Statistica Sinica Preprint No: SS-2017-0532						
Title	On supervised reduction and its dual					
Manuscript ID	SS-2017-0532					
URL	http://www.stat.sinica.edu.tw/statistica/					
DOI	10.5705/ss.202017.0532					
Complete List of Authors	Peirong Xu and					
	Tao Wang					
Corresponding Author	Tao Wang					
E-mail	neowangtao@sjtu.edu.cn					

ON SUPERVISED REDUCTION AND ITS DUAL

Peirong Xu and Tao Wang

Shanghai Normal University and Shanghai Jiao Tong University

Abstract: Existing regression dimension-reduction 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. We propose a principled method for estimating a sufficient reduction in the dual sample-based space, based on a supervised inverse regression model. The reduction is performed without needing 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 show how to apply the method by means of a 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 transform the predictors, and then fit a least squares model using the transformed variables. For example, the widely used principal component regression extracts the first few principal components of the predictors, and then uses these components as the predictors in a linear model. However, one of the main concerns with this approach is that the directions in which the predictors show the most variation are not necessarily the directions associated with the response. Many methods have been proposed to deal with this issue, including the partial least squares and sliced inverse regression (Li, 1991) methods. A common goal of such 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 $X \in \mathbb{R}^p$. Formally, the aim is to estimate a reduction $\mathcal{R} : \mathbb{R}^p \to \mathbb{R}^d$, for $d \leq p$, such that

$$Y \perp \!\!\!\perp \boldsymbol{X} \mid \mathcal{R}(\boldsymbol{X}), \tag{1.1}$$

where $\perp\!\!\!\perp$ indicates independence. Here, $\mathcal{R}(\boldsymbol{X})$ is called a sufficient reduction for the regression of Y onto \boldsymbol{X} (Cook, 1998). Sufficient dimension reduction has been an active research area since the introduction of the sliced inverse regression and sliced average variance estimation methods (Cook and Weisberg, 1991). In this study 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 X, there are three paradigms for determining a sufficient reduction: forward reduction, inverse reduction, and joint reduction, which are equivalent when Y and X are jointly distributed (Cook, 2007). Without requiring a pre-specified model for $Y \mid X$, inverse reduction is promising in regressions with many predictors. To estimate a reduction inversely, methods such as the sliced inverse regression exploit properties of the conditional moments of $X \mid Y$. These inverse momentbased methods impose constraints on the marginal distribution of X. Alternatively, inverse model-based approaches directly specify a model for the inverse regression of X onto Y. Much of the existing work relies on normal models. See Adragni and Cook (2009) for a recent review of inverse reduction methods.

Sufficient reduction permits us to restrict attention to a few new predictors $\eta^{\top} X$, 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 η . 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) proposed an inverse model-based method after studying 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 an inverse moment-based method for estimating sparse reductions using a novel representation of a sliced inverse regression.

We propose a new approach for estimating a sufficient reduction, motivated by the well-known duality between principal component analysis and principal coordinate analysis (Gower, 1966), also known as classical multidimensional scaling (CMDS; Hastie et al., 2009). Let $\mathbf{X} = (\mathbf{x}_1, \ldots, \mathbf{x}_n)^\top \in \mathbb{R}^{n \times p}$ be a data matrix of predictor values. Without loss of generality, assume that each column of \mathbf{X} has been centered to have mean zero. A singular value decomposition offers a way of expressing a principal component analysis. Let $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ be the singular value decomposition of \mathbf{X} ; that is, $\mathbf{U} =$ $(\mathbf{U}_1, \ldots, \mathbf{U}_p)$ is $n \times p$ with orthonormal columns, $\mathbf{V} = (\mathbf{V}_1, \ldots, \mathbf{V}_p)$ is $p \times p$ orthogonal, and \mathbf{D} is a $p \times p$ diagonal matrix with diagonal entries $d_1 \ge d_2 \ge$ $\cdots \ge d_p \ge 0$. Here, \mathbf{V}_j is called the *j*th principal component direction, and $d_j \mathbf{U}_j$ is the *j*th principal component score vector. In the terminology of regression dimension reduction, $\mathbf{UD} = \mathbf{XV}$ are linear reductions used in a 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, a principal component analysis can be viewed as an ordination method. Indeed, it is equivalent to CMDS, 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 a moment-based inverse reduction, Zhang et al. (2008) calculated projection coordinates by applying CMDS to slice means, and then interpolated the projection of a new predictor vector using these coordinates. This method can be thought of as a dual version of a sliced inverse regression. At the population level, however, it is not clear what quantity is being treated as the target. Here, we propose a principled method for estimating a sufficient reduction under the inverse model-based reduction scheme. The reduction is performed without needing to estimate the subspace. To the best of our knowledge, this study is the first time to examine 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. Because the coordinates are unconstrained, sufficient reduction is achieved using CMDS, or a principal component analysis by duality. To perform a supervised reduction, we model the coordinates in a parametric way in Section 3, extending the method of Section 2 for a known error structure. A reduction with an unknown error structure is considered in Section 4, and our theoretical conclusions are presented. Simulation results and a real-data application are presented in Section 5. Section 6 concludes the paper. All proofs are available in the online Supplementary Material.

