



Minerva Access is the Institutional Repository of The University of Melbourne

Author/s:

Prendergast, LA;Garnham, AL

Title:

Response and predictor folding to counter symmetric dependency in dimension reduction

Date:

2016-12-01

Citation:

Prendergast, L. A. & Garnham, A. L. (2016). Response and predictor folding to counter symmetric dependency in dimension reduction. Australian and New Zealand Journal of Statistics, 58 (4), pp.515-532. <https://doi.org/10.1111/anzs.12170>.

Persistent Link:

<https://hdl.handle.net/11343/292157>

Author Manuscript

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/anzs.12170](https://doi.org/10.1111/anzs.12170)

This article is protected by copyright. All rights reserved

RESPONSE AND PREDICTOR FOLDING TO COUNTER SYMMETRIC DEPENDENCY IN DIMENSION REDUCTION

L. A. PRENDERGAST^{1*} AND A. L. GARNHAM²

La Trobe University, Walter and Eliza Hall Institute of Medical Research

Summary

In the regression setting, dimension reduction allows for complicated regression structures to be detected via visualisation in a low-dimensional framework. However, some popular dimension reduction methodologies fail to achieve this aim when faced with a problem often referred to as symmetric dependency. In this paper we show how vastly superior results can be achieved when carrying out response and predictor transformations for methods such as least squares and sliced inverse regression. These transformations are simple to implement and utilise estimates from other dimension reduction methods that are not faced with the symmetric dependency problem. We highlight the effectiveness of our approach via simulation and an example. Furthermore, we show that ordinary least squares can effectively detect multiple dimension reduction directions. Methods robust to extreme response values are also considered.

Key words: cumulative slicing estimation; ordinary least squares, principal Hessian directions, robust M -estimation, sliced inverse regression, sliced average variance estimates

1. Introduction

Let $Y \in \mathbb{R}$ denote a random univariate response and $\mathbf{x} \in \mathbb{R}^p$ a random p -dimensional vector of predictors. [Li & Duan \(1989\)](#) considered the model

$$Y = f(\boldsymbol{\beta}^\top \mathbf{x}, \varepsilon) \quad (1)$$

where $\boldsymbol{\beta}$ is an unknown p -dimensional vector of predictor coefficients, f is the unknown link function and ε is the error term that is assumed to be independent of \mathbf{x} . Of interest is the regression function $E(Y|\mathbf{x})$ where, ideally, a plot of Y versus \mathbf{x} can reveal the form of f . However, we are limited in this sense when p is large due to our inability to visualise objects in high dimensions. Importantly, Y depends on \mathbf{x} only through $\boldsymbol{\beta}^\top \mathbf{x}$ so that if we could determine $\boldsymbol{\beta}$ then we could replace the p -dimensional \mathbf{x} with the one-dimensional $\boldsymbol{\beta}^\top \mathbf{x}$.

* Author to whom correspondence should be addressed.

¹Department of Mathematics and Statistics, La Trobe University, Melbourne, Australia, 3086. luke.prendergast@latrobe.edu.au

² Division of Bioinformatics, Walter and Eliza Hall Institute of Medical Research

Acknowledgements. The authors are very thankful to two anonymous referees and an Associate Editor for their useful comments and suggestions that lead to a clearer and more accessible manuscript.

Our ability to explore possibilities for f would then be enhanced due to the resulting lower-dimensional framework.

In the sample setting, let $\{y_i, \mathbf{x}_i\}_{i=1}^n$ be n sample realisations of Y and \mathbf{x} where the relationship between Y and \mathbf{x} is assumed to be of the form given in (1). Suppose that β can be estimated and let this estimate be denoted $\hat{\beta}$. Then the y_i s can be plotted against the $\hat{\beta}^\top \mathbf{x}_i$ s to visually determine f . Such a plot is called an *Estimated Sufficient Summary Plot* (ESSP, see, e.g., Cook 1998b). The focus of our work here will be to obtain good ESSPs in settings for which estimation of β is difficult.

While our initial focus will be on the model in (1), later we will also look at a further generalised model studied by Li (1991) and given as

$$Y = f(\beta_1^\top \mathbf{x}, \dots, \beta_K^\top \mathbf{x}, \varepsilon) \quad (2)$$

which allows for more than one vector of predictor coefficients.

Li & Duan (1989) extended earlier works by Brillinger (1977, 1983) to show that ordinary least squares (OLS), and robust versions, can be used to estimate the direction of β when the model is of the form (1) and under some mild conditions on \mathbf{x} . We will provide a brief review of these results in Section 2. However, for some forms of f OLS is not expected to find β . Consequently we also discuss another approach, Principal Hessian Directions (PHD) introduced by Li (1992), which may often be more suitable for models with these forms of f . In Section 3 we propose a simple transformation of the response based on an initial PHD estimate that can be used to ensure that OLS can provide a good ESSP. Simulation results, described in Section 4, demonstrate that this approach can be used to obtain vastly superior estimates. Extensions to other approaches are discussed in Section 5. These have applications to the multi-index model given in (2). Finally, an example is provided in Section 6 and the paper is concluded with a discussion in Section 7.

