

Statistica Sinica Preprint No: SS-2016-0080R3

Title	Flexible dimension reduction in regression
Manuscript ID	SS-2016-0080R3
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202016.0080
Complete List of Authors	Lixing Zhu and Tao Wang
Corresponding Author	Lixing Zhu
E-mail	lzhu@hkbu.edu.hk
Notice: Accepted version subject to English editing.	

Flexible dimension reduction in regression*

Tao Wang^{1,2} and Lixing Zhu³

¹ Department of Bioinformatics and Biostatistics, Shanghai Jiao Tong University, Shanghai, China

² SJTU-Yale Joint Center for Biostatistics, Shanghai Jiao Tong University, Shanghai, China

³ Department of Mathematics, Hong Kong Baptist University, Hong Kong

Abstract

Sliced inverse regression is a valuable tool for dimension reduction. One can replace the predictor vector with a few linear combinations of its components without loss of information on the regression. This paper is about richer nonlinear dimension reduction. Each direction of sliced inverse regression is simply a slope vector of multiple linear regression applied to an optimally transformed response. Using this connection, we propose a nonlinear version of sliced inverse regression by replacing linear function by an additive function of the predictors. Our procedure has a clear interpretation as sliced inverse regression on a set of adaptively chosen transformations of the predictors. The flexibility of our method is illustrated via a simulation study and a data application.

Key words: Canonical correlation; Optimal scoring; Sufficient dimension reduction.

*The research described herewith was supported by National Natural Science Foundation of China (11601326), and the University Grants Council of Hong Kong, Hong Kong. The authors thank the Editor, the Associate Editor, and two anonymous referees for their constructive comments and suggestions that led to an improvement of an early manuscript.

1 Introduction

Transformation is a commonly used technique in statistics. Regarding regression, Box & Cox (1964) stated that “*In regression problems, where both dependent and independent variables can be transformed, there are more possibilities to be considered.*” On one hand, transformation of the response variable can often achieve simultaneously normally distributed errors with a constant variance and a linear regression function. The methodology pioneered by Box & Cox (1964) is perhaps the most common. On the other hand, transformation of the predictors can convert a complex nonlinear regression relation into a simple and often linear one. For example, if it appears in a problem that many predictors are likely to need transformation, then iterative fitting of additive models should be considered. More generally, nonlinear multivariate techniques, such as alternating conditional expectations of Breiman & Friedman (1985) and monotone spline regression of Ramsay (1988), allow transformation on both the response variable and the predictors. However, only one transformation on the response variable is allowed in these algorithms for the purpose of model fitting.

A fairly common practice is to apply first either marginal transformations or a joint transformation to multivariate data to induce normality, elliptical symmetry, or other appropriate distributions, and then carry out model-based regression, in most cases linear regression, on the transformed data. Since this strategy is statistically unjustified, it often leads to model mis-specification giving biased estimators. In principle, any use of transformations requires that the effects of them on the error structure be understood. The story is very different, however, in the context of dimension reduction.

Regarding dimension reduction in regression, transformation of variables plays two roles. On one hand, response transformations, such as slicing (Duan & Li 1991, Li 1991) and spline transformation (He & Shen 1997, Fung et al. 2002), serve as an intermediate tool for finding interesting patterns in the data, instead of being used to improve the goodness of model

fitting in a traditional way. See also Zhu et al. (2010), Wu & Li (2011), Yin & Li (2011), and references therein. On the other hand, predictor transformations are useful for reducing the structural dimension of the regression when the data are concentrated on a nonlinear low-dimensional space (Cook 1998, Chapter 14). In many cases, it is preferable to use predictor transformations rather than response transformations, because transformation of the predictors does not change the interpretation of the response variable. This is particularly the case when we are not assuming a model for the regression.

Related to the present work is Wang et al. (2014), who present a framework for dimension reduction in regression that lies between linear and fully nonlinear dimension reduction. Suppose Y is a univariate response variable and $\mathbf{X} = (X_1, \dots, X_p)^\top$ is a p -vector of predictors. The main idea is to transform first each of the raw predictors marginally and monotonically, in the form $\mathbf{f}(\mathbf{X}) = \{f_1(X_1), \dots, f_p(X_p)\}^\top$, and then search for a low-dimensional projection in the space defined by the transformed predictors. Toward this end, they assume that, given $\mathbf{B}^\top \mathbf{f}(\mathbf{X})$, Y is independent of $\mathbf{f}(\mathbf{X})$, where \mathbf{B} is a $p \times d$ matrix with $d \leq p$. The aim of the analysis is then to characterize the subspace spanned by the columns of \mathbf{B} , and for this purpose they propose a two-step procedure by combining probability integral transformation and sliced inverse regression (Li 1991). In particular, assuming that the distribution of $\mathbf{f}(\mathbf{X})$ is multivariate Gaussian, they proceed as follows: replacing the observations, for each predictor, by their corresponding normal scores, and then applying sliced inverse regression to the transformed data. Although $\mathbf{f}(\mathbf{X})$ or its distribution is user-specified, probability integral transformed sliced inverse regression suffers from mis-specification of transformations.

In this paper we propose a new nonlinear dimension-reduction method to overcome this problem. This is primarily motivated by observing the close connection between sliced inverse regression and multiple linear regression. Our procedure estimates predictor transformations in a data-driven way, and thus can be regarded as an adaptive version of probability integral transformed sliced inverse regression.

2 Methodology

2.1 Sliced inverse regression by optimal scoring

The basic idea of sliced inverse regression by optimal scoring (Wang & Zhu 2013) is to linearise a response transformation $T(Y)$ by $\boldsymbol{\phi}(Y)^\top \boldsymbol{\theta}$, where $\boldsymbol{\phi} = (\phi_1, \dots, \phi_K)^\top$ is a K -vector of basis functions, and $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^\top$ is a K -vector of unknown scores. The scored data is then predicted by linear regression on \mathbf{X} . The simultaneous estimation of the scores and the linear parameters constitutes the optimal scoring problem.

Suppose that $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ is a random sample on (\mathbf{X}, Y) . Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$ and $\boldsymbol{\Phi} = \{\boldsymbol{\phi}(y_1), \dots, \boldsymbol{\phi}(y_n)\}^\top$ be the two data matrices containing the predictor values and the basis function values respectively. Write $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_p)$ with $\mathbf{X}_k = (x_{1k}, \dots, x_{nk})^\top$ as its k -th column. Without loss of generality, we assume that the columns of \mathbf{X} are centered. In the sample, the criterion of sliced inverse regression by optimal scoring takes the form

$$\begin{aligned} & \underset{\boldsymbol{\theta}_i \in \mathbf{R}^K, \boldsymbol{\beta}_i \in \mathbf{R}^p}{\text{minimize}} && \|\boldsymbol{\Phi} \boldsymbol{\theta}_i - \mathbf{X} \boldsymbol{\beta}_i\|_2^2 \\ & \text{subject to} && \boldsymbol{\theta}_i^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\theta}_i = n, \boldsymbol{\theta}_i^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\theta}_j = 0, \quad j = 1, \dots, i-1, \end{aligned} \quad (2.1)$$

where $i = 1, \dots, d \leq K$.

There are various choices for the basis functions. We in this paper concentrate on slice indicator functions. Since the columns of \mathbf{X} are centered to have mean zero, one can see that the constant score vector $\mathbf{1}$ of length K is trivial, and hence there are at most $K - 1$ nontrivial solutions to (2.1).

REMARK 2.1. *In the population, sliced inverse regression by optimal scoring sequentially solves*

$$\begin{aligned} & \underset{T_i, a_i \in \mathbf{R}, \mathbf{b}_i \in \mathbf{R}^p}{\text{minimize}} && E(T_i - a_i - \mathbf{X}^\top \mathbf{b}_i)^2 \\ & \text{subject to} && \text{var}\{T_i(Y)\} = 1, \text{cov}\{T_i(Y), T_j(Y)\} = 0, \quad j = 1, \dots, i-1. \end{aligned} \quad (2.2)$$

The i -th optimal transformation $T_i(Y)$ is identical up to a scalar multiplication to $E(\boldsymbol{\eta}_i^\top \mathbf{X}|Y)$ and $\mathbf{b}_i(T_i)$ is proportional to $\boldsymbol{\eta}_i$, where $\boldsymbol{\eta}_i$ is the i -th sliced inverse regression direction such that

$$\text{cov}\{E(\mathbf{X}|Y)\}\boldsymbol{\eta}_i = \lambda_i \text{cov}(\mathbf{X})\boldsymbol{\eta}_i.$$

Here $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_p$ are eigenvectors satisfying $\boldsymbol{\eta}_i^\top \text{cov}(\mathbf{X})\boldsymbol{\eta}_j = 1$ if $i = j$, and 0 if $i \neq j$, and $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ are the corresponding eigenvalues. Details can be found in Chen & Li (1998), Wang & Zhu (2013).

To study the relationship between a set of predictors and a categorical response, it is known that Fisher's linear discriminant analysis, canonical correlation analysis, and optimal scoring are equivalent from the dimension-reduction point of view (Hastie et al. 1995). These data-analytic tools, which are popular for classification, are also useful in the regression setting. First, the original sliced inverse regression method is formally equivalent to Fisher's linear discriminant analysis (Kent 1991), and it is tantamount to sliced inverse regression by optimal scoring by using the slice indicator functions obtained from slicing the response variable. Second, Fung et al. (2002) proposed a dimension-reduction method based on canonical correlation, which can be viewed as a variant of sliced inverse regression. The latter estimates the kernel matrix $\text{cov}\{E(\mathbf{X}|Y)\}$ using the slice indicator functions, while the former produces a spline-based estimate by using the B-spline basis functions generated for the response variable. Clearly, optimal scoring and canonical correlation analysis are equivalent without regard to classification or regression. Therefore, sliced inverse regression by optimal scoring includes the original sliced inverse regression method and the canonical correlation method as special cases.

