048)	0.852 (0.034)	0.923 (0.018)	0.731 (0.049)	0.867 (0.052)	0.933 (0.025)
	$\theta = 0.5$	0.730 (0.049)	0.859 (0.042)	0.928 (0.021)	0.380 (0.074)	0.040 (0.037)	0.290 (0.066)

Table 4: Means and standard deviations (in parentheses) of MCC, VCC, and TCC, over 200 data applications. Each component of  $\epsilon$  is uniformly distributed on  $[-\sqrt{3}, \sqrt{3}]$  and  $\mathbf{f}_y$  is misspecified.

		SRIR			PC		
		MCC	VCC	TCC	MCC	VCC	TCC
$p = 10$	$\theta = 0$	0.776 (0.033)	0.934 (0.022)	0.967 (0.011)	0.777 (0.033)	0.944 (0.028)	0.972 (0.014)
	$\theta = 0.5$	0.769 (0.036)	0.944 (0.026)	0.971 (0.012)	0.488 (0.060)	0.111 (0.078)	0.411 (0.066)
$p = 20$	$\theta = 0$	0.753 (0.037)	0.853 (0.034)	0.924 (0.018)	0.757 (0.039)	0.881 (0.049)	0.941 (0.024)
	$\theta = 0.5$	0.753 (0.045)	0.869 (0.037)	0.933 (0.019)	0.449 (0.065)	0.049 (0.045)	0.354 (0.058)

We also compared our method with principal components (PC) and principal fitted components (PFC) in Cook and Forzani (2008). The results of PC are shown in the last three columns of Tables 1-4. SRIR appears to

## 5.1 Simulations

dominate PC in most cases, especially when  $\theta = 0.5$ . The results of PFC are the same as for SRIR and are omitted. For subspace estimation, this is expected from Theorem 4, but for prediction this comes as somewhat of a surprise. We provide a theoretical support in the Supplementary Materials.

So far the value of the structural dimension is assumed to be known. Using Example 2, we evaluated the performance of the BIC-type criterion (4.8). Tables 5 and 6 report the frequencies of decisions over 200 replications. We see that the proportion of correctly identifying the true dimension is larger than 80% in each configuration. We also see that misspecification can have some impact.

Table 5: Selection frequencies of BIC over 200 data replications.  $\epsilon$  is from multivariate  $t$ -distribution with 5 degrees of freedom.

		Correctly specified		Misspecified	
		$\hat{d} = 1$	$\hat{d} = 2$	$\hat{d} = 1$	$\hat{d} = 2$
$p = 10$	$\theta = 0$	2	198	39	161
	$\theta = 0.5$	2	198	31	169
$p = 20$	$\theta = 0$	5	195	33	167
	$\theta = 0.5$	4	196	30	170

## 5.2 Boston housing data

Table 6: Selection frequencies of BIC over 200 data replications. Each component of  $\epsilon$  is uniformly distributed on  $[-\sqrt{3}, \sqrt{3}]$ .

		Correctly specified		Misspecified	
		$\hat{d} = 1$	$\hat{d} = 2$	$\hat{d} = 1$	$\hat{d} = 2$
$p = 10$	$\theta = 0$	0	200	29	171
	$\theta = 0.5$	1	199	32	168
$p = 20$	$\theta = 0$	3	197	26	174
	$\theta = 0.5$	1	199	30	170

## 5.2 Boston housing data

We applied SRIR to the Boston housing data (Harrison and Rubinfeld, 1978), which is available in the **MASS** library in **R**. This data set has 14 variables and 506 observations, with each observation representing a census tract in Boston Standard Metropolitan Statistical Areas. The variable of primary interest is the median value, in thousands of dollars, of owner occupied homes. The 13 explanatory variables include per capita crime rate by town, average number of rooms per house, percent of households with low socioeconomic status, and so on.

Fitting the supervised inverse regression model (3.5), with the cubic polynomial basis, resulted in BIC choosing  $d = 2$ , suggesting that two linear combinations of the 13 predictors are sufficient. The top panel of Figure 3 shows the 2-dimensional plot of the 506 observations, with coordinates computed using the formula (4.6). We see a horseshoe-like pattern in

the data cloud. We also see an association between the response and the coordinates, similar to the one in the toy example. For comparison, we also carried out classical multidimensional scaling. In the bottom panel, the ordination of the first two CMDS coordinates is shown. The unsupervised method failed to show any useful relationship.

Figure 4 gives plots of the response versus the SRIR coordinates. The upper panel shows a strong linear relation between the response and the first SRIR coordinate. In the lower panel, we see a nonlinear association between the response and the second SRIR coordinate.