2. Methods

In the context of the multi-index model in (2), Li (1991) defined the *Linear Design Condition* to be:

Condition 1. For any $\mathbf{c} \in \mathbb{R}^p$, $E(\mathbf{c}^\top \mathbf{x} | \beta_1^\top \mathbf{x}, \dots, \beta_K^\top \mathbf{x})$ is linear in $\beta_1^\top \mathbf{x}, \dots, \beta_K^\top \mathbf{x}$.

While this condition holds when \mathbf{x} belongs to the family of elliptically symmetric distributions, Hall & Li (1993) showed that Condition 1 often approximately holds in practice when p is large. One also has the possibility to utilise predictor transformations to ensure that it approximately holds (see, e.g., Fox & Weisberg 2011).

Li & Duan (1989) had earlier considered Condition 1 in the context of the model in (1) (i.e. when $K = 1$ in (2)). This will be our focus for the remainder of this section and we will revisit the case of general K again later.

2.1. Least squares and similar approaches

When Condition 1 and the model in (1) hold, Li & Duan (1989) showed that the OLS slope vector, denoted $\mathbf{b} = \text{var}(\mathbf{x})^{-1}\text{Cov}(\mathbf{x}, Y)$, is equal to $c\boldsymbol{\beta}$ for some $c \in \mathbb{R}$. Consequently, OLS can recover the direction of $\boldsymbol{\beta}$ when $c \neq 0$ and a plot of Y versus $\mathbf{b}^\top \mathbf{x}$ can be used to seek f . It should be pointed out that any non-zero c is adequate since any \mathbf{b} in the direction of $\boldsymbol{\beta}$ is suitable for finding an appropriate link function. In practice, OLS can be used to obtain $\hat{\mathbf{b}}$, the usual OLS slope estimate, and an ESSP created using the y_i s and the $\hat{\mathbf{b}}^\top \mathbf{x}_i$ s. While OLS is one simple approach (e.g. Brillinger 1983), Li and Duan generalised Brillinger's results to include estimators satisfying

$$\underset{a, \mathbf{b}}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \rho(a + \mathbf{b}^\top \mathbf{x}_i, y_i) \quad (3)$$

provided that the function ρ is convex in its first argument and that a solution exists. Hence, other possibilities can be robust estimators such as M -estimators with the Huber weight function (Huber 1973). While the temptation would be to consider a robust approach only when possible errors are present in the data set, Prendergast & Sheather (2013) showed that for some models the robust estimators can outperform OLS even when data is sampled without error. Similarly, Prendergast (2008) used trimming of influential observations to also improve estimates.

If an estimator is expected to find the direction of $\boldsymbol{\beta}$, then it is required that $c \neq 0$. Some discussion of the case $c = 0$ can be found in Li (1991) and Cook & Weisberg (1991). While these discussions are for a different method they can similarly be applied to OLS. That is, when the link function f is symmetric about the mean of $\boldsymbol{\beta}^\top \mathbf{x}$, then $\mathbf{b} = \mathbf{0}$. To highlight this, we consider two simulated examples. The first does not have the symmetric dependency issue while the second does. The models we use are

Model 1. $Y = \sin(0.7\boldsymbol{\beta}^\top \mathbf{x}) + 0.35\varepsilon$

Model 2. $Y = \cos(0.7\boldsymbol{\beta}^\top \mathbf{x}) + 0.35\varepsilon$

where, for both models, \mathbf{x} is a 10 dimensional random vector which is distributed as $N(\mathbf{0}, \mathbf{I})$, $\varepsilon \sim N(0, 1)$ and $\boldsymbol{\beta} = (1, -2, 0, \dots, 0)^\top$.

***** Figure 1 about here *****

In Figure 1 we provide true views (where the y_i s are plotted against the ideally dimension reduced x_i s, i.e. the $\beta^\top x_i$ s) and ESSPs where OLS has been used as the estimator. Plots A and B are for Model 1 and Plots C and D for Model 2 where, in both cases, $n = 200$ observations have been randomly generated. If OLS has performed well, then we would expect the ESSP to look similar to the true views, although with possible differences in scale on the horizontal axis, since OLS is targeting $c\beta$ for a c that is not necessarily one. For Model 1, we can see that OLS has performed exceptionally well, producing an excellent ESSP as seen in Plot B. However, OLS has failed for Model 2 with an ESSP in Plot D that does not provide any evidence of a relationship between the responses and dimension reduced predictors. The true view shows that there is certainly something to find. Recall that Model 2 exhibits symmetric dependency and OLS is trying to estimate $0 \times \beta$.

2.2. Principal Hessian Directions

Li (1992) introduced Principal Hessian Directions (PHD), a method that does not suffer from the symmetric dependency problem and one that is also capable of finding multiple vectors of predictor coefficients. That is, the model can be assumed to be of the form shown in (2) where $K \geq 1$, in which case it is desirable to find a basis for the span of the β_k s.

Consider the following condition required by PHD.