In typical applications, the number of basis functions K needed is very small (He & Shen 1997, Fung et al. 2002), and it is well known that the original sliced inverse regression method is not overly sensitive to the number of slices.

REMARK 2.2. *Due to the equivalence between optimal scoring and canonical correlation anal-*

ysis, sliced inverse regression by optimal scoring should have the same theoretical properties as those of the canonical correlation method. Further, sliced inverse regression by optimal scoring can be interpreted as a method for estimating the canonical correlations between the columns of \mathbf{X} and the columns of Φ , and thus the directions found by it are useful in their own right in identifying some important features of the regression of Y on \mathbf{X} . In particular, if the linearity condition of Li (1991) holds, the directions $\boldsymbol{\eta}_i$ in Remark 2.1 belong to the central dimension-reduction subspace, which is an essential concept of sufficient dimension reduction (Cook 1998).

2.2 Flexible dimension reduction

We first introduce some notation and definitions. We call $\text{span}(\mathbf{B})$ a transformed dimension-reduction subspace with respect to \mathbf{f} , if

$$Y \perp\!\!\!\perp \mathbf{X} \mid \mathbf{B}^\top \mathbf{f}(\mathbf{X}). \quad (2.3)$$

$\text{span}(\mathbf{B})$ is called a transformed central dimension-reduction subspace, if it satisfies (2.3) for some \mathbf{f} , and at the same time it has the smallest dimension. The following proposition shows that the transformed central dimension-reduction subspace is well defined.

PROPOSITION 2.1. *Assume that there is a p -vector $\mathbf{x}^* = (x_1^*, \dots, x_p^*)^\top$, such that for all $j = 1, \dots, p$, $f_j'(x_j^*) = 1$ and $f_j'(x_j)$ is continuous at x_j^* , where $f_j'(x_j)$ denotes the derivative of $f_j(x_j)$ with respect to x_j . Then the transformed central dimension-reduction subspace is identifiable.*

REMARK 2.3. *While Wang et al. (2014) requires the transformations f_j to be pre-specified, we propose to estimate f_j and perform sufficient dimension reduction simultaneously. One consequence of this difference is that f_j may not be identifiable, even if the structural dimension is given. For example, if $Y = X_1^3 + \epsilon$, then $f_1(x)$ can be the identity map x or x^3 , or any other monotone function. Note that sufficient dimension reduction is often considered*

as the first step in statistical analyses. After that, graphical tools or nonparametric methods can be used to further investigate the relationship between the response variable and the reduced set of variables. Suppose $\mathbf{g} = \{g_1, \dots, g_p\}$ is an alternative set of transformations such that $\text{span}(\mathbf{B})$ is a transformed central subspace with respect to \mathbf{g} . Theoretically, $\mathbf{B}^\top \mathbf{f}(\mathbf{X})$ and $\mathbf{B}^\top \mathbf{g}(\mathbf{X})$ can be treated equally, because each of them is a minimal and sufficient reduced predictor. But, numerically, it is possible that the identifiability issue of transformation functions affects speed and convergence as seen in the next section.

Sliced inverse regression by optimal scoring is a linear dimension-reduction method. However, our interest in this approach is that it includes linear regression as a building block. Many techniques exist for generalizing linear regression to more flexible, nonlinear and/or nonparametric forms of regression. This in turn leads to more flexible forms of dimension reduction. One simple and effective approach is to augment the set of predictors to include quadratic and bilinear terms, and then carry out sliced inverse regression by optimal scoring in the enlarged space, which in effect results in nonlinear dimension reduction in the original predictor space. More flexible approaches use adaptive nonparametric regression or kernels (Hastie et al. 1994, 2009).

The procedure that we have in mind is less “nonlinear and/or nonparametric” than these, but it has a clear interpretation as a linear dimension-reduction method on a set of marginally and adaptively transformed predictors. Specifically, we transform the raw predictors $\mathbf{X} = (X_1, \dots, X_p)^\top$ in the form $\mathbf{f}(\mathbf{X}) = \{f_1(X_1), \dots, f_p(X_p)\}^\top$, or simply $\mathbf{f} = (f_1, \dots, f_p)^\top$, for a set of smooth univariate functions $\{f_1, \dots, f_p\}$, and at the same time consider the optimal scoring problem with the transformed data as follows:

$$\begin{aligned} & \underset{\mathbf{f}, \{\boldsymbol{\theta}_i\}, \{\boldsymbol{\beta}_i\}}{\text{minimize}} && \sum_{i=1}^H \|\boldsymbol{\Phi} \boldsymbol{\theta}_i - \mathbf{X}_f \boldsymbol{\beta}_i\|_2^2 \\ & \text{subject to} && \boldsymbol{\theta}_i^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\theta}_i = n, \boldsymbol{\theta}_i^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\theta}_j = 0, \quad j = 1, \dots, i-1, i = 1, \dots, H, \end{aligned} \tag{2.4}$$

where $1 \leq H \leq K$ is a working dimension, $\mathbf{X}_f = (\mathbf{X}_{f_1}, \dots, \mathbf{X}_{f_p})$ with $\mathbf{X}_{f_k} = \{f_k(x_{1k}), \dots, f_k(x_{nk})\}^\top$

as its k -th column. Since the f_k 's in (2.4) are up to scale and shift, we demand that $\sum_{i=1}^n f_k(x_{ik}) = 0$ and $\sum_{i=1}^n \{f_k(x_{ik})\}^2 = n$.

The restriction to component-wise monotone transformations is not always necessary. Our procedure thus permits as much flexibility in the estimated transformations as the data require, and can be interpreted as a method for estimating the canonical correlations between the transformed predictors and the basis functions generated for the response variable. This makes it very different from probability integral transformed sliced inverse regression in Wang et al. (2014).

REMARK 2.4. *One advantage of the proposed method over the original sliced inverse regression is that it can capture the marginal symmetry between \mathbf{X} and Y , such as in the model $Y = X_1^2 + X_2^2 + \epsilon$, where ϵ is the error independent of \mathbf{X} , and \mathbf{X} follows a symmetric distribution. As will be demonstrated in simulation studies, this is related to the robustness of sliced inverse regression against violation of the linearity condition, given that \mathbf{f} is correctly specified. However, our method also inherits one drawback of sliced inverse regression: it fails when the functional relation is joint symmetric in the sense that $E(\mathbf{f} | Y)$ is zero, such that in the model $Y = (X_1^3 + X_2^3)^2 + \epsilon$.*

The proposed procedure allows multiple transformations on the response variable for the purpose of dimension reduction, in contrast to alternating conditional expectations of Breiman & Friedman (1985) and monotone spline regression of Ramsay (1988) which focus on model fitting. It is closely related to monotone spline canonical correlation of Ramsay (1988), which is a tool for comparison between two sets of marginally and monotonically transformed variables via canonical correlation. The difference is that our procedure is concerned with dimension reduction in regression based on both response transformations and predictor transformations and in particular, the response transformations, or more precisely the set of basis functions, and the predictor transformations, are not restricted to be monotone.

2.3 Algorithm

To estimate \mathbf{f} , $\{\boldsymbol{\theta}_i\}$, and $\{\boldsymbol{\beta}_i\}$ in (2.4), we use an iterative approach: we first fix $\{\boldsymbol{\theta}_i\}$ and $\{\boldsymbol{\beta}_i\}$ and estimate \mathbf{f} , then we fix \mathbf{f} and estimate $\{\boldsymbol{\theta}_i\}$ and $\{\boldsymbol{\beta}_i\}$, and we iterate between these two steps until the algorithm converges.

When $\{\boldsymbol{\theta}_i\}$ and $\{\boldsymbol{\beta}_i\}$ are fixed, minimizing (2.4) with respect to $\mathbf{f} = (f_1, \dots, f_p)^\top$ is similar to fitting an additive model (Wood 2006). To this end, we use a variant of the back-fitting procedure. Write $\boldsymbol{\beta}_i = (\beta_{i1}, \dots, \beta_{ip})^\top$. For each $k = 1, \dots, p$, consider the one-dimensional smoothing problem

$$\underset{f_k}{\text{minimize}} \quad \sum_{i=1}^H \left\| \boldsymbol{\Phi} \boldsymbol{\theta}_i - \sum_{l \neq k} \beta_{il} \mathbf{X}_{f_l} - \beta_{ik} \mathbf{X}_{f_k} \right\|_2^2. \quad (2.5)$$

Let $\tilde{\mathbf{x}}_k = (\tilde{\mathbf{x}}_{1k}^\top, \dots, \tilde{\mathbf{x}}_{Hk}^\top)^\top$, $\mathbf{w}_k = (\mathbf{w}_{1k}^\top, \dots, \mathbf{w}_{Hk}^\top)^\top$, and $\tilde{\mathbf{y}}_k = (\tilde{\mathbf{y}}_{1k}^\top, \dots, \tilde{\mathbf{y}}_{Hk}^\top)^\top$, where

$$\tilde{\mathbf{x}}_{ik} = (x_{1k}, \dots, x_{nk})^\top, \mathbf{w}_{ik} = (\beta_{ik}^2, \dots, \beta_{ik}^2)^\top \in \mathbb{R}^n,$$

and

$$\tilde{\mathbf{y}}_{ik} = \frac{\boldsymbol{\Phi} \boldsymbol{\theta}_i - \sum_{l \neq k} \beta_{il} \mathbf{X}_{f_l}}{\beta_{ik}}.$$

Let $\tilde{n} = Hn$ and write $\tilde{\mathbf{x}}_k = (\tilde{x}_{1k}, \dots, \tilde{x}_{\tilde{n}k})^\top$, $\mathbf{w}_k = (w_{1k}, \dots, w_{\tilde{n}k})^\top$, and $\tilde{\mathbf{y}}_k = (\tilde{y}_{1k}, \dots, \tilde{y}_{\tilde{n}k})^\top$.