2. A naive inverse regression model

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

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

denote the subspace spanned by the centered inverse regression curves. We have the following proposition.

Proposition 1. Assume (C1) $S_{E(X|Y)} \subseteq Var(X)S_{Y|X}$ and (C2) Var(X | Y) is positive definite and nonrandom. Then,

 $\operatorname{Var}(\boldsymbol{X} \mid Y) \mathcal{S}_{Y \mid \boldsymbol{X}} = \operatorname{Var}(\boldsymbol{X}) \mathcal{S}_{Y \mid \boldsymbol{X}}.$

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

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

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

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

$$\boldsymbol{X}_{y} = \boldsymbol{\mu} + \boldsymbol{\Gamma} \boldsymbol{v}_{y} + \boldsymbol{\epsilon}, \qquad (2.2)$$

where $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^\top \in \mathbb{R}^p, \ \boldsymbol{v}_y \in \mathbb{R}^d$ is an unknown vector-valued function of $y, \boldsymbol{\epsilon}$ is a *p*-dimensional random vector with mean vector zero and covariance matrix Δ , and ϵ is independent of Y. Because Γ is not usually identifiable, we require that $\Delta^{-1/2}\Gamma$ be a $p \times d$ matrix with orthonormal columns; that is, $\Gamma^{\top}\Delta^{-1}\Gamma$ is the $d \times d$ identity matrix. Let $\mu_y = E(X_y) =$ $E(X \mid Y = y)$. Then, $v_y = \Gamma^{\top}\Delta^{-1}(\mu_y - \mu)$. We assume that $Var(v_Y)$ is positive definite.

2.1 Reduction via CMDS

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

Suppose the data consist of n independent observations, $\boldsymbol{x}_{y_1}, \ldots, \boldsymbol{x}_{y_n}$. For two observations indexed by y and y', define $d_{yy'} = \|\boldsymbol{\mu}_y - \boldsymbol{\mu}_{y'}\|_2^2$. We have

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

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

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

where $\mathbf{1}_n$ is an *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 **V** are centered; that is, $\sum_{i=1}^{n} \boldsymbol{v}_{y_i}$ is the *d*-vector of zeros.

Because $d_{yy'}$ is actually unknown, we replace it with $\hat{d}_{yy'} = \|\boldsymbol{x}_y - \boldsymbol{x}_{y'}\|_2^2 - 2p$. It is easy to show that $\hat{d}_{yy'}$ is an unbiased estimate of $d_{yy'}$. Let $\hat{\mathbf{D}} = (\hat{d}_{yy'})$ and $\mathbf{X} = (\boldsymbol{x}_{y_1}, \dots, \boldsymbol{x}_{y_n})^\top \in \mathbb{R}^{n \times p}$. 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 $\lambda_1 \geq \cdots \geq \lambda_n \geq 0$ are the eigenvalues, and $\alpha_1, \ldots, \alpha_n$ are the corresponding eigenvectors. By the Eckart–Young theorem, a solution for **V** is given by

$$ilde{\mathbf{V}} = (\lambda_1^{1/2} oldsymbol{lpha}_1, \dots, \lambda_d^{1/2} oldsymbol{lpha}_d)$$