Condition 2. $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

When Condition 2 holds, imposing normality of the predictor, Condition 1 also holds. As a consequence, if PHD is applicable due to this condition being met, then so too is OLS. There are slightly weaker conditions needed for PHD, namely that $\text{var}(\mathbf{x}|\beta^\top \mathbf{x})$ is constant in conjunction with assuming that Condition 1 holds. Recent work by Leeb (2013) showed that this will often hold approximately in practice.

While OLS returns an estimate of $\boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma}_{xy}$ (where $\boldsymbol{\Sigma}_{xy}$ is the covariance vector between \mathbf{x} and Y) as an estimate of the direction of β , PHD instead carries out an eigen-decomposition of an estimate of $\bar{\mathbf{H}} = \boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma}_{yxx}\boldsymbol{\Sigma}^{-1}$, where

$$\boldsymbol{\Sigma}_{yxx} = E \{ (Y - E(Y))(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \}.$$

When Condition 2 holds, for many models satisfying (1), the rank of $\bar{\mathbf{H}}$ is one and the eigenvector corresponding to the non-zero eigenvalue is in the same direction as β . We have emphasised “many models” here since PHD will not be able to find the direction of β when there is an odd symmetric dependency between Y and the mean of $\beta^\top \mathbf{x}$. Here odd symmetric dependency refers to the type of symmetry seen for Model 1 (see Figure 1). Li (1991) also noted that Y can be replaced by the OLS residual without changing $\bar{\mathbf{H}}$ where, notationally,

we replace Σ_{yxx} by Σ_{rxx} to distinguish between the two approaches. While \bar{H} does not change, the estimator is influenced. Empirical and theoretical results have suggested that this residual-based PHD approach is often a better estimator of β (Cook 1998a; Prendergast & Smith 2010). Consequently the residual-based PHD will be our method of choice.

Since Σ_{rxx} is moment-based, estimation is straightforward. Let r_1, \dots, r_n denote the usual OLS residuals for the regression of the y_i s on the x_i s and also let \bar{x} be the sample mean of the x_i s. Then the estimate of Σ_{rxx} is

$$\hat{\Sigma}_{rxx} = \frac{1}{n} \sum_{i=1}^n r_i (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top.$$

Lue (2001) has previously shown that, similarly to OLS estimation, trimming can improve PHD estimation.

3. Predictor and response transformations to remove symmetric dependency

Garnham & Prendergast (2013) showed that response transformations can greatly improve OLS estimates. Their results, however, do not solve the issue of the symmetric dependency that troubles OLS. The aim here is to introduce two transformation functions, one for the response and the other for the predictor, that can be useful in the symmetric dependency setting.

3.1. Theory

We start by introducing the transformations by discussing their motivation. The response transformation that we will focus on is

$$t_y(Y; \mathbf{v}) = \begin{cases} Y, & \mathbf{v}^\top \mathbf{x} > \mathbf{v}^\top \boldsymbol{\mu} \\ Y - 2(Y - \mathbb{E}(Y | \mathbf{v}^\top \mathbf{x} = \mathbf{v}^\top \boldsymbol{\mu})), & \mathbf{v}^\top \mathbf{x} \leq \mathbf{v}^\top \boldsymbol{\mu} \end{cases} \quad (4)$$

where \mathbf{v} needs to be chosen. If $\mathbf{v} = c_1 \boldsymbol{\beta}$ for a non-zero scalar c_1 , then $t_y(Y; \mathbf{v})$ and \mathbf{x} still satisfy the model in (1) and Condition 1. An estimator of $\boldsymbol{\beta}$ is then the OLS slope vector estimator for the regression of $t_y(Y; \mathbf{v})$ on \mathbf{x} . Soon we will show that, in the empirical setting, good estimates of $\boldsymbol{\beta}$ used in the transformation function can generate vastly improved results. That is, if we can find a good estimate of $\boldsymbol{\beta}$, then we can use this estimate to transform the response and obtain a vastly improved estimate of $\boldsymbol{\beta}$ for some models. We also consider the following predictor transformation function

$$t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v}) = \text{sign}(\mathbf{v}^\top \mathbf{x} - \mathbf{v}^\top \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu}). \quad (5)$$

***** Figure 2 about here *****

In Figure 2 we show the effects of the transformations on the curves defining Y in Model 2 under the assumption of zero error. In Plot A the original curve for Model 2 is displayed as the thick grey line. It is clear that this curve is symmetric around $\beta^\top \mu = 0$, which is marked by the vertical line. We can also see that $E(Y|\beta^\top \mathbf{x} = 0) = 1$, and this is shown as the horizontal line on the plot. The black line is the curve resulting from the transformation in (4). The curve in the lower-left quadrant is folded up to the top-left quadrant thus removing the symmetric dependency. Similarly, in Plot B we show the effect of employing the predictor transformation in (5). Here, the curve in the lower-left quadrant is folded over to the bottom-right quadrant. Again, the symmetric dependency is removed. For further clarity, in Plots C and D we show how the transformations work with respect to 500 observations generated according to Model 2. The original data are the grey points and the transformed data are the black points. As can be seen in these plots, the transformed data do not suffer from symmetric dependency

In Theorem 1 below, we show that for certain choices of \mathbf{v} , the transformation in (5) can be useful for finding directions hidden by symmetric dependency. The result is provided in the more general context of the multi-index model in (2). The proof is in the Appendix.