One can show that (2.5) is equivalent to

$$\underset{f_k}{\text{minimize}} \quad \sum_{s=1}^{\tilde{n}} w_{sk} \{\tilde{y}_{sk} - f_k(\tilde{x}_{sk})\}^2. \quad (2.6)$$

As will be seen below, the criterion adopted here is slightly different from that used in the back-fitting procedure. Specifically, in each (inner) iteration we identify and update the function that most reduces the objective function, but that function is kept fixed and will not be further updated in the next iteration. We stop the algorithm until convergence or until all the functions are updated. In our implementation, each f_k is represented using penalized regression splines with smoothing parameters selected by restricted maximum likelihood. Specifically, we use the function **gam** in the R package **mgcv**.

When \mathbf{f} is fixed, (2.4) becomes a standard optimal scoring problem

$$\begin{aligned} & \underset{\boldsymbol{\theta}_i \in \mathbf{R}^K, \boldsymbol{\beta}_i \in \mathbf{R}^p}{\text{minimize}} && \|\Phi \boldsymbol{\theta}_i - \mathbf{X}_f \boldsymbol{\beta}_i\|_2^2 \\ & \text{subject to} && \boldsymbol{\theta}_i^\top \mathbf{D} \boldsymbol{\theta}_i = 1, \boldsymbol{\theta}_i^\top \mathbf{D} \boldsymbol{\theta}_j = 0, \quad j = 1, \dots, i-1, \end{aligned} \quad (2.7)$$

where $\mathbf{D} = n^{-1} \Phi^\top \Phi$ and $i = 1, \dots, H$. The standard way of solving an optimal scoring problem is by way of a suitable eigenvalue decomposition. However, we propose to update $\{\boldsymbol{\theta}_i\}$ and $\{\boldsymbol{\beta}_i\}$ separately as follows. For fixed $\{\boldsymbol{\beta}_i\}$, the optimal scores $\{\boldsymbol{\theta}_i\}$ sequentially solve

$$\begin{aligned} & \underset{\boldsymbol{\theta}_i \in \mathbf{R}^K}{\text{minimize}} && \|\Phi \boldsymbol{\theta}_i - \mathbf{X}_f \boldsymbol{\beta}_i\|_2^2 \\ & \text{subject to} && \boldsymbol{\theta}_i^\top \mathbf{D} \boldsymbol{\theta}_i = 1, \boldsymbol{\theta}_i^\top \mathbf{D} \boldsymbol{\theta}_j = 0, \quad j = 1, \dots, i-1. \end{aligned} \quad (2.8)$$

Let $\mathbf{Q}_i = (\mathbf{1}, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{i-1})$ denote the $K \times i$ matrix consisting of the previous $i-1$ solutions including the trivial score vector of all ones. One can show that the i -th solution is given by

$$\boldsymbol{\theta}_i = c_i (\mathbf{I}_K - \mathbf{Q}_i \mathbf{Q}_i^\top \mathbf{D}) \mathbf{D}^{-1} \Phi^\top \mathbf{X}_f \boldsymbol{\beta}_i,$$

where \mathbf{I}_K is the $K \times K$ identity matrix and c_i is a constant such that $\boldsymbol{\theta}_i^\top \mathbf{D} \boldsymbol{\theta}_i = 1$. For fixed $\{\boldsymbol{\theta}_i\}$, we obtain H linear least squares problems

$$\underset{\boldsymbol{\beta}_i \in \mathbf{R}^p}{\text{minimize}} \quad \|\Phi \boldsymbol{\theta}_i - \mathbf{X}_f \boldsymbol{\beta}_i\|_2^2. \quad (2.9)$$

The solutions are easily seen to be $\boldsymbol{\beta}_i = (\mathbf{X}_f^\top \mathbf{X}_f)^{-1} \mathbf{X}_f^\top \Phi \boldsymbol{\theta}_i$.

In summary, the proposed algorithm for solving (2.4) proceeds as follows.

1. Standardization and initialization: Center and normalize \mathbf{X} . Let $\mathbf{X}_f = \mathbf{X}$. Initialize $\boldsymbol{\theta}_i$ and $\boldsymbol{\beta}_i$ with some plausible values. For example, we can use the solutions to (2.1). Let $\mathbf{w}_{ik} = (\beta_{ik}^2, \dots, \beta_{ik}^2)^\top \in \mathbf{R}^n$ and $\mathbf{w}_k = (\mathbf{w}_{1k}^\top, \dots, \mathbf{w}_{Hk}^\top)^\top$. Write $\mathbf{w}_k = (w_{1k}, \dots, w_{\tilde{n}k})^\top$.
2. Iterate until convergence in terms of $\{\boldsymbol{\beta}_i\}$ or until a maximum number of iterations is reached.
 - 2.1. Update $\mathbf{f} = (f_1, \dots, f_p)^\top$: Set $\mathcal{C} = \{1, \dots, p\}$ and $\check{\mathbf{y}}_i = \Phi \boldsymbol{\theta}_i$.

2.1.1. For each $k \in \mathcal{C}$, let

$$\tilde{\mathbf{y}}_{ik} = \frac{\check{\mathbf{y}}_i - \sum_{l \in \mathcal{C}, l \neq k} \beta_{il} \mathbf{X}_{f_l}}{\beta_{ik}}$$

and

$$\tilde{\mathbf{y}}_k = (\tilde{\mathbf{y}}_{1k}^\top, \dots, \tilde{\mathbf{y}}_{Hk}^\top)^\top = (\tilde{y}_{1k}, \dots, \tilde{y}_{\tilde{n}k})^\top.$$

Solve

$$\underset{f_k}{\text{minimize}} \quad \sum_{s=1}^{\tilde{n}} w_{sk} \{\tilde{y}_{sk} - f_k(\tilde{x}_{sk})\}^2.$$

2.1.2. Choose the function $f_{k'}$, $k' \in \mathcal{C}$, that most reduces the objective function.

Update

$$f_{k'} \leftarrow f_{k'}, \quad \check{\mathbf{y}}_i \leftarrow \check{\mathbf{y}}_i - \beta_{ik'} \mathbf{X}_{f_{k'}} \quad \text{and} \quad \mathcal{C} \leftarrow \mathcal{C} \setminus \{k'\}.$$

2.1.3. Continue in this way until the change in the objective function falls below a prespecified threshold or until all p functions have been updated. Center \mathbf{X}_f .

2.2. Update $\{\boldsymbol{\theta}_i\}$: Let $\mathbf{Q}_i = (\mathbf{1}, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{i-1})$ and

$$\boldsymbol{\theta}_i = (\mathbf{I}_K - \mathbf{Q}_i \mathbf{Q}_i^\top \mathbf{D}) \mathbf{D}^{-1} \boldsymbol{\Phi}^\top \mathbf{X}_f \boldsymbol{\beta}_i.$$

Normalize $\boldsymbol{\theta}_i$ so that $\boldsymbol{\theta}_i^\top \mathbf{D} \boldsymbol{\theta}_i = 1$. Normalize \mathbf{X}_f .

2.3. Update $\{\boldsymbol{\beta}_i\}$: Let

$$\boldsymbol{\beta}_i = (\mathbf{X}_f^\top \mathbf{X}_f)^{-1} \mathbf{X}_f^\top \boldsymbol{\Phi} \boldsymbol{\theta}_i.$$

2.4. Update $\{\mathbf{w}_k\}$: Let $\mathbf{w}_{ik} = (\beta_{ik}^2, \dots, \beta_{ik}^2)^\top \in \mathbb{R}^n$ and $\mathbf{w}_k = (\mathbf{w}_{1k}^\top, \dots, \mathbf{w}_{Hk}^\top)^\top$. Write $\mathbf{w}_k = (w_{1k}, \dots, w_{\tilde{n}k})^\top$.

Since the value of the objective function decreases at each step, the algorithm converges to a unique solution. However, it has no guarantee that the solution minimizes the objective function, because the overall problem is not convex. In general, the algorithm converges to a local minimum. This is the price paid for the computational efficiency, in view of the difficulty of carrying out the simultaneous minimization of (2.4). Our limited experience in simulation studies shows that the algorithm converges fast, and it works well empirically.

2.4 Theoretical properties

In order to simplify the exposition, we assume that for each j , there is a known finite-dimensional envelope $\{g_{j1}, \dots, g_{jq_j}\}$ such that f_j is a linear combination of g_{j1}, \dots, g_{jq_j} . Let $\mathbf{g}_j(X_j) = \{g_{j1}(X_j), \dots, g_{jq_j}(X_j)\}^\top$. Under this assumption, $f_j = \mathbf{c}_j^\top \mathbf{g}_j$ for some $\mathbf{c}_j \in \mathbb{R}^{q_j}$. We can write $\mathbf{f} = \mathbf{C}\mathbf{g}$, where $\mathbf{C} = \text{diag}(\mathbf{c}_1^\top, \dots, \mathbf{c}_p^\top)$ and $\mathbf{g} = (\mathbf{g}_1^\top, \dots, \mathbf{g}_p^\top)^\top$. Let $\mathbf{c}_{0j} \in \mathbb{R}^{q_j}$, $\mathbf{C}_0 = \text{diag}(\mathbf{c}_{01}^\top, \dots, \mathbf{c}_{0p}^\top)$, and $\mathbf{f}_0 = \mathbf{C}_0\mathbf{g}$. To study the asymptotic behavior of our procedure, we require that

$$Y = G(\mathbf{B}_0^\top \mathbf{f}_0, \epsilon), \quad (2.10)$$

where $G(\cdot, \cdot)$ is an unknown function, $\mathbf{B}_0 \in \mathbb{R}^{p \times d_0}$ with $d_0 \geq 1$, and ϵ is independent of \mathbf{g} .