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

We can view \tilde{v}_y as the vector of coordinates of x_y with respect to the orthonormal basis Γ . From the viewpoint of regression dimension reduction, $\tilde{\mathbf{V}}$ then contains all 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, $\boldsymbol{x}_{y^*}, y^* \in \mathbb{Y}$. This can be done using the classical method of adding a point to vector diagrams (Gower, 1968; Zhang et al., 2008). For each $i \in$ $\{1, \ldots, n\}$, we define $\tilde{s}_i = \|\tilde{\boldsymbol{v}}_{y_i}\|_2^2 - \|\boldsymbol{x}_{y^*} - \boldsymbol{x}_{y_i}\|_2^2$. Let $\tilde{\boldsymbol{s}} = (\tilde{s}_1, \ldots, \tilde{s}_n)^\top \in \mathbb{R}^n$. Then, the predicted coordinates $\tilde{\boldsymbol{v}}_{y^*}$ of \boldsymbol{x}_{y^*} are given by

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

In a classical sufficient reduction, one is interested mainly in the matrix Γ , or the subspace $S_{Y|X}$ spanned by it. The above procedure operates in the space of coordinates of x_y with respect to the orthonormal basis Γ . The approach is appealing because it achieves dimension reduction while avoiding the need to estimate Γ .

2.2 Subspace estimation

Once \boldsymbol{v}_y has been determined, it becomes natural to use the least squares method to estimate $\boldsymbol{\Gamma}$ in model (2.2), if desired. Specifically, we estimate Γ 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 is given by

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

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

After some further manipulations, $\tilde{\Gamma}_j$ can be shown to equal the *j*th principal component direction of $\mathbf{P}_n \mathbf{X}$. Thus, the first *d* principal component score vectors of $\mathbf{P}_n \mathbf{X}$ produce 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 the former assumes $\boldsymbol{\epsilon}$ is normally distributed, and it employs a maximum likelihood estimation.

2.3 A toy example

Before we proceed, we consider a simple simulation with p = 5 and d = 2. Observations on (\boldsymbol{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, $\boldsymbol{X}_y = \boldsymbol{x}_y$ was generated according to $\boldsymbol{X}_y = \boldsymbol{\Gamma} \boldsymbol{v}_y + \boldsymbol{\epsilon}$, where $\boldsymbol{\Gamma} = (\boldsymbol{\Gamma}_1, \boldsymbol{\Gamma}_2)^{\top}$, with $\boldsymbol{\Gamma}_1 = (1, 0, 0, 0, 0)^{\top}$ and

2.3 A toy example

 $\Gamma_2 = (0, 1, 0, 0, 0)^{\top}$, and $\boldsymbol{v}_y = (y, y^2/3)^{\top}$. The error vector was sampled from a normal distribution with mean vector **0** and covariance matrix $\boldsymbol{\Delta} =$ diag(1, 1, 5, 5, 5).

In the upper plot of Figure 1, the two-dimensional coordinates of 200 CMDS samples are displayed, with each sample indexed 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 left to right, and that the middle and extreme response 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 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 model the coordinate vectors as

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

where $\boldsymbol{\beta} \in \mathbb{R}^{d \times r}$ has rank $d \leq \min(r, p)$, and $\boldsymbol{f}_y \in \mathbb{R}^r$ is a known vectorvalued function of y. Usually, \boldsymbol{f}_y is required to contain a reasonably flexible set of basis functions, such as slice indicator functions or B-spline basis functions. This parameterization is widely used in model-based reduction; see, for example, Cook and Forzani (2008), Cook et al. (2012), and Wang and Zhu (2013). Replacing \boldsymbol{v}_y in model (2.2) with $\boldsymbol{\beta} \boldsymbol{f}_y$, we have the following model:

$$\boldsymbol{X}_{y} = \boldsymbol{\mu} + \boldsymbol{\Gamma} \boldsymbol{\beta} \boldsymbol{f}_{y} + \boldsymbol{\epsilon}. \tag{3.5}$$

We refer to this model as a supervised inverse regression model. Without loss of generality, we assume that the sample mean vector of f_y is zero.