Theorem 1. *Consider the predictor transformation given in (5) and let $\mathbf{v} \in \mathcal{S}$ where \mathcal{S} is the column space of the β_1, \dots, β_K . Under the model in (2) and Condition 1,*

$$\text{var}(\mathbf{x})^{-1} \text{cov}(t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v}), Y) \in \mathcal{S} \quad (6)$$

$$\text{var}(t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v})) - 1 \text{cov}(t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v}), Y) \in \mathcal{S}. \quad (7)$$

The estimator in (7) is simply the OLS slope from the regression of Y on $t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v})$. The estimator in (6) is similar although it utilises the variance estimator for the original \mathbf{x} .

3.2. Application in practice

Recall that our sample of n observations is denoted $\{y_i, \mathbf{x}_i\}_{i=1}^n$. Throughout let \bar{y} , $\bar{\mathbf{x}}$, \mathbf{S}_x and \mathbf{S}_{xy} denote the sample mean of the y_i s, sample mean of the \mathbf{x}_i s, sample covariance matrix of the \mathbf{x}_i s and the sample covariance between the \mathbf{x}_i s and y_i s respectively. Also let \mathbf{X} denote the $n \times p$ design matrix whose i -th row is \mathbf{x}_i .

There are two points that need clarification prior to application in practice. Firstly, how to choose an appropriate vector \mathbf{v} ? Our simulations indicate that for many models, OLS is a better estimator of β in (1) than PHD. However, PHD is preferred when symmetric dependency is evident, in which case OLS can struggle to find β . Consequently, in practice we propose to set $\mathbf{v} = \hat{\mathbf{b}}_{phd}$ - the PHD estimate of β . In the next section, our results show

that reasonable PHD estimates of β can lead to much improved estimates of β when OLS is employed following the transformations in Section 3.1.

Secondly, the response transformation requires estimation of $E(Y|\mathbf{v}^\top \mathbf{x} = \mathbf{v}^\top \boldsymbol{\mu})$. We propose to find an approximation to this estimate as

$$\bar{y}(\mathbf{v}) = \hat{f}(\mathbf{v}^\top \bar{\mathbf{x}}) \quad (8)$$

where \hat{f} is a fitted spline for the y_i s versus the $\mathbf{v}^\top \mathbf{x}_i$ s. We use a cubic smoothing spline (for details see, e.g., Faraway 2006). If computational issues arise with fitting the spline, such as can happen when n is small, then a simple option is to use $\sum_{j \in I_m} y_j / m$ with, for example, $m = 10$ and where I_m is the set of indices for the closest m $\mathbf{v}^\top \mathbf{x}_i$ s to $\mathbf{v}^\top \bar{\mathbf{x}}$.

The transformations we employ are then

$$y_i^* = t_{y_i}(y_i; \hat{\mathbf{b}}_{phd}) = \begin{cases} y_i, & \hat{\mathbf{b}}_{phd}^\top \mathbf{x}_i > \hat{\mathbf{b}}_{phd}^\top \bar{\mathbf{x}} \\ y_i - 2(y_i - \bar{y}(\hat{\mathbf{b}}_{phd})), & \hat{\mathbf{b}}_{phd}^\top \mathbf{x}_i \leq \hat{\mathbf{b}}_{phd}^\top \bar{\mathbf{x}} \end{cases} \quad (9)$$

and

$$\mathbf{x}_i^* = t_{x_i}(\mathbf{x}_i - \bar{\mathbf{x}}; \hat{\mathbf{b}}_{phd}) = \text{sign}(\hat{\mathbf{b}}_{phd}^\top \mathbf{x}_i - \hat{\mathbf{b}}_{phd}^\top \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}}) \quad (10)$$

where we use the notation y_i^* and \mathbf{x}_i^* for convenience. The methods we use are:

Method 1. *The OLS slope vector for the regression of the y_i^* s on the \mathbf{x}_i s.*

Method 2. *The OLS slope vector for the regression of the y_i s on the \mathbf{x}_i^* s.*

Potentially, a combination of the transformations could also be used. However, our simulations revealed that better results are achieved by using only one at a time. Another possibility is to use $S_x^{-1} S_{xy}^*$. However, our simulations also revealed that this approach was very typically inferior to the other two. For brevity, we therefore do not consider this approach further.

3.3. An iterative approach

A potential problem with the transformation methods is that the initial estimated direction is poor. However, one approach to alleviate this is to apply an iterative scheme which starts with the initial estimate, obtains a new estimate after transformation and iteratively uses the new estimate as the initial estimate until convergence. A general algorithm for this approach is:

Step 0.1: Estimate the direction of β using PHD and denote this as $\hat{\mathbf{b}}^{(1)}$.

Step 0.2: Set $i = 1$ and `tol.met = FALSE`.