REMARK 2.5. *Note that $\text{span}(\mathbf{B}_0)$ and \mathbf{f}_0 are generally not identifiable, since $\mathbf{B}_0^\top \mathbf{f}_0 = (\mathbf{\Gamma}\mathbf{B}_0)^\top \mathbf{\Gamma}\mathbf{f}_0$ for any $p \times p$ diagonal matrix $\mathbf{\Gamma}$ containing 1 or -1 in its diagonal. This makes our procedure different from sliced inverse regression assuming \mathbf{f}_0 is known. We also pointed out that if the ultimate goal is the reduction $\mathbf{B}_0^\top \mathbf{f}_0$ itself, then it is unnecessary to impose identifiability constraints. In particular, the structural dimension d_0 remains estimable.*

Let $\hat{\mathbf{B}} \in \mathbb{R}^{p \times d_0}$, $\hat{\mathbf{C}} \in \mathbb{R}^{p \times q}$ and $\hat{\mathbf{f}} = \hat{\mathbf{C}}\mathbf{g}$ denote our estimates of \mathbf{B}_0 , \mathbf{C}_0 and \mathbf{f}_0 , respectively.

THEOREM 2.1. *Suppose $d_0 \leq H \leq K$. Assume that (A1) $E(\mathbf{g}|\mathbf{B}_0^\top \mathbf{C}_0\mathbf{g})$ is linear in $\mathbf{B}_0^\top \mathbf{C}_0\mathbf{g}$ and (A2) $\text{span}[\text{cov}\{E(\mathbf{f}_0|Y)\}] = \text{span}(\mathbf{B}_0)$. Then, for some $p \times p$ diagonal matrix $\mathbf{\Gamma}_n$ containing 1 or -1 in its diagonal, we have $\mathbf{\Gamma}_n \hat{\mathbf{C}} - \mathbf{C}_0 = O_P(n^{-1/2})$ and $\text{span}\{\hat{\mathbf{B}}\}$ is a \sqrt{n} -consistent estimator of $\text{span}\{\mathbf{\Gamma}\mathbf{B}_0\}$.*

The linearity condition (A1) and the coverage condition (A2), which are assumed to hold at the true parameters, are common in regression studies based on methods such as sliced inverse regression. In particular, the linearity condition holds to a reasonable approximation

in many problems (Hall & Li 1993). Nevertheless, these conditions are restrictive due to the effects of predictor transformations, and are imposed mainly for ease of and for simplification of theory investigation. Currently, we are not able to develop the theory in the very general case. As mentioned in Section 2.2, without the linearity condition, our procedure can be interpreted as a method for estimating the canonical correlations between the transformed predictors and the basis functions generated for the response variable, and thus the directions found by it are useful in their own right in identifying some important features of the regression. We plan to conduct theoretical research along this direction.

Let $\{(\hat{\boldsymbol{\theta}}_i, \hat{\boldsymbol{\beta}}_i), i = 1, \dots, H\}$ denote the solution to (2.4). Our procedure thus produces two sets of linear combinations: $\{\boldsymbol{\Phi}\hat{\boldsymbol{\theta}}_i\}$ for the response and $\{\mathbf{X}_{\hat{f}}\hat{\boldsymbol{\beta}}_i\}$ for the predictors. From the dimension-reduction viewpoint, it is always desirable to have a reduced rank solution. That is, we prefer to retain only the ‘‘important’’ projections. However, this is in general a challenging task. Fortunately, under (2.10), the problem reduces to that of determining the structural dimension d_0 . To this end, we propose a BIC-type criterion. More specifically, let r_i^2 denote the squared correlation between $\boldsymbol{\Phi}\hat{\boldsymbol{\theta}}_i$ and $\mathbf{X}_{\hat{f}}\hat{\boldsymbol{\beta}}_i$. Define

$$\text{BIC}_d = \frac{\sum_{i=1}^d r_i^4}{\sum_{i=1}^H r_i^4} - \frac{\log n}{n} \times \frac{d(d+1)}{2}.$$

The estimated structural dimension is then

$$\hat{d} = \arg \max_{1 \leq d \leq H} \text{BIC}_d.$$

REMARK 2.6. *We continue our discussion in Remark 2.1. Consider the linear regression of a transformed response $T(Y)$ on \mathbf{X} . Let $R^2(T) = [\text{cor}\{T(Y), \mathbf{X}^\top \mathbf{b}(T)\}]^2$ be the R -squared value, where $\mathbf{b}(T)$ denotes the least squares solution. Theorem 3.2 of Chen & Li (1998) showed that, at the population level, $R^2(T_i) = \lambda_i, i = 1, \dots, d$. The same applies to the regression of Y on \mathbf{f}_0 . Consequently, r_i^2 as defined above is an estimator of δ_i , the i -th largest eigenvalue of $\text{cov}\{E(\mathbf{f}_0|Y)\}$. Using the equivalence of the optimal scoring problem and the eigen-decomposition problem, it is easily verified that our BIC-type criterion is equivalent to*

that of Zhu et al. (2010) when \mathbf{f}_0 is assumed to be known. The degrees of freedom should be modified accordingly when \mathbf{f}_0 is actually unknown. However, the use of the proposed BIC-type criterion remains feasible, since the estimation of \mathbf{f}_0 is independent of the dimension d .

We have the following corollary.

COROLLARY 2.1. *Under the conditions of Theorem 2.1, $r_i^2 - \delta_i = O_P(n^{-1/2})$ for $i = 1, \dots, H$. Consequently, \hat{d} converges to d_0 in probability.*

Throughout the numerical studies, we initialize β_i and θ_i with solutions to (2.1). That is, we take the ordinary sliced inverse regression estimates as the initial values. As we will see in Section 3, the BIC criterion works well, even in cases where sliced inverse regression fails. Nevertheless, taking the identifiability of transformations into account, we recommend to use multiple initial values, and choose the $\hat{\mathbf{f}}$ and $\hat{\mathbf{B}}$ that correspond to the smallest \hat{d} .

3 Simulation study

In this section, we present some simulation results to illustrate the performance of the proposed nonlinear dimension-reduction method.

Let $\mathbf{0}_p$ be a p -vector of zeros and $\Sigma = (\Sigma_{ij})$ with $\Sigma_{ij} = \rho^{|i-j|}$ for $1 \leq i, j \leq p$. Throughout the simulation study, we take $n = 400, p = 10$, and $\rho = 0.5$. Five simulation examples are considered.

EXAMPLE 3.1. *Let $\beta_{01} = (1, 1, 0, \dots, 0)^\top \in \mathbb{R}^p$ and $\beta_{02} = (0, 0, 1, 1, 0, \dots, 0)^\top \in \mathbb{R}^p$. We consider the model*

$$Y = \frac{\mathbf{f}_0^\top \beta_{01}}{(\mathbf{f}_0^\top \beta_{02} + 1.5)^2 + 0.5} + 0.5 \times \epsilon,$$

where $\mathbf{f}_0 \sim N(\mathbf{0}_p, \Sigma)$, $\epsilon \sim N(0, 1)$, and \mathbf{f}_0 and ϵ are independent. To generate \mathbf{X} , we explore two cases: (i) $X_j = \text{sign}(f_{0j}) \times f_{0j}^2$ and (ii) $X_j = F_C^{-1}\{\Phi(f_{0j})\}$, where $F_C(\cdot)$ is the standard Cauchy distribution function and $\Phi(\cdot)$ is the standard normal distribution function.

EXAMPLE 3.2. This example has the same setting as in Example 3.1 except $\mathbf{f}_0 \sim t_5(\mathbf{0}_p, \Sigma)$, that is, \mathbf{f}_0 has a multivariate t -distribution with location vector zero, scale matrix Σ , and 5 degrees of freedom. We explore two cases: (i) $X_j = \text{sign}(f_{0j}) \times f_{0j}^2$ and (ii) $X_j = F_C^{-1}\{F_{t_5}(f_{0j})\}$, where $F_{t_5}(\cdot)$ denotes the common marginal distribution function of \mathbf{f}_0 .

EXAMPLE 3.3. This example has the same setting as in Example 3.1 except that $\mathbf{f}_0 = \Sigma^{1/2} \mathbf{U}$, where \mathbf{U} is uniformly distributed in the hypercube $[-\sqrt{3}, \sqrt{3}]^p$. To generate \mathbf{X} , we set $X_j = \text{sign}(f_{0j}) \times f_{0j}^2$.

EXAMPLE 3.4. Let $\beta_{01} = (1, 0, \dots, 0)^\top \in \mathbb{R}^p$ and $\beta_{02} = (0, 1, 0, \dots, 0)^\top \in \mathbb{R}^p$. In this example, we first generate $\mathbf{X} \sim N(\mathbf{0}_p, \Sigma)$ and then set $f_{01} = (X_1^2 - 1)/\sqrt{2}$ and $f_{0j} = X_j$ for $j = 2, \dots, p$. Consider again the model

$$Y = \frac{\mathbf{f}_0^\top \beta_{01}}{(\mathbf{f}_0^\top \beta_{02} + 1.5)^2 + 0.5} + 0.5 \times \epsilon,$$

where $\epsilon \sim N(0, 1)$ is independent of \mathbf{X} .