## 6. Discussion

Linear reduction methods aim to construct a few linear combinations of the original predictors that are useful for subsequent analyses. Nearly all existing methods estimate a subspace in the primal predictor-based space, and then obtain the set of reduced predictors by projecting the original predictor vector onto this subspace. In this paper, we have proposed a principled reduction method in the dual sample-based space, on the basis of a supervised inverse regression model. Instead of estimating the subspace, our method directly estimate the projection coordinates of the predictor vector onto the subspace. The results extend the well-known

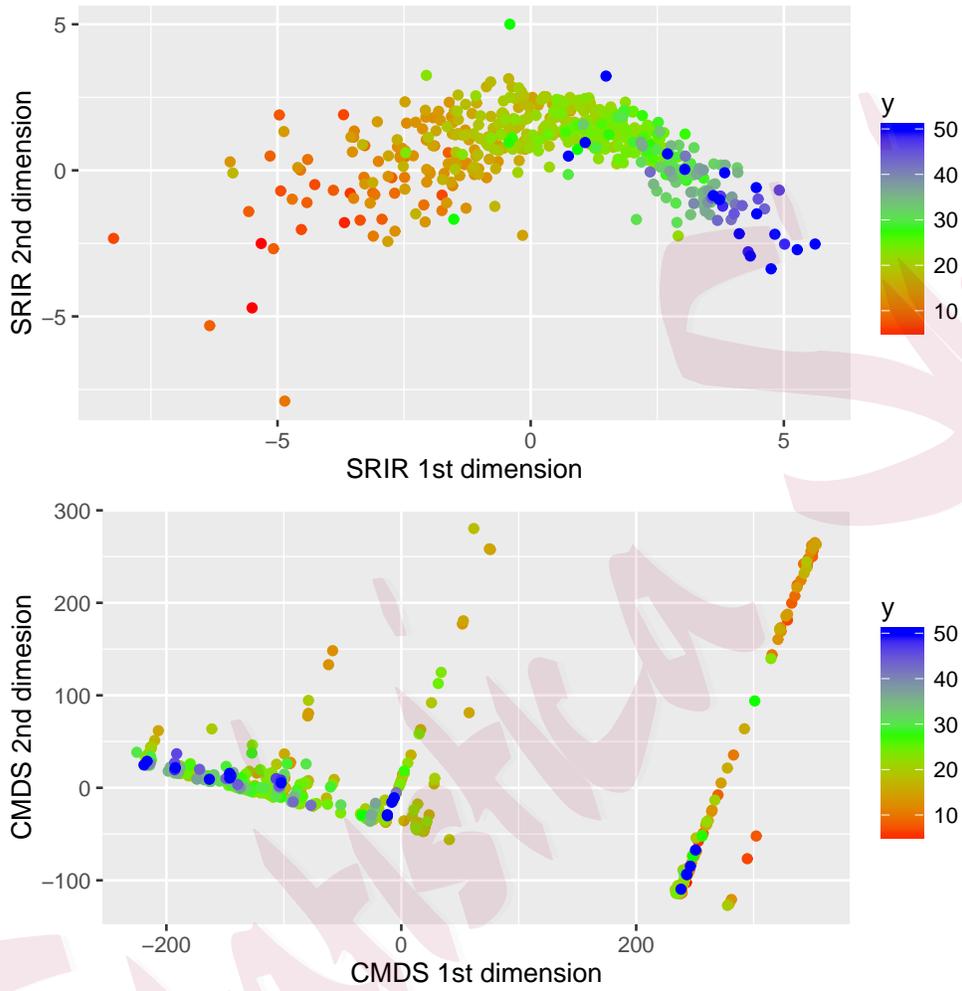


Figure 3: 2-dimensional plots for the Boston housing data. Top: The axes represent the first and second SRIR coordinates. Down: The axes represent the first and second CMDS coordinates.