Step i : While `tol.met` is FALSE do

Step $i.1$: Apply **Method j** using $\hat{\mathbf{b}}^{(i)}$ as the direction for transformation and obtain a new estimated direction $\hat{\mathbf{b}}^{(i+1)}$.

Step $i.2$: If $1 - \text{cor}^2(\mathbf{X}\hat{\mathbf{b}}^{(i)}, \mathbf{X}\hat{\mathbf{b}}^{(i+1)}) < \text{tol}$ then set `tol.met` to TRUE.

Step $i.3$: Increment $i = i + 1$.

Step $i + 1$: Return $\hat{\mathbf{b}}^{(i)}$ as the final estimate to the direction of β .

We have used $\text{cor}^2(\mathbf{X}\hat{\mathbf{b}}^{(i)}, \mathbf{X}\hat{\mathbf{b}}^{(i+1)})$ in the criterion for exiting the iterating loop since, when there is correlation present amongst the columns of \mathbf{X} , notably different directions can result in very similar values of the ESSP, which is the targeted estimate. However, substantial differences in the squared correlation will similarly be accompanied by changes in the ESSP. Such an assessment is often used; for example, [Li \(1991\)](#) used the squared trace correlation, which is a multi-index version, to compare collections of estimated directions in dimension reduction.

4. Simulations

In this section we consider the performance of the transformation approaches defined earlier as Methods 1 and 2. Comparisons are also made with standard OLS and PHD estimation before we consider other methods in the next section.

***** Table 1 about here *****

In [Table 1](#) we provide simulated average squared correlations (with standard deviations in italics) for [Model 2](#) over 10,000 runs between $\mathbf{X}\beta$ and $\mathbf{X}\hat{\mathbf{b}}$, the true and estimated dimension reduced predictors. The estimators considered are OLS, PHD and the transformation approaches. Both $p = 10$ and $p = 20$ were considered. As expected OLS performs poorly due the symmetric dependency evident in the model, while PHD performs well. However, Methods 1 and 2 perform exceptionally well, having successfully drawn on the good PHD estimates to remove the symmetric dependency. The iterative estimation methods also provide improved estimates. We exited the iterative procedure when $\text{cor}^2(\mathbf{X}\hat{\mathbf{b}}^{(i)}, \mathbf{X}\hat{\mathbf{b}}^{(i+1)}) \geq 0.999$ or when ten iterations were reached. The small standard deviations for the methods indicate consistently excellent results.

[Prendergast & Sheather \(2013\)](#) found that robust least squares regression methods can provide improved single-index model estimates even when the data are well-behaved in the sense that they were sampled from a single index model and with a normal \mathbf{x} . Similarly, [Prendergast \(2008\)](#) found that trimming observations in the estimation step can also improve outcomes. We further explore this by considering the following model:

Model 3. $Y = 1/|\beta^\top \mathbf{x}| + 0.2\varepsilon$ with $p = 20$ and $\beta = (1, -2, 0, \dots, 0)^\top$.

For this model we will assume that \mathbf{x} is a random vector of dimension 10, distributed as $N(\mathbf{0}, \mathbf{I})$ so that this model also includes the symmetric dependency that troubles OLS. Moreover data simulated from this model can result in exceptionally extreme responses, since the denominator in the first term on the right hand side can be very close to zero.

***** Table 2 about here *****

In Table 2 we report the average squared correlations between the true and estimated dimension reduced predictors for 10,000 simulated runs from Model 3, with standard deviations in italics. For PHD, the very large response values often generated can result in extremely poor results. However, for OLS Garnham & Prendergast (2013) showed that using the rank of the response instead of the response itself could provide improved results. Consequently, we used the rank of the response values for PHD and this approach provides vast improvements. Therefore the PHD results presented in this table are based on this estimation approach. As well as employing OLS, PHD based on ranks and Methods 1 and 2, we also consider other variations that can be used to limit the influence of very large response values. RR, RM1 and RM2 refer to the usual OLS, Methods 1 and 2 but where OLS has been replaced with the M -estimation robust version (Huber 1964, 1973) with the Huber weight function. To do this we used the `r1m` function from the MASS package (Venables & Ripley 2002) in R (R Core Team 2013). M1-trim and M2-trim refer to Methods 1 and 2 where 10% of observations with the largest Cook's distance have been trimmed prior to the least squares step. This is one of the trimming procedures from Prendergast (2008). The iterative estimation scheme for this model and methods did not provide improved results so for simplicity the results are not presented here. Not surprisingly, OLS and the M -estimation equivalent completely fail even for large n due to symmetric dependency. Methods 1 and 2 perform much better but can still struggle as is made evident by the moderate average squared correlations and large standard deviations. On the other hand PHD performs well, in particular for the larger sample size settings. For Methods 1 and 2 coupled with M -estimation, we see improved performance over PHD for both methods. These results suggest that, by using the good PHD results in the transformation step to remove the symmetric dependency problem and then M -estimation to protect against large response values, excellent results can be achieved. For the trimming approaches, improvements have been found in comparison with standard Methods 1 and 2. However, the results are a little worse than PHD and much worse than the transformation plus M -estimation methods.