EXAMPLE 3.5. This example has the same setting as in Example 3.4 except

$$Y = (\mathbf{f}_0^\top \beta_{01} + 1) \times \mathbf{f}_0^\top \beta_{02} + 0.5 \times \epsilon.$$

In all the examples considered here, $\mathbf{B}_0 = (\beta_{01}, \beta_{02})$ and the structural dimension $d_0 = 2$. In Examples 3.1-3.3, each component $f_{0j}(X_j)$ is a monotone increasing function of X_j . In Examples 3.4 and 3.5, $f_{01}(X_1)$ is symmetric in X_1 and \mathbf{B}_0 is a sub-matrix of \mathbf{I}_p .

Examples 3.1-3.3 have also been used in Wang et al. (2014), in which probability integral transformed sliced inverse regression (T-SIR) was shown to perform comparatively to the oracle sliced inverse regression (O-SIR) assuming \mathbf{f}_0 is known. Specifically, T-SIR is based

on the assumption that \mathbf{f}_0 is multivariate Gaussian, that is, $f_{0j} = \Phi^{-1}\{F_j(X_j)\}$ with $F_j(\cdot)$ being the marginal distribution function of X_j . This assumption holds in Example 3.1, but is violated in Examples 3.2-3.5. Since probability integral transformation is monotonic, we expect T-SIR to fail in Examples 3.4 and 3.5. Finally, the linearity condition on the distribution of \mathbf{f}_0 holds in Examples 3.1 and 3.2, but is not satisfied in Examples 3.3-3.5.

For comparison purposes, we also examine the performance of T-SIR and O-SIR. For each competitor, we adopt slice indicator functions as transformation functions of the response and fix the number of slices, K , at 10. The numerical results in Wang and Zhu (2013) show that slice indicator functions are quite competitive with B-spline basis functions. In our flexible dimension-reduction method (FDR), two values of the working dimension H are explored: $H = 4$ and $H = K$. The corresponding methods are denoted by P-FDR and F-FDR, respectively. F-FDR is computationally more intensive than P-FDR. For each example, we generate 200 replications.

For any matrix \mathbf{A} , let \mathbf{A}_o be the orthonormalized version of \mathbf{A} . Assume, for the moment, that $d_0 = 2$ is known. To evaluate the accuracy of an estimator $\text{span}(\hat{\mathbf{B}})$ of $\text{span}(\mathbf{B}_0)$, we use both the vector correlation coefficient, defined by $(\prod_{l=1}^2 \phi_l^2)^{1/2}$, and the trace correlation coefficient, defined by $(2^{-1} \sum_{l=1}^2 \phi_l^2)^{1/2}$, where $1 \geq \phi_1^2 \geq \phi_2^2 \geq 0$ are the eigenvalues of the matrix $\hat{\mathbf{B}}_o^\top \mathbf{B}_o \mathbf{B}_o^\top \hat{\mathbf{B}}_o$. To assess how well T-SIR and FDR estimate \mathbf{f}_0 , we calculate the sample Pearson's correlation coefficient between \hat{f}_j and f_{0j} . For simplicity, we concentrate on the set of active transformations, that is, $\{f_{01}, f_{02}, f_{03}, f_{04}\}$ in Examples 3.1-3.3, and $\{f_{01}, f_{02}\}$ in Examples 3.4 and 3.5. The absolute value of this measure is reported for f_{01} in Examples 3.4 and 3.5. Finally, to assess the accuracy of each method in terms of dimension reduction, we use the multiple correlation coefficient (Li & Dong 2009). Suppose \mathbf{V}_1 and \mathbf{V}_2 are two d -dimensional random vectors. Let $\Sigma_{\mathbf{V}_i} = \text{cov}(\mathbf{V}_i), i = 1, 2$ and $\Sigma_{\mathbf{V}_1 \mathbf{V}_2} = \text{cov}(\mathbf{V}_1, \mathbf{V}_2)$. The multiple correlation coefficient between \mathbf{V}_1 and \mathbf{V}_2 is defined as $\{d^{-1} \text{trace}(\Sigma_{\mathbf{V}_1}^{-1} \Sigma_{\mathbf{V}_1 \mathbf{V}_2} \Sigma_{\mathbf{V}_2}^{-1} \Sigma_{\mathbf{V}_1 \mathbf{V}_2}^\top)\}^{1/2}$. We employ the sample version of this measure

based on $\mathbf{X}_{\hat{f}}\hat{\mathbf{B}}$ and $\mathbf{X}_f\mathbf{B}_0$.

The means and standard deviations of various measures are summarized in Tables 1-5. First, we can see, not surprisingly, that O-SIR performs very well in all five examples. The relative insensitivity of sliced inverse regression to the non-linearity of the distribution of the predictors has long been observed in the sufficient dimension-reduction literature. Second, we see that in terms of subspace estimation, the performance of FDR is similar to that of O-SIR. However, FDR is outperformed by O-SIR, with the latter showing consistently higher multiple correlation coefficients. This is because FDR involves the estimation of p component functions, some of which are even redundant. Generally, the estimates of the active component functions are more accurate for monotonic functions than for symmetric ones. Third, T-SIR performs comparably to O-SIR in Examples 3.1-3.3. The reason for this is that the degree of the similarity between the true transformations and the corresponding probability integral transformations is very high: from Tables 1-3 we see that the average Pearson's correlation coefficients of T-SIR are very close to 1. However, T-SIR fails in Examples 3.4 and 3.5, where probability integral transformations, which are by definition monotonic, severely mis-specify transformations that are symmetric. We see from Tables 4 and 5 that the Pearson's correlation coefficients of T-SIR are almost zero for symmetric component functions. Essentially, T-SIR inherits the drawback of sliced inverse regression of failing to capture any symmetric structure. Finally, we see that F-FDR generally outperforms P-FDR, but the differences are often small. Also, unreported results show that the estimation accuracy of FDR improves as the sample size increases.

The performance of each competitor relies heavily on the correct determination of the structural dimension d_0 . Tables 6-10 shows the results of applying the BIC-type criterion. We can see that O-SIR-BIC and FDR-BIC are roughly equally powerful, with O-SIR-BIC performing slightly better. On the other hand, T-SIR-BIC has an alarmingly low rate of correctly identifying the true dimension d_0 in Examples 3.4 and 3.5. Again, this indicates

the disadvantage of using monotonic transformations.

4 Real data analysis

We now apply T-SIR and F-FDR to the Ozone data used in Breiman and Friedman (1985). The goal is to study the relationship between atmospheric ozone concentration and meteorology in the Los Angeles basin. The dataset, which is available from the R package **mlbench**, consists of daily measurements of maximum one-hour-average ozone reading (Y) and eight meteorological variables for $n = 330$ days in 1976. Specifically, the $p = 8$ predictors used in the study are Sandburg Air Force Base temperature (X_1), inversion base height (X_2), Daggett pressure gradient (X_3), visibility (X_4), Vandenburg 500 millibar height (X_5), humidity (X_6), inversion base temperature (X_7), and wind speed (X_8). All the predictors have been scaled to the unit interval $[0, 1]$. Let $\mathbf{X} = (X_1, \dots, X_8)^\top$.

Using the BIC-type criterion, both T-SIR and F-FDR find two directions. The coefficient estimates from either method, denoted by $\hat{\beta}_i$ for $i = 1, 2$, are shown in Table 11. The relative important predictors are X_1, X_2, X_3, X_5, X_6 and X_7 . Let $Z_i = \hat{\beta}_i^\top \mathbf{X}$. The first T-SIR predictor and the first F-FDR predictor are almost identical to the predictor found by the linear least squares fitting (not shown). A linear trend is visible from Figure 1. To check whether the second predictor Z_2 from each method has a significant effect on the response Y , we consider the following model

$$E(Y|Z_1, Z_2) = g_1(Z_1) + g_2(Z_2) + g_{12}(Z_1, Z_2),$$

where $g_1(\cdot)$ and $g_2(\cdot)$ are smooth main effect functions, and $g_{12}(\cdot, \cdot)$ is a smooth interaction term. We use the **gam** function from the R package **mgcv** to fit this model, and find that for both T-SIR and F-FDR, $g_1(Z_1)$ is highly significant (p -values $< 10^{-15}$) and $g_2(Z_2)$ is not significant (p -values > 0.1). Furthermore, at the 0.01 significance level, the interaction

term $g_{12}(Z_1, Z_2)$ for F-FDR is very significant (p -value = 0.004), while that for T-SIR is not (p -value = 0.021). The adjusted R^2 values are 71.6% and 75% for T-SIR and F-FDR, respectively. Using a full two-dimensional smooth does not improve the fit by much (not shown). Thus the second T-SIR direction is likely to be spurious. For F-FDR, the presence of the interaction term is the key to the estimated structural dimension being 2. The adjusted R^2 value from fitting an additive model is 73.8%. We note that, given the transformations of predictors, additive models are essentially one-dimensional. This is also true for alternating conditional expectations of Breiman and Friedman (1985). Therefore, as a dimension-reduction method, our method is far more flexible. Finally, Figure 2 shows the estimated transformations from F-FDR of the six important predictors. We can see that the transformation function of X_5 (or X_6) is symmetric over at least some of the range of observation. This suggests that there may be some symmetric pattern that T-SIR fails to handle. Li (1992) used the same dataset to demonstrate how other dimension-reduction methods are needed to complement the original SIR in symmetric cases. Note that we use $K = 10$ slices here for T-SIR and F-FDR, but other choices yield almost identical scenes.