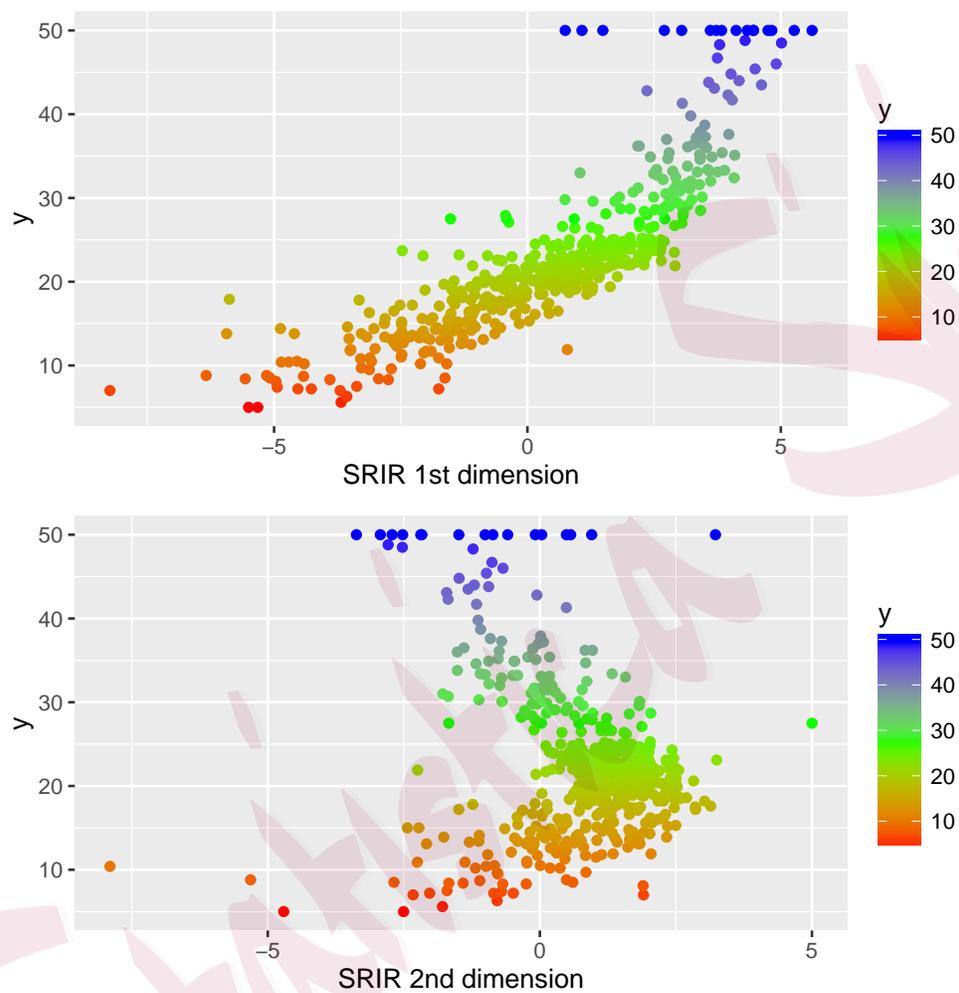


Figure 4: Boston housing data. Top: Response versus the first SRIR coordinate. Down: Response versus the second SRIR coordinate.

duality between principal component analysis and classical multidimensional scaling.

Computation of SRIR has the same order of computation as the maximum likelihood estimation (Cook and Forzani, 2008). However, in terms of generating the reduction for the observed data, our method has a smaller computational cost than the method of maximum likelihood. Specifically, the computational complexity of the former is  $O(d \times n \times r^2)$ , while that of the latter is  $O(d \times n \times p^2)$ .

We have studied the theoretical properties of our method, and used simulation results to support them. As with most reduction methods, we have adopted the traditional asymptotic reasoning by letting the sample size  $n \rightarrow \infty$ , with the number of predictors  $p$  fixed. Our method requires the inverse of the residual sample covariance matrix, and hence is problematic in situations where  $p$  is comparable to or even larger than  $n$ . Regularized versions in the dual space have a strong practical appeal and are currently under investigation.

Our method is related to a nonparametric multivariate analysis procedure in ecological studies (Mcardle and Anderson, 2001). This procedure, known as permutation multivariate analysis of variance, partitions the variability in multivariate ecological data according to factors in an experimental design.

The underlying intuition is the duality between  $\mathbf{X}^\top \mathbf{X}$ , an inner product matrix in the primal space, and  $\mathbf{X}\mathbf{X}^\top$ , an outer product matrix in the dual space, in the sense that  $\text{trace}(\mathbf{X}^\top \mathbf{X}) = \text{trace}(\mathbf{X}\mathbf{X}^\top)$ . This equivalence is important because an outer product matrix can be obtained from any symmetric distance matrix  $\mathbf{D} = (d_{ij}) \in \mathbb{R}^{n \times n}$  (Gower, 1966). In particular, for a  $p \times p$  positive definite matrix  $\mathbf{B}$ , if we let  $d_{ij}(\mathbf{B}) = (\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{B} (\mathbf{x}_i - \mathbf{x}_j)$ , then  $\mathbf{X}\mathbf{B}\mathbf{X}^\top = -\mathbf{P}_n \mathbf{D} \mathbf{P}_n / 2$ , where  $\mathbf{P}_n$  is the centering matrix. Similar to permutation multivariate analysis of variance, we can extend our supervised reduction method, based solely on measures of distance or dissimilarity between pairs of observations, even without assuming the inverse regression model. Alternatively, under a notion of nonlinear sufficient reduction (Zhang et al., 2008), it is possible to derive a kernel extension of the proposed method. Work along these lines is in progress.