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

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

Let $\boldsymbol{\pi} = (\boldsymbol{f}_{y_{1}}^{\top}\boldsymbol{\beta}^{\top}\boldsymbol{\beta}\boldsymbol{f}_{y_{1}}, \dots, \boldsymbol{f}_{y_{n}}^{\top}\boldsymbol{\beta}^{\top}\boldsymbol{\beta}\boldsymbol{f}_{y_{n}})^{\top} \in \mathbb{R}^{n}$ and $\mathbf{F} = (\boldsymbol{f}_{y_{1}}, \dots, \boldsymbol{f}_{y_{n}})^{\top} \in$



3 A SUPERVISED INVERSE REGRESSION MODEL

Figure 1: Two-dimensional plots of 200 samples from the inverse regression model in (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

 $\mathbb{R}^{n \times r}$. In matrix form,

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

Because $\mathbf{P}_n \mathbf{F} = \mathbf{F}$, 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 **D** with $\hat{\mathbf{D}}$, we have

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

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

$$(\mathbf{F}^{ op}\mathbf{F})^{-1/2}\mathbf{F}^{ op}\mathbf{X}\mathbf{X}^{ op}\mathbf{F}(\mathbf{F}^{ op}\mathbf{F})^{-1/2} = \sum_{j=1}^r
ho_j \phi_j \phi_j^{ op}$$

where $\rho_1 \geq \cdots \geq \rho_r \geq 0$ are the eigenvalues, and ϕ_1, \ldots, ϕ_r are the corresponding eigenvectors. A solution for β is then given by

$$ilde{oldsymbol{eta}} = (
ho_1^{1/2} oldsymbol{\phi}_1, \dots,
ho_d^{1/2} oldsymbol{\phi}_d)^ op (\mathbf{F}^ op \mathbf{F})^{-1/2}.$$

Furthermore,

$$\tilde{\boldsymbol{v}}_y = \boldsymbol{eta} \boldsymbol{f}_y$$

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

4. Reduction when Δ is unknown

4.1 The proposed method

In practice, Δ is seldom known in advance, and thus has to be estimated from the data. Throughout this paper, we estimate Δ using the residual sample covariance matrix from the multivariate linear regression of X_y on f_y :

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

where $\mathbf{P}_{\mathbf{F}} = \mathbf{F}(\mathbf{F}^{\top}\mathbf{F})^{-1}\mathbf{F}^{\top}$ is the hat matrix. The asymptotic properties of $\hat{\Delta}$ can be found in Lemmas 1 and 2 in the Supplementary Material. 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.

We fix Δ at $\hat{\Delta}$, and base the analysis on the standardized data $\mathbf{X}\hat{\Delta}^{-1/2}$. For simplicity, we focus on the supervised inverse regression model (3.5). Replacing \mathbf{X} with $\mathbf{X}\hat{\Delta}^{-1/2}$, we compute

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

and its eigen-decomposition

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

We estimate $\boldsymbol{\beta}$ and \boldsymbol{v}_y as

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

and

$$\hat{m{v}}_y = \hat{m{eta}} m{f}_y$$

Define $\hat{s}_i = \|\hat{\boldsymbol{v}}_{y_i}\|_2^2 - \|\hat{\boldsymbol{\Delta}}^{-1/2}(\boldsymbol{x}_{y^*} - \boldsymbol{x}_{y_i})\|_2^2$. Let $\hat{\boldsymbol{s}} = (\hat{s}_1, \dots, \hat{s}_n)^\top$ and $\hat{\mathbf{V}} = (\hat{\boldsymbol{v}}_{y_1}, \dots, \hat{\boldsymbol{v}}_{y_n})^\top$. Then, the vector of coordinates of a new observation \boldsymbol{x}_{y^*} is predicted by

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

We refer to this method as supervised reduction via inverse regression (SRIR).

As mentioned earlier, the advantage of working with coordinate vectors is that a reduction can be performed without needing to estimate Γ or $S_{Y|X}$. Nevertheless, there are situations in which the inferential target is $S_{Y|X}$, as is the case in a traditional sufficient dimension reduction. To conduct a reduction in the original predictor space, we have to determine both Δ and Γ . In general, it is infeasible to find a closed-form expression for these estimators, and so we usually need an alternating procedure. Fortunately, the estimate $\hat{\Delta}$ has nothing to do with Γ , suggesting a one-step estimate for Γ . Specifically, we estimate Γ by minimizing the residual sum-of-squares

$$RRS(\mathbf{\Gamma}) = \|\mathbf{P}_n \mathbf{X} \hat{\boldsymbol{\Delta}}^{-1/2} - \hat{\mathbf{V}} \mathbf{\Gamma}^\top \hat{\boldsymbol{\Delta}}^{-1/2} \|_F^2$$
$$= \operatorname{trace}\{(\mathbf{P}_n \mathbf{X} - \hat{\mathbf{V}} \mathbf{\Gamma}^\top) \hat{\boldsymbol{\Delta}}^{-1} (\mathbf{P}_n \mathbf{X} - \hat{\mathbf{V}} \mathbf{\Gamma}^\top)^\top\}.$$