5. Inverse regression methods and multiple direction OLS

5.1. Inverse regression approaches

The transformations discussed in Section 3.2 are certainly not limited to OLS and PHD. Here we briefly discuss the use of sliced inverse regression (SIR, Li 1991) and sliced average variance estimates (SAVE, Cook & Weisberg 1991). For brevity, we only briefly discuss these methods here and the reader is directed to the aforementioned articles for more detail. Let S_1, \dots, S_H denote H non-overlapping yet collectively exhaustive intervals covering the range of Y . Let $\boldsymbol{\mu}_h = E(\boldsymbol{x}|Y \in S_h)$ ($h = 1, \dots, H$) denote slice means and consider the matrix

$$\mathbf{V} = \boldsymbol{\Sigma}^{-1/2} \sum_{h=1}^H p_h (\boldsymbol{\mu}_h - \boldsymbol{\mu})(\boldsymbol{\mu}_h - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1/2}$$

where $p_h = \Pr(Y \in S_h)$. Let $\boldsymbol{\gamma}$ be an eigenvector of \mathbf{V} corresponding to a non-zero eigenvalue. Li (1991) showed that $\boldsymbol{\Sigma}^{-1/2}\boldsymbol{\gamma}$ is an element of the span of the β_k s from the model in (2), provided that a K -direction version of Condition 1 holds. Consequently, if \mathbf{V} is of rank K , then SIR can recover a complete basis for the dimension reduction directions. However, SIR suffers from the same problems with symmetric dependency as OLS (Li 1991; Cook & Weisberg 1991) and is therefore a candidate for the same type of transformation.

Cook & Weisberg (1991) introduced SAVE, which does not suffer from symmetric dependency issues. It does, however, require that $\text{var}(\boldsymbol{x}|\beta_1^\top \boldsymbol{x}, \dots, \beta_K^\top \boldsymbol{x})$ is constant in addition to Condition 1. Both conditions are satisfied when \boldsymbol{x} is normally distributed, although both will often approximately hold in practice (Hall & Li 1993; Leeb 2013). Let $\boldsymbol{\Sigma}_h = \text{var}(\boldsymbol{x}|Y \in S_h)$ ($h = 1, \dots, H$) denote slice covariance matrices. Then SAVE is carried out similarly to SIR but where $\mathbf{M} = \sum_{h=1}^H p_h (\mathbf{I} - \boldsymbol{\Sigma}^{-1/2}\boldsymbol{\Sigma}_h\boldsymbol{\Sigma}^{-1/2})^2$ is used instead of \mathbf{V} . For some models SAVE requires large sample sizes to achieve good results. However, we have found that a variation of SAVE that was proposed by Zhu, Zhu & Feng (2010) called Cumulative Variance Estimation (CUVE) often provides excellent results. For $\boldsymbol{z} = \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{x} - \boldsymbol{\mu})$, CUVE estimates $E((P(Y \leq \tilde{Y})\mathbf{I} - \text{var}(\boldsymbol{z}I(Y \leq \tilde{Y})))^2)$ where \tilde{Y} is an independent copy of Y and $I(\cdot)$ is the indicator function taking the value 1 if its argument is true or zero otherwise. Similarly, Zhu, Zhu & Feng (2010) provided Cumulative Mean Estimation (CUME) which is a variation of SIR based on $E(E(\boldsymbol{z}I(Y \leq \tilde{Y}))E(\boldsymbol{z}I(Y \leq \tilde{Y}))^\top)$. Neither method requires choosing H , and Shaker & Prendergast (2011) showed that they can be successfully combined to obtain excellent results. Consequently, we will also consider CUME, following a transformation based on the first CUVE direction.

For SIR and SAVE the user must specify the subranges of Y used for ‘slicing’. One simple approach is to set H equally-probable slices. In practice this is equivalent to ordering the data by the magnitude of the response and allocating an (approximately) equal number of

observations to each slice. For estimation we use the `R` `dr` package (Weisberg 2002) which, by default, uses $\max(8, p + 3)$ for H . However, another option would be, for example, to seek an optimal choice for H (and simultaneously a value for K) using a bootstrap approach as given by Liquet & Saracco (2012).

***** Table 3 about here *****

An advantage of SIR, SAVE, CUME and CUVE is that the y_i s are used only to allocate x_i s. Consequently, we would not expect extremely large y_i s to have the same detrimental effect on estimation as they do for OLS. To highlight this we reconsider Model 3 and adapt Method 2 as follows. SAVE is used to get the initial estimate to the direction of β . SIR is then used on either the transformed y_i s (Method 1) or x_i s (Method 2), where the SAVE direction has been used to facilitate the transformations. Similarly, we use CUME following a transformation using the CUVE estimated direction. We also consider iterative estimation procedures for both. In Table 3 we provide the results from 10,000 simulations where $p = 20$ and $n = 100$ or 200 . Due to symmetry, both SIR and CUME fail to estimate the direction of β . SAVE also has trouble estimating the direction of β , especially for $n = 100$, although improvements are found for $n = 200$ and we observed good results for $n = 500$ (not shown). In contrast CUVE performs well, even for $n = 100$. The results also indicate that transformation Method 2 results in improved estimation, although the combination of SAVE and SIR is only moderately successful for the larger sample size. The combination of CUVE and CUME, however, provides excellent results even for $n = 100$. For both approaches, the iterative estimation scheme also provides improvements, as evidenced by the increased mean squared correlations.