**Supplementary Materials** The supplementary file contains the proofs.

**Acknowledgements** Tao Wang is the corresponding author. Peirong Xu was supported by the Natural Science Foundation of Shanghai (19ZR1437000). Tao Wang was supported in part by the National Natural Science Foundation of China (11601326), National Key R&D Program of China (2018YFC0910500),

## REFERENCES

---

Shanghai Municipal Science and Technology Major Project (2017SHZDZX01),  
and Neil Shen's SJTU Medical Research Fund.

### References

Adragni, K. P. and R. D. Cook (2009). Sufficient dimension reduction and prediction in regression. *Philosophical Transactions of the Royal Society A* 367(1906), 4385–4405.

Cook, R. D. (1998). *Regression Graphics: Ideas for Studying Regressions Through Graphics*. Wiley, New York.

Cook, R. D. (2007). Fisher lecture: Dimension reduction in regression. *Statistical Science* 22(1), 1–26.

Cook, R. D. and L. Forzani (2008). Principal fitted components for dimension reduction in regression. *Statistical Science* 23(4), 485–501.

Cook, R. D., L. Forzani, and A. J. Rothman (2012). Estimating sufficient reductions of the predictors in abundant high-dimensional regressions. *The Annals of Statistics* 40(1), 353–384.

Cook, R. D. and S. Weisberg (1991). Comment. *Journal of the American Statistical Association* 86(414), 328–332.

## REFERENCES

---

- Gower, J. C. (1966). Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika* 53, 325–338.
- Gower, J. C. (1968). Adding a point to vector diagrams in multivariate analysis. *Biometrika* 55(3), 582–585.
- Hall, W. J. and D. J. Mathiason (1990). On large-sample estimation and testing in parametric models. *International Statistical Review* 58(1), 77–97.
- Harrison, D. and D. L. Rubinfeld (1978). Hedonic housing prices and the demand for clean air. *Journal of Environmental Economics and Management* 5(1), 81–102.
- Hastie, T., R. Tibshirani, and J. Friedman (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, New York.
- Li, B. and Y. Dong (2009). Dimension reduction for nonelliptically distributed predictors. *The Annals of Statistics* 37(3), 1272–1298.
- Li, B. and S. Wang (2007). On directional regression for dimension reduction. *Journal of the American Statistical Association* 102, 997–1008.

---

REFERENCES

- Li, K.-C. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association* 86(414), 316–327.
- Mcardle, B. H. and M. J. Anderson (2001). Fitting multivariate models to community data: A comment on distance-based redundancy analysis. *Ecology* 82(1), 290–297.
- Wang, T., M. Chen, H. Zhao, and L. Zhu (2018). Estimating a sparse reduction for general regression in high dimensions. *Statistics and Computing* 28(1), 33–46.
- Wang, T. and L. Zhu (2013). Sparse sufficient dimension reduction using optimal scoring. *Computational Statistics & Data Analysis* 57(1), 223–232.
- Ye, Z. and R. E. Weiss (2003). Using the bootstrap to select one of a new class of dimension reduction methods. *Journal of the American Statistical Association* 98(464), 968–979.
- Zhang, Z., D. Yeung, J. T. Kwok, and E. Y. Chang (2008). Sliced coordinate analysis for effective dimension reduction and nonlinear extensions. *Journal of Computational and Graphical Statistics* 17(1), 225–242.
- Zhu, L., B. Miao, and H. Peng (2006). On sliced inverse regression

## REFERENCES

---

with high-dimensional covariates. *Journal of the American Statistical Association* 101(474), 630–643.

Zhu, L., L. Zhu, and Z. Feng (2012). Dimension reduction in regressions through cumulative slicing estimation. *Journal of the American Statistical Association* 105(492), 1455–1466.

College of Mathematics and Sciences, Shanghai Normal University, Shanghai 200234, China

E-mail: prxu@shnu.edu.cn

Department of Bioinformatics and Biostatistics, Shanghai Jiao Tong University, Shanghai 200240, China

SJTU-Yale Joint Center for Biostatistics and Data Science, Shanghai Jiao Tong University, Shanghai 200240, China

MoE Key Lab of Artificial Intelligence, Shanghai Jiao Tong University, Shanghai 200240, China

E-mail: neowangtao@sjtu.edu.cn