The solution is

$$\hat{\boldsymbol{\Gamma}} = \mathbf{X}^{\top} \hat{\mathbf{V}} (\hat{\mathbf{V}}^{\top} \hat{\mathbf{V}})^{-1} = \mathbf{X}^{\top} \mathbf{F} \hat{\boldsymbol{\beta}}^{\top} (\hat{\boldsymbol{\beta}} \mathbf{F}^{\top} \mathbf{F} \hat{\boldsymbol{\beta}}^{\top})^{-1}.$$
(4.

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

4.2 Theoretical properties

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

Theorem 1. Assume that $v_Y = \beta f_Y$ has finite sixth moments, and that ϵ has finite fourth moments. Then, for some $d \times d$ rotation matrix \mathbf{R} ,

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

For two *d*-dimensional random vectors V_1 and V_2 , let Σ_1, Σ_2 , and Σ_{12} denote the covariance matrix of V_1 , the covariance matrix of V_2 , and the covariance matrix between V_1 and V_2 , respectively. To assess the prediction accuracy, we use the multiple correlation coefficient, which is

defined as the positive square root of

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

This measure takes the maximum value of one when V_1 and V_2 are linearly related, and takes the minimum zero when 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,

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

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

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

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

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

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

$$\hat{oldsymbol{v}}_{y^*} = \mathbf{R}(oldsymbol{c} + \mathbf{A}oldsymbol{v}_{y^*} + \mathbf{A}oldsymbol{\Gamma}^ op \mathbf{\Omega}^{-1}oldsymbol{\epsilon}_{y^*}) + O_P\left(rac{1}{\sqrt{n}}
ight),$$

for some $d \times d$ rotation matrix **R**, where

$$\boldsymbol{c} = \frac{1}{2} \{ \boldsymbol{\Phi} \operatorname{Var}(\boldsymbol{f}_Y) \boldsymbol{\Phi}^\top \}^{-1} \{ \operatorname{E}(\boldsymbol{\Phi} \boldsymbol{f}_Y \boldsymbol{f}_Y^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{f}_Y) - \operatorname{E}(\boldsymbol{\Phi} \boldsymbol{f}_Y \boldsymbol{v}_Y^\top \boldsymbol{v}_Y) \},$$
$$\boldsymbol{A} = \{ \boldsymbol{\Phi} \operatorname{Var}(\boldsymbol{f}_Y) \boldsymbol{\Phi}^\top \}^{-1} \boldsymbol{\Phi} \operatorname{Cov}(\boldsymbol{f}_Y, \boldsymbol{v}_Y),$$

and

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

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 f_Y and v_Y are sufficiently correlated. From a dimension reduction point of view, we can treat v_{y^*} and $c + Av_{y^*}$ as the same reduction.

The following theorem gives the consistency of the subspace estimation.

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

4.3 Choice of d

Let $\hat{\Sigma} = \mathbf{X}^{\top} \mathbf{P}_{\mathbf{F}} \mathbf{X}$. Define $\mathcal{S}_d(\hat{\Delta}, \hat{\Sigma})$ as the span of $\hat{\Delta}^{-1/2}$ multiplied by the first d eigenvectors of $\hat{\Delta}^{-1/2} \hat{\Sigma} \hat{\Delta}^{-1/2}$. One connection between our onestep subspace estimate and the maximum likelihood estimate is captured in the following theorem.

Theorem 4. Assume that f_y is correctly specified. Then, $\operatorname{span}(\hat{\Delta}^{-1}\hat{\Gamma}) = S_d(\hat{\Delta}, \hat{\Sigma})$. Consequently, if ϵ is normally distributed, then $\operatorname{span}(\hat{\Delta}^{-1}\hat{\Gamma})$ is the maximum likelihood estimate of $S_{Y|X}$.