5 Appendix

PROOF OF PROPOSITION 2.1. Without loss of generality, assume that $\mathbf{B} = (\mathbf{I}_d, \mathbf{B}_l^\top)^\top$, where \mathbf{B}_l is the lower $(p-d) \times d$ sub-matrix of \mathbf{B} . By definition, the conditional distribution of Y given $\mathbf{B}^\top \mathbf{f}(\mathbf{X})$ equals that of Y given \mathbf{X} , that is,

$$F\{y \mid \mathbf{B}^\top \mathbf{f}(\mathbf{X})\} = F(y \mid \mathbf{X}).$$

If $\text{span}(\mathbf{B})$ is not identifiable, then there is another set of univariate functions $\{g_1, \dots, g_p\}$ and a $p \times d$ matrix \mathbf{M} , such that

$$F\{y \mid \mathbf{M}^\top \mathbf{g}(\mathbf{X})\} = F(y \mid \mathbf{X}).$$

Write $\mathbf{B} = (B_{ij})$ and $\mathbf{M} = (M_{ij})$. Then

$$F\{y \mid \mathbf{B}^\top \mathbf{f}(\mathbf{X})\} = h_B\{f_1(x_1) + B_{(1+d)1}f_{1+d}(x_{1+d}) + \cdots + B_{p1}f_p(x_p), \dots, \\ f_d(x_d) + B_{(1+d)d}f_{1+d}(x_{1+d}) + \cdots + B_{pd}f_p(x_p), y\}$$

and

$$F\{y \mid \mathbf{M}^\top \mathbf{g}(\mathbf{X})\} = h_M\{g_1(x_1) + M_{(1+d)1}g_{1+d}(x_{1+d}) + \cdots + M_{p1}g_p(x_p), \dots, \\ g_d(x_d) + M_{(1+d)d}g_{1+d}(x_{1+d}) + \cdots + M_{pd}g_p(x_p), y\}.$$

Here, h_B and h_M are $(d+1)$ -dimensional functions. Consequently,

$$h_B\{f_1(x_1) + B_{(1+d)1}f_{1+d}(x_{1+d}) + \cdots + B_{p1}f_p(x_p), \dots, \\ f_d(x_d) + B_{(1+d)d}f_{1+d}(x_{1+d}) + \cdots + B_{pd}f_p(x_p), y\} \\ = h_M\{g_1(x_1) + M_{(1+d)1}g_{1+d}(x_{1+d}) + \cdots + M_{p1}g_p(x_p), \dots, \\ g_d(x_d) + M_{(1+d)d}g_{1+d}(x_{1+d}) + \cdots + M_{pd}g_p(x_p), y\}.$$

Taking derivatives with respect to x_1, \dots, x_p on both sides, we get

$$h'_{Bj}f'_j = h'_{Mj}g'_j,$$

for $j = 1, \dots, d$, and

$$(h'_{B1}B_{(j+d)1} + \cdots + h'_{Bd}B_{(j+d)d})f'_{j+d} = (h'_{M1}M_{(j+d)1} + \cdots + h'_{Md}M_{(j+d)d})g'_{j+d},$$

for $j = 1, \dots, p-d$, where h'_{Bj} stands for the derivative of h_B with respect to the j th argument, and similarly for h'_{Mj} . Write $\mathbf{h}'_B = (h'_{B1}, \dots, h'_{Bd})^\top$ and $\mathbf{h}'_M = (h'_{M1}, \dots, h'_{Md})^\top$.

Define $\mathbf{D}_{f1} = \text{diag}\{f'_1, \dots, f'_d\}$, $\mathbf{D}_{f2} = \text{diag}\{f'_{1+d}, \dots, f'_p\}$, $\mathbf{D}_{g1} = \text{diag}\{g'_1, \dots, g'_d\}$, and $\mathbf{D}_{g2} = \text{diag}\{g'_{1+d}, \dots, g'_p\}$. Then

$$\mathbf{D}_{f1}\mathbf{h}'_B = \mathbf{D}_{g1}\mathbf{h}'_M$$

and

$$\mathbf{D}_{f2}\mathbf{B}_l\mathbf{h}'_B = \mathbf{D}_{g2}\mathbf{M}_l\mathbf{h}'_M.$$

By the condition of Proposition 2.1, for \boldsymbol{x} in a neighborhood of \boldsymbol{x}^* ,

$$\mathbf{B}_l \boldsymbol{h}'_B = \mathbf{D}_{f_2}^{-1} \mathbf{D}_{g_2} \mathbf{M}_l \mathbf{D}_{g_1}^{-1} \mathbf{D}_{f_1} \boldsymbol{h}'_B.$$

Hence,

$$(\mathbf{B}_l - \mathbf{M}_l) \boldsymbol{h}'_B + (\mathbf{M}_l - \mathbf{D}_{f_2}^{-1} \mathbf{D}_{g_2} \mathbf{M}_l \mathbf{D}_{g_1}^{-1} \mathbf{D}_{f_1}) \boldsymbol{h}'_B = \mathbf{0}.$$

Again by the condition of Proposition 2.1, as $\boldsymbol{x} \rightarrow \boldsymbol{x}^*$,

$$(\mathbf{M}_l - \mathbf{D}_{f_2}^{-1} \mathbf{D}_{g_2} \mathbf{M}_l \mathbf{D}_{g_1}^{-1} \mathbf{D}_{f_1}) \rightarrow \mathbf{0},$$

and hence

$$(\mathbf{B}_l - \mathbf{M}_l) \boldsymbol{h}'_B = \mathbf{0}.$$

Since $d = d_0$, none of the components of \boldsymbol{h}'_B are zero, and there is at most one constant component. This implies that $\mathbf{B}_l = \mathbf{M}_l$. The proof is complete.

PROOF OF THEOREM 2.1. Without loss of generality, assume that $\text{cov}(\boldsymbol{f}) = \mathbf{C} \text{cov}(\boldsymbol{g}) \mathbf{C}^\top$ is equal to \mathbf{I}_p . By the equivalence of linear discriminant analysis and optimal scoring (Hastie et al. 1995), our method amounts to estimating $\text{cov}\{E(\boldsymbol{g}|Y)\}$ and then maximizing

$$\text{trace}[\text{cov}\{E(\boldsymbol{f}|Y)\}] = \text{trace}[\mathbf{C} \text{cov}\{E(\boldsymbol{g}|Y)\} \mathbf{C}^\top]$$

with respect to \mathbf{C} .

From (2.10) we know that $Y \perp\!\!\!\perp \boldsymbol{g} | \mathbf{B}_0^\top \mathbf{C}_0 \boldsymbol{g}$. Then

$$E(\boldsymbol{f}|Y) = E\{E(\boldsymbol{f}|Y, \mathbf{B}_0^\top \mathbf{C}_0 \boldsymbol{g})|Y\} = \mathbf{C} E\{E(\boldsymbol{g} | \mathbf{B}_0^\top \mathbf{C}_0 \boldsymbol{g})|Y\}.$$

Under the linearity condition (A1), $E(\boldsymbol{g} | \mathbf{B}_0^\top \mathbf{C}_0 \boldsymbol{g}) = \text{cov}(\boldsymbol{g}) \mathbf{C}_0^\top \mathbf{B}_0 (\mathbf{B}_0^\top \mathbf{B}_0)^{-1} \mathbf{B}_0^\top \mathbf{C}_0 \boldsymbol{g}$. Hence,

$$E(\boldsymbol{f}|Y) = \mathbf{C} \text{cov}(\boldsymbol{g}) \mathbf{C}_0^\top \mathbf{B}_0 (\mathbf{B}_0^\top \mathbf{B}_0)^{-1} \mathbf{B}_0^\top E(\boldsymbol{f}_0|Y).$$

Setting $\mathbf{C} = \mathbf{C}_0$ yields $E(\boldsymbol{f}_0|Y) = \mathbf{B}_0 (\mathbf{B}_0^\top \mathbf{B}_0)^{-1} \mathbf{B}_0^\top E(\boldsymbol{f}_0|Y)$. This implies that

$$\text{cov}\{E(\boldsymbol{f}|Y)\} = \mathbf{C} \text{cov}(\boldsymbol{g}) \mathbf{C}_0^\top \text{cov}\{E(\boldsymbol{f}_0|Y)\} \mathbf{C}_0 \text{cov}(\boldsymbol{g}) \mathbf{C}^\top.$$

It is easy to check that

$$\begin{aligned} & \text{trace}[\mathbf{C}\text{cov}(\mathbf{g})\mathbf{C}_0^\top \text{cov}\{E(\mathbf{f}_0|Y)\}\mathbf{C}_0\text{cov}(\mathbf{g})\mathbf{C}^\top] \\ & \leq \text{trace}[\{\text{cov}(\mathbf{g})\}^{1/2}\mathbf{C}_0^\top \text{cov}\{E(\mathbf{f}_0|Y)\}\mathbf{C}_0\{\text{cov}(\mathbf{g})\}^{1/2}] \\ & = \text{trace}[\text{cov}\{E(\mathbf{f}_0|Y)\}]. \end{aligned}$$

Consequently, one population solution of \mathbf{C} is $\mathbf{\Gamma}\mathbf{C}_0$, and by the coverage condition (A2), the corresponding solution of $\text{span}\{\mathbf{B}\}$ is $\text{span}\{\mathbf{\Gamma}\mathbf{B}_0\}$. At the sample level, following Li (1991), Fung et al. (2002), it is easily verified that $\widehat{\text{cov}}\{E(\mathbf{g}|Y)\} - \text{cov}\{E(\mathbf{g}|Y)\} = O_P(n^{-1/2})$. This completes the proof.