5.2. Detecting multiple directions with OLS

We show here that a simple two-step estimation procedure for OLS can work exceptionally well when OLS is faced with the task of finding two directions, one of which is expected to be non-detectable due to symmetric dependency. For comparison we also consider PHD|OLS which is an iterative version of PHD and OLS considered by Shaker & Prendergast (2011). Here the first direction estimated is the OLS slope. Then PHD is used, conditional on this OLS slope estimate already detected, so that only new information is found in the second direction. The model that we focus on, given below, was also considered by Shaker & Prendergast (2011).

Model 4. $Y = \sin(0.5\beta_1^\top \mathbf{x}) + \cos(0.5\beta_2^\top \mathbf{x}) + 0.3\epsilon$ where $\beta_1 = (1, 2, -3, 0, \dots, 0)^\top$ and $\beta_2 = (1, 1, 0, -2, 0, \dots, 0)^\top$.

For the model above we will consider the performance of SIR, PHD, PHD|OLS and three new approaches based on Methods 1 and 2. For these new approaches we will use the OLS slope vector as the first estimated direction and then use the transformation methods to estimate a second direction.

***** Table 4 about here *****

In Table 4 we provide the simulated average first and second canonical correlations between $\mathbf{X}(\beta_1, \beta_2)$ and $\mathbf{X}\hat{\mathbf{B}}$ where $\hat{\mathbf{B}}$ is a $p \times 2$ matrix consisting of the first and second estimated directions. A large average first canonical correlation, \bar{r}_1 , indicates that the approach successfully detects the first direction. Similarly, a large \bar{r}_2 is indicative of good performance in detecting the second direction. SIR and PHD are both capable of finding one of the directions: for SIR the direction it is expected to find is β_1 and for PHD it is β_2 . However, these methods do not perform well at finding the other direction. PHD|OLS is expected to find both and the average canonical correlations indicate this, although the method has some trouble in estimating the second direction for $n = 100$. The transformation methods provide improvements in estimating both directions, certainly in the case of SIR and PHD and give marginally better results than even PHD|OLS.

We could similarly use robust M -estimation regression methods here too. Rather than repeat the simulation for similar results, we choose this approach for the example considered in Section 6.

5.3. Estimating K

So far, we have not discussed how to choose, or estimate, a suitable K . One approach is exploratory, in which ESSPs can be viewed to determine whether associated estimated vectors of predictor coefficients provide any information regarding the response. For this to work properly, the true K should be one or two, so that two- or three-dimensional ESSPs can be insightful. However it may be preferable to apply tests concerning the value of K so as to avoid confusion between what may be simply artifacts the data, as opposed to valid inference about K . This approach is also more suitable when $K > 2$, since a series of ESSPs may not be adequately informative. With respect to the methods we have employed throughout and which allow for the estimation of multiple directions (e.g., SIR, PHD, CUVE etc), there are several choices for the testing of K based on the method's respective eigenvalues. However, as we will show here, tests for K need to be treated with caution.

For SIR, SAVE and PHD we employ tests for a candidate choice of K , denoted K^* . These tests are based on asymptotic chi-squared distributions for sums of estimates (squared estimates in the case of PHD) of p hypothesised K^* zero-valued eigenvalues (Li 1991, 1992).

For CUME and CUVE we use the approach due to [Zhu, Zhu & Feng \(2010\)](#) which opts for the K^* that maximises a criterion function. If $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p \geq 0$ is the estimated eigenvalues for either CUME or CUVE then the criterion function is $n \sum_{k=1}^{K^*} \hat{\lambda}_k^2 / \sum_{k=1}^p \hat{\lambda}_k^2 - c_n K^* (K^* + 1)/2$ where the second term provides a penalty for larger K^* and where we follow the advice of [Zhu, Zhu & Feng \(2010\)](#) and choose $c_n = 2n^{3/4}/p$. It would also be possible to employ this approach for SIR, SAVE and PHD but for simplicity we choose not to do so. However, it should be noted that the minimum choice for K using the criterion is one, whereas zero can be chosen via the testing approach. In all cases we let \hat{K} denote the chosen K based on either the tests or the criterion.

***** Table 5 about here *****

When faced with symmetric dependency, care needs to be taken when estimating K . This is highlighted in [Table 5](#) where we simulated data 1000 times from each of Models 2, 3 and 4 and estimated K using the approaches described above and for two different sample sizes choices of $n = 200$ and $n = 500$. For each of these models the rank of the SIR matrix is equal to $K - 1$ and, consequently, K is typically underestimated. The test for PHD performs well for Model 2, but poorly for Model 3 where the dimension is often underestimated. This is likely due to a tendency for some extremely large response values to be generated from this model. These values are harmful to PHD but not to the slicing approaches. PHD also often underestimates K for Model 4, which is expected since the rank of the Hessian matrix can be shown to equal to one for this model. Overall, CUVE appears to perform very well for Model 2 both when $n = 200$ and 500 and for Model 3 when $n = 500$. For Model 2, the rank of the CUVE matrix is one so that K is expected to be underestimated. However, for the smaller sample size of $n = 200$ the correct $K = 2$ was chosen close to half of the time, presumably due to the zero-valued eigenvalues being estimated comparatively poorly (compared to when $n = 500$) and the penalty not being large enough.