4.3 Choice of d

In practice, the structural dimension d is unknown, and thus 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 uses an information criterion (Zhu et al., 2006). Let d_0 denote the true dimension. To estimate d_0 , we propose using the Bayesian information criterion (Zhu et al., 2012). With

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

the estimated dimension is

$$\hat{d} = \arg \max_{1 \le d \le r} \text{BIC}_d. \tag{4.8}$$

5 NUMERICAL STUDIES

Theorem 5. Assume the conditions of Theorem 1 or Theorem 2 hold. 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 in an 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 $\boldsymbol{\Delta} = (\theta^{|i-j|})$, with θ taking 0 or 0.5. Let $\boldsymbol{\Gamma}_{01} = (1, 1, -1, -1, 0, \dots, 0)^{\top}/2$, $\boldsymbol{\Gamma}_{02} = (1, 0, 1, 0, 1, 0, \dots, 0)^{\top}/\sqrt{3}$, and $\boldsymbol{\Gamma}_{0} = (\boldsymbol{\Gamma}_{01}, \boldsymbol{\Gamma}_{02})$. Set $\boldsymbol{\Gamma} = \boldsymbol{\Gamma}_{0}(\boldsymbol{\Gamma}_{0}^{\top}\boldsymbol{\Delta}^{-1}\boldsymbol{\Gamma}_{0})^{-1/2}$.

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, $\{(\boldsymbol{x}_{y_i^*}, y_i^*)\}$, of size 100. Let $\hat{\boldsymbol{v}}_{y_i^*}$ be the predicted vector of coordinates of $\boldsymbol{x}_{y_i^*}$. To measure the closeness between $\hat{\boldsymbol{v}}_{y_i^*}$ and $\boldsymbol{v}_{y_i^*}$, we used the sample version of the multiple correlation coefficient (MCC). For each configuration, the number of repetitions was 200.

Example 1. To gain insight into the operating characteristics of the

5.1 Simulations

proposed method, consider the model

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

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

$$ho^2(\hat{\boldsymbol{v}}_{Y^*}, \boldsymbol{v}_{Y^*}) = g^2(\sigma) + O_P\left(rac{1}{\sqrt{n}}
ight).$$

Here, $g(\sigma) = \sqrt{\sigma^2/(2\sigma^2 + 2) + \sigma^4/(2\sigma^4 + 1)}$ is an increasing function of σ . Six values of σ 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 there is an excellent agreement between the theoretical prediction and the empirical behavior.

Next, we examine the behavior of our method in further detail. In addition to the prediction accuracy, we also assessed its 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, respectively. Denote by $\rho_1 \geq \cdots \geq \rho_d$ the ordered eigenvalues of $\hat{\mathbf{B}}^{\top} \mathbf{B} \mathbf{B}^{\top} \hat{\mathbf{B}}$. VCC is defined as the positive square root of $q^2(\hat{\mathbf{B}}, \mathbf{B}) = \prod_{i=1}^d \rho_i$, and TCC is defined as the positive square root of $r^2(\hat{\mathbf{B}}, \mathbf{B}) = d^{-1} \sum_{i=1}^d \rho_i$.

5.1 Simulations



Figure 2: The estimated MCC curves, based on 200 data replications, for $\theta = 0$ (dotted) and $\theta = 0.5$ (solid), along with the theoretical MCC curve $g(\sigma)$ (longdashed). The error bars indicate one standard deviation.

Example 2. Consider the model

$$\boldsymbol{X}_{y} = \boldsymbol{\Gamma} \boldsymbol{\beta} \boldsymbol{f}_{y} + \boldsymbol{\epsilon},$$

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 nonGaussian errors were explored, with covariance structures as in the previous example. In the first, $\boldsymbol{\epsilon}$ is drawn from a multivariate *t*-distribution with five degrees of freedom. In the second, each component of $\boldsymbol{\epsilon}$ is uniformly distributed on $[-\sqrt{3}, \sqrt{3}]$. For the coefficient matrix $\boldsymbol{\beta}$, we set

$$\boldsymbol{\beta} = \left(\begin{array}{ccc} 1 & 0 & 0 \\ & & \\ 0 & 0 & 1 \end{array} \right) \quad \text{or} \quad \left(\begin{array}{ccc} 1 & -0.5 & 0 \\ & & \\ 0 & 0.5 & 1 \end{array} \right).$$