PROOF OF COROLLARY 2.1. The first part follows immediately from Theorem 2.1. For a proof of the second part, see Zhu et al. (2010).

REFERENCES

- Box, G. E. & Cox, D. R. (1964), ‘An analysis of transformations’, *Journal of the Royal Statistical Society, Series B* **26**(2), 211–252.
- Breiman, L. & Friedman, J. H. (1985), ‘Estimating optimal transformations for multiple regression and correlation’, *Journal of the American statistical Association* **80**(391), 580–598.
- Chen, C.-H. & Li, K.-C. (1998), ‘Can SIR be as popular as multiple linear regression?’, *Statistica Sinica* **8**(2), 289–316.
- Cook, R. D. (1998), *Regression graphics: Ideas for studying regressions through graphics*, John Wiley & Sons, New York.
- Duan, N. & Li, K.-C. (1991), ‘Slicing regression: a link-free regression method’, *The Annals of Statistics* **19**(2), 505–530.

- Fung, W. K., He, X., Liu, L. & Shi, P. (2002), ‘Dimension reduction based on canonical correlation’, *Statistica Sinica* **12**(4), 1093–1114.
- Hall, P. & Li, K.-C. (1993), ‘On almost linearity of low dimensional projections from high dimensional data’, *The Annals of Statistics* **21**(2), 867–889.
- Hastie, T., Buja, A. & Tibshirani, R. (1995), ‘Penalized discriminant analysis’, *The Annals of Statistics* **23**(1), 73–102.
- Hastie, T., Tibshirani, R. & Buja, A. (1994), ‘Flexible discriminant analysis by optimal scoring’, *Journal of the American Statistical Association* **89**(428), 1255–1270.
- Hastie, T., Tibshirani, R. & Friedman, J. (2009), *The Elements of Statistical Learning*, Springer, New York.
- He, X. & Shen, L. (1997), ‘Linear regression after spline transformation’, *Biometrika* **84**(2), 474–481.
- Kent, J. T. (1991), ‘Comment’, *Journal of the American Statistical Association* **86**(414), 336–337.
- Li, B. & Dong, Y. (2009), ‘Dimension reduction for nonelliptically distributed predictors’, *The Annals of Statistics* **37**(3), 1272–1298.
- Li, K.-C. (1991), ‘Sliced inverse regression for dimension reduction’, *Journal of the American Statistical Association* **86**(414), 316–327.
- Li, K.-C. (1992), ‘On principal hessian directions for data visualization and dimension reduction: another application of stein’s lemma’, *Journal of the American Statistical Association* **87**(420), 1025–1039.
- Ramsay, J. O. (1988), ‘Monotone regression splines in action’, *Statistical Science* **3**(4), 425–441.

- Wang, T., Guo, X., Xu, P. & Zhu, L. (2014), ‘Transformed sufficient dimension reduction’, *Biometrika* **101**(4), 815–829.
- Wang, T. & Zhu, L. (2013), ‘Sparse sufficient dimension reduction using optimal scoring’, *Computational Statistics & Data Analysis* **57**(1), 223–232.
- Wood, S. (2006), *Generalized Additive Models: An Introduction with R*, Chapman & Hall, Boca Raton.
- Wu, Y. & Li, L. (2011), ‘Asymptotic properties of sufficient dimension reduction with a diverging number of predictors’, *Statistica Sinica* **2011**(21), 707–730.
- Yin, X. & Li, B. (2011), ‘Sufficient dimension reduction based on an ensemble of minimum average variance estimators’, *The Annals of Statistics* **39**(6), 3392–3416.
- Zhu, L., Wang, T., Zhu, L. & Ferré, L. (2010), ‘Sufficient dimension reduction through discretization-expectation estimation’, *Biometrika* **97**(2), 295–304.

Table 1: Means and standard deviations (in parentheses) of the vector correlation coefficient (VCC) and the trace correlation coefficient (TCC) for subspace estimation, the Pearson’s correlation coefficient (PCC) for component estimation, and the multiple correlation coefficient (MCC) for dimension reduction, based on 200 data replications, are reported for Example 3.1

		VCC	TCC	PCC (\hat{f}_1)	PCC (\hat{f}_2)	PCC (\hat{f}_3)	PCC (\hat{f}_4)	MCC
Case (i)	O-SIR	0.817 (0.081)	0.910 (0.037)					0.956 (0.021)
	T-SIR	0.825 (0.077)	0.914 (0.036)	0.998 (0.001)	0.998 (0.001)	0.998 (0.001)	0.997 (0.001)	0.956 (0.020)
	P-FDR	0.815 (0.078)	0.909 (0.036)	0.960 (0.026)	0.932 (0.043)	0.872 (0.104)	0.867 (0.104)	0.878 (0.032)
	F-FDR	0.804 (0.093)	0.904 (0.040)	0.959 (0.030)	0.929 (0.048)	0.889 (0.094)	0.882 (0.093)	0.882 (0.032)
Case (ii)	O-SIR	0.825 (0.083)	0.914 (0.038)					0.959 (0.019)
	T-SIR	0.830 (0.080)	0.916 (0.037)	0.998 (0.001)	0.998 (0.001)	0.998 (0.001)	0.998 (0.001)	0.959 (0.018)
	P-FDR	0.824 (0.086)	0.913 (0.039)	0.961 (0.024)	0.931 (0.042)	0.880 (0.097)	0.858 (0.099)	0.880 (0.032)
	F-FDR	0.822 (0.090)	0.912 (0.041)	0.960 (0.026)	0.929 (0.045)	0.898 (0.087)	0.877 (0.089)	0.887 (0.032)

Table 2: Means and standard deviations (in parentheses) of the vector correlation coefficient (VCC) and the trace correlation coefficient (TCC) for subspace estimation, the Pearson’s correlation coefficient (PCC) for component estimation, and the multiple correlation coefficient (MCC) for dimension reduction, based on 200 data replications, are reported for Example 3.2

		VCC	TCC	PCC (\hat{f}_1)	PCC (\hat{f}_2)	PCC (\hat{f}_3)	PCC (\hat{f}_4)	MCC
Case (i)	O-SIR	0.709 (0.128)	0.859 (0.056)					0.926 (0.036)
	T-SIR	0.782 (0.095)	0.893 (0.043)	0.979 (0.015)	0.979 (0.018)	0.977 (0.018)	0.979 (0.013)	0.926 (0.029)
	P-FDR	0.788 (0.094)	0.895 (0.043)	0.913 (0.040)	0.890 (0.046)	0.804 (0.111)	0.781 (0.131)	0.813 (0.043)
	F-FDR	0.786 (0.101)	0.895 (0.045)	0.912 (0.041)	0.890 (0.045)	0.827 (0.103)	0.804 (0.117)	0.821 (0.041)
Case (ii)	O-SIR	0.719 (0.127)	0.863 (0.055)					0.924 (0.038)
	T-SIR	0.789 (0.091)	0.896 (0.041)	0.977 (0.015)	0.976 (0.018)	0.978 (0.015)	0.979 (0.014)	0.924 (0.028)
	P-FDR	0.788 (0.089)	0.895 (0.042)	0.909 (0.033)	0.881 (0.052)	0.782 (0.129)	0.803 (0.118)	0.808 (0.040)
	F-FDR	0.784 (0.108)	0.894 (0.048)	0.909 (0.033)	0.884 (0.049)	0.819 (0.117)	0.814 (0.115)	0.818 (0.040)

Table 3: Means and standard deviations (in parentheses) of the vector correlation coefficient (VCC) and the trace correlation coefficient (TCC) for subspace estimation, the Pearson’s correlation coefficient (PCC) for component estimation, and the multiple correlation coefficient (MCC) for dimension reduction, based on 200 data replications, are reported for Example 3.3

	VCC	TCC	PCC (\hat{f}_1)	PCC (\hat{f}_2)	PCC (\hat{f}_3)	PCC (\hat{f}_4)	MCC
O-SIR	0.844 (0.069)	0.923 (0.032)					0.962 (0.017)
T-SIR	0.801 (0.086)	0.902 (0.040)	0.988 (0.002)	0.991 (0.001)	0.991 (0.002)	0.991 (0.002)	0.949 (0.019)
P-FDR	0.836 (0.086)	0.919 (0.038)	0.983 (0.019)	0.947 (0.042)	0.895 (0.084)	0.885 (0.088)	0.899 (0.031)
F-FDR	0.831 (0.094)	0.917 (0.042)	0.982 (0.021)	0.949 (0.043)	0.907 (0.075)	0.904 (0.079)	0.904 (0.031)

Table 4: Means and standard deviations (in parentheses) of the vector correlation coefficient (VCC) and the trace correlation coefficient (TCC) for subspace estimation, the Pearson’s correlation coefficient (PCC) for component estimation, and the multiple correlation coefficient (MCC) for dimension reduction, based on 200 data replications, are reported for Example 3.4

	VCC	TCC	PCC (\hat{f}_1)	PCC (\hat{f}_2)	MCC
O-SIR	0.807 (0.095)	0.906 (0.042)			0.935 (0.030)
T-SIR	0.562 (0.209)	0.787 (0.088)	0.064 (0.048)	0.998 (0.001)	0.643 (0.034)
P-FDR	0.823 (0.099)	0.914 (0.043)	0.742 (0.070)	0.906 (0.051)	0.776 (0.050)
F-FDR	0.834 (0.102)	0.919 (0.044)	0.739 (0.084)	0.908 (0.049)	0.782 (0.049)

Table 5: Means and standard deviations (in parentheses) of the vector correlation coefficient (VCC) and the trace correlation coefficient (TCC) for subspace estimation, the Pearson’s correlation coefficient (PCC) for component estimation, and the multiple correlation coefficient (MCC) for dimension reduction, based on 200 data replications, are reported for Example 3.5

	VCC	TCC	PCC (\hat{f}_1)	PCC (\hat{f}_2)	MCC
O-SIR	0.959 (0.017)	0.979 (0.008)			0.987 (0.005)
T-SIR	0.691 (0.216)	0.867 (0.073)	0.073 (0.053)	0.997 (0.001)	0.691 (0.008)
P-FDR	0.929 (0.049)	0.965 (0.022)	0.665 (0.140)	0.945 (0.028)	0.824 (0.043)
F-FDR	0.935 (0.049)	0.967 (0.022)	0.661 (0.155)	0.945 (0.031)	0.826 (0.047)

Table 6: Percentages of estimated structural dimensions \hat{d} smaller than, equal to, and larger than d_0 , based on 200 data replications, are reported for Example 3.1

		O-SIR	T-SIR	P-FDR	F-FDR
Case (i)	$\hat{d} < d_0$	0	0	0	0
	$\hat{d} = d_0$	0.99	0.995	0.94	0.98
	$\hat{d} > d_0$	0.01	0.005	0.06	0.02
Case (ii)	$\hat{d} < d_0$	0	0	0	0
	$\hat{d} = d_0$	0.99	0.995	0.93	0.975
	$\hat{d} > d_0$	0.01	0.005	0.07	0.025

Table 7: Percentages of estimated structural dimensions \hat{d} smaller than, equal to, and larger than d_0 , based on 200 data replications, are reported for Example 3.2

		O-SIR	T-SIR	P-FDR	F-FDR
Case (i)	$\hat{d} < d_0$	0.01	0.005	0	0
	$\hat{d} = d_0$	0.98	0.995	0.97	0.995
	$\hat{d} > d_0$	0.01	0	0.03	0.005
Case (ii)	$\hat{d} < d_0$	0.005	0.005	0	0
	$\hat{d} = d_0$	0.985	0.985	0.925	0.98
	$\hat{d} > d_0$	0.010	0.010	0.075	0.02

Table 8: Percentages of estimated structural dimensions \hat{d} smaller than, equal to, and larger than d_0 , based on 200 data replications, are reported for Example 3.3

	O-SIR	T-SIR	P-FDR	F-FDR
$\hat{d} < d_0$	0	0	0	0
$\hat{d} = d_0$	1	0.995	0.965	0.98
$\hat{d} > d_0$	0	0.005	0.035	0.02

Table 9: Percentages of estimated structural dimensions \hat{d} smaller than, equal to, and larger than d_0 , based on 200 data replications, are reported for Example 3.4

	O-SIR	T-SIR	P-FDR	F-FDR
$\hat{d} < d_0$	0	0	0	0
$\hat{d} = d_0$	0.98	0.265	0.88	0.92
$\hat{d} > d_0$	0.02	0.735	0.12	0.08

Table 10: Percentages of estimated structural dimensions \hat{d} smaller than, equal to, and larger than d_0 , based on 200 data replications, are reported for Example 3.5

	O-SIR	T-SIR	P-FDR	F-FDR
$\hat{d} < d_0$	0	0.89	0	0
$\hat{d} = d_0$	1	0.11	0.995	0.995
$\hat{d} > d_0$	0	0	0.005	0.005

Table 11: Ozone data. Estimated directions by T-SIR and F-FDR

	T-SIR		F-FDR	
	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_1$	$\hat{\beta}_2$
X_1	0.454	0.472	0.349	0.619
X_2	-0.102	-0.302	-0.230	-0.352
X_3	0.122	0.082	0.231	0.052
X_4	-0.081	-0.043	-0.110	-0.058
X_5	-0.059	-0.093	0.065	-0.173
X_6	0.129	-0.214	0.042	-0.218
X_7	0.242	-0.732	0.188	-0.784
X_8	0.012	-0.049	-0.014	-0.007

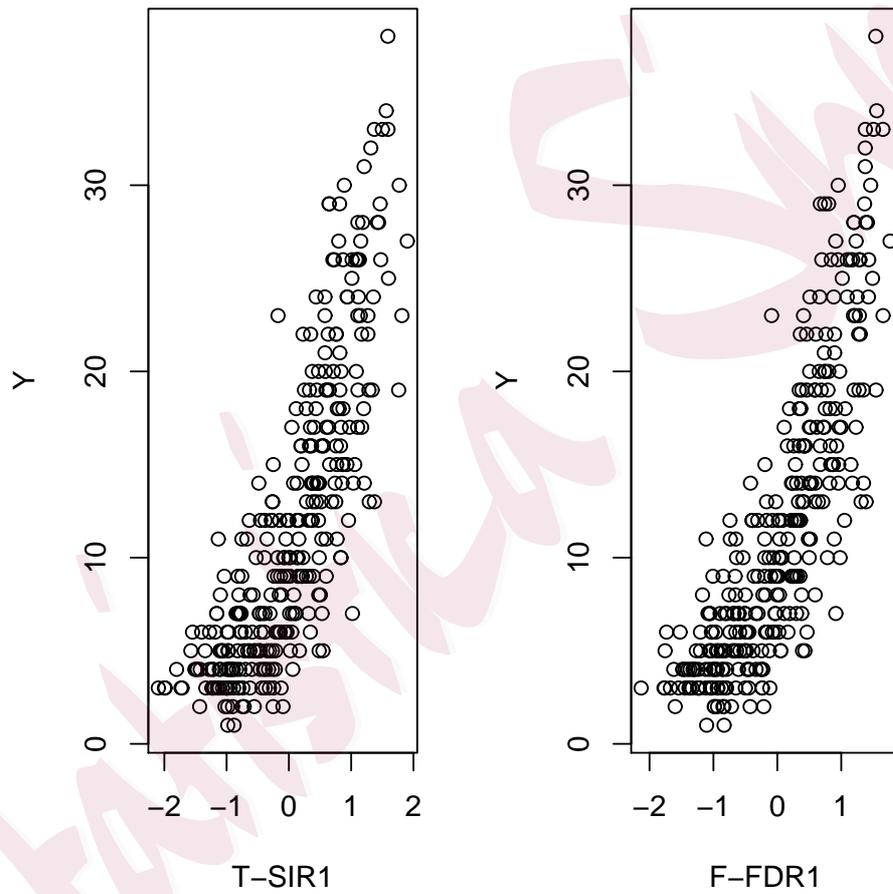


Figure 1: Ozone concentration against the first T-SIR predictor (left panel) and Ozone concentration against the first F-FDR predictor (right panel)

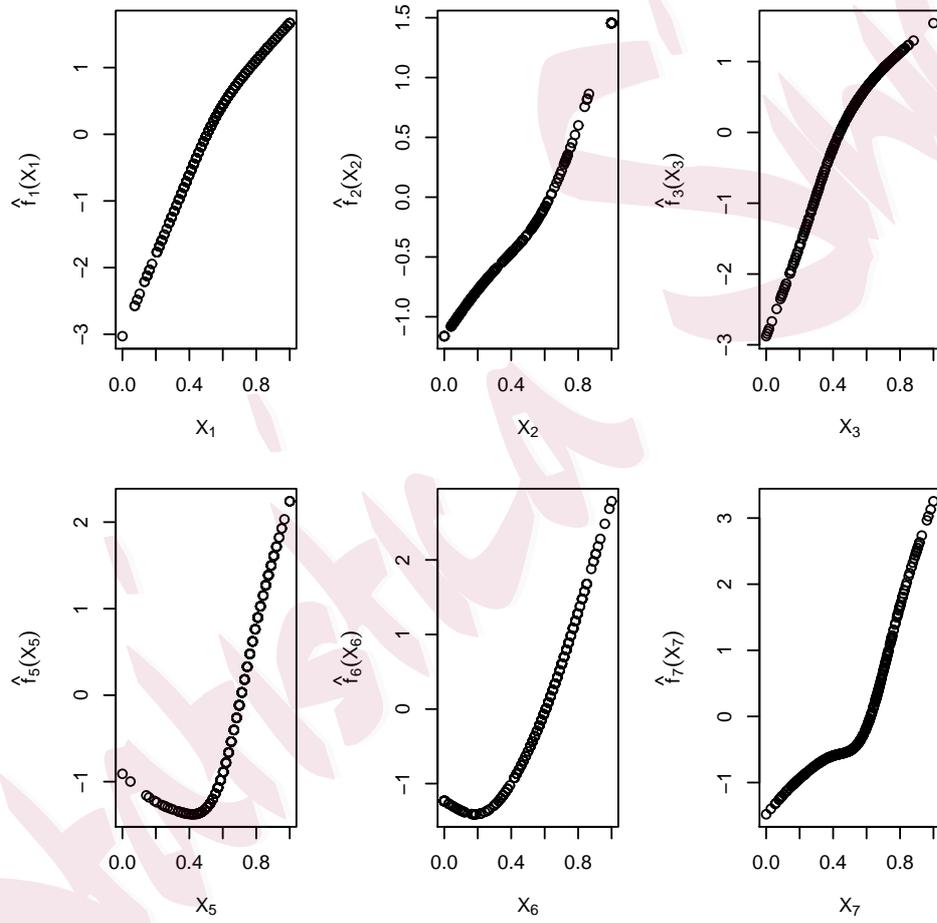


Figure 2: Plots of the estimated transformation functions by F-FDR