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 the nonGaussian errors is comparable with that under the Gaussian assumption (Figure 2, $\sigma = 1$). Furthermore, our method performs well in terms of subspace estimation. In general, the performance improves as the number of predictors decreases. From Tables 2 and 4, we see that our method is

5.1 Simulations

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. ϵ is drawn from a multivariate t-

distribution with five degrees of freedom, and \boldsymbol{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(0.0)	(0.66) 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 (0.0	(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 (0.0	(0.279) 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 (0.0	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. ϵ is drawn from a multivariate tdistribution with five degrees of freedom, and 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)	

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

distributed on $[-\sqrt{3}, \sqrt{3}]$, and \boldsymbol{f}_y is misspecified.

			SRIR				
		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 the principal components (PC) and principal fitted components (PFC) methods of Cook and Forzani (2008). The PC results are shown in the last three columns of Tables 1–4. SRIR appears to dominate PC in most cases, especially when $\theta = 0.5$. The PFC results are the same as the SRIR results and, thus, are omitted. For

5.1 Simulations

subspace estimation, this is expected, from Theorem 4, but for prediction, this comes as somewhat of a surprise. We provide theoretical support for this conclusion in the Supplementary Material.

Thus far, we have assumed that the value of the structural dimension is 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 greater than 80% in each configuration. We also see that a misspecification can have a significant impact.

Table 5: Selection frequencies of BIC over 200 data replications. ϵ is drawn from a multivariate *t*-distribution with five degrees of freedom.

		Correctl	y specified	Misspe	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	

component of $\boldsymbol{\epsilon}$ 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	

Table 6: Selection frequencies of BIC over 200 data replications. Each

5.2 Boston housing data

We applied SRIR to Boston housing data (Harrison and Rubinfeld, 1978), 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 the per capita crime rate by town, average number of rooms per house, and percentage of households with low socioeconomic status, among others.

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 a two-dimensional plot of the 506 observations, with coordinates computed using (4.6). We see a horseshoe-like pattern in the data cloud.

6 DISCUSSION

We also see an association between the response and the coordinates, similar to that in the toy example. For comparison purposes, we also carried out CMDS. 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 shows 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 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. We have proposed a principled reduction method in a dual sample-based space, based on a supervised inverse regression model. Instead of estimating the subspace, our method directly estimates the projection coordinates of the predictor vector onto the subspace. The results extend the well-known duality between principal component analysis and CMDS.

6 DISCUSSION



Figure 3: Two-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.

6 DISCUSSION



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

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

We have presented the theoretical properties of our method, supported by simulation results. As with most reduction methods, we have adopted the traditional asymptotic reasoning of letting the sample size $n \to \infty$, with the number of predictors p fixed. Our method requires the inverse of the residual sample covariance matrix, and hence is problematic when pis 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), known as a permutation multivariate analysis of variance. This procedure 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 trace($\mathbf{X}^{\top}\mathbf{X}$) = 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 $d_{ij}(\mathbf{B}) = (\boldsymbol{x}_i - \boldsymbol{x}_j)^\top \mathbf{B}(\boldsymbol{x}_i - \boldsymbol{x}_j)$, then $\mathbf{XBX}^\top = -\mathbf{P}_n \mathbf{DP}_n/2$, where \mathbf{P}_n is the centering matrix. Similarly to the permutation multivariate analysis of variance, we can extend our supervised reduction method, based solely on measures of distance or dissimilarity between pairs of observations, 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. These topics are left to future research.

Supplementary Material

The online Supplementary Material contains proofs of the relevant lemmas and theorems.

Acknowledgements

Tao Wang is the corresponding author. Peirong Xu was supported by the National Natural Science Foundation of China (11971018) and Natural Science Foundation of Shanghai (19ZR1437000). Tao Wang was supported, in part, by the National Natural Science Foundation of China (11601326, 11971017), National Key R&D Program of China (2018YFC0910500), Shanghai

REFERENCES

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.

- 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.

- 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 distancebased 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

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.

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