As can be seen above, estimation of K needs to be treated with some caution. One approach would be to compute several estimates of K and to explore any discrepancies. Additionally, it may be possible that estimates of K are identical using two different methods, yet the associated directions are different. For example, for Model 4 both SIR and PHD together estimated K to be just one (underestimates of the desired $K = 2$) in 52.3% of the trials. However, among these 52.3% of trials, the average squared correlations between the dimension reduced predictors (i.e. between the $\hat{\mathbf{b}}^\top \mathbf{x}_i$ s for each of the methods) was just 0.057, indicating that different information regarding the β_k s was typically found from the two methods. It is not difficult to compute the correlations from different methods with a single data set, and this therefore provides a logical way in which to proceed. However in future

work we may investigate some form of meta-estimates of K that take into consideration the directions estimated using different methods.

6. The Ozone data example

Li (1992) considered the Ozone data from Breiman & Friedman (1985) which consists of 330 observations and eight predictors (e.g. wind speed, humidity etc.; refer to Table 4 of Li 1992 for the full list of predictors). The response is atmospheric ozone concentration. Li (1992) noted that SIR found a quadratic relationship (although not one that includes symmetric dependency) between the response and eight predictors and that an almost identical relationship could be found using least squares. Using PHD, another direction was found that eluded SIR, which provided an ESSP that exhibited symmetric dependency. Conflicting estimates of K were also obtained based on the SIR and PHD eigenvalues, where SIR chose $K = 1$ and PHD chose $K = 2$.

We now consider the Ozone data example further. Let Y denote the response variable and x denote the eight-dimensional vector of predictor variables. As previously we let the sample data be denoted by $\{y_i, x_i\}_{i=1}^{330}$. We base our model on \sqrt{Y} which, as we will see shortly, allows methods such as OLS, M -estimator regression methods and SIR to detect a linear relationship between the response and predictors. We also use M -estimation with the Huber weight function as a robust least-squares method. For convenience we will refer to this method as RR.

***** Figure 3 about here *****

Let \hat{b}_1 be the estimated slope for the RR regression of the $\sqrt{y_i}$ s on the x_i s. A plot of the $\sqrt{y_i}$ s versus the $\hat{b}_1^\top x_i$ s in Plot A of Figure 3 shows a linear relationship between the response and the dimension reduced predictors (labelled "1st RR dr predictor" on the plot). Using least squares regression, a simple linear model using the $\hat{b}_1^\top x_i$ s explains 72.1% of the variation in the square root of ozone. We now use transformation Method 1 with the first PHD direction but with RR replacing OLS, and let \hat{b}_2 denote this new estimate. Plot B shows that RR has now found another direction exhibiting a quadratic (or higher) relationship which is further emphasised by a fitted spline curve (black line). A least squares analysis based on a quadratic polynomial fit of just the $\hat{b}_2^\top x_i$ s results in a highly significant model that explains 33.1% of the variation in the response.

We now use OLS to fit a model to the $\hat{b}_1^\top x_i$ s, the $\hat{b}_2^\top x_i$ s and the square of each of these (we did not include the product of the two for a full quadratic model since this made little contribution). The estimated model is

$$\hat{Y}^{1/2} = -99.33 + 1.15 \times (\hat{\mathbf{b}}_1^\top \mathbf{x}) - 7.72 \times (\hat{\mathbf{b}}_2^\top \mathbf{x}) + 0.17 \times (\hat{\mathbf{b}}_1^\top \mathbf{x})^2 - 0.14 \times (\hat{\mathbf{b}}_2^\top \mathbf{x})^2.$$

The above fitted model explains approximately 76.4% of the variation in the square root of the response, indicating a good fit, and all of the terms in the model were highly significant. In Plots C and D we provide the residuals versus fits plot for the fit and also the quantile-quantile plot to check to see whether one could assume something close to a normal error term for the underlying model. These plots are excellent. They evince no evidence that the assumption of normally distributed error terms with homogeneous variance does not hold, at least approximately. In summary, we have successfully used RR twice to find two directions that can be used to construct a simple model with simple error term properties.

7. Discussion

This paper shows that simple response and predictor transformations can be used to remove the problem of symmetric dependency that affects some dimension reduction methods. While we initially show that OLS and PHD can be successfully employed in tandem for improved estimates, our approaches need not be limited to these methods. To highlight this, we also show that the popular robust M -estimation methods can be used, as well as sliced inverse regression in conjunction with sliced average variance estimates and associated cumulative slicing approaches. These approaches are particularly useful when the data include very large response values that can be detrimental to OLS estimation. Another interesting discovery presented in this paper is the ability of OLS, and robust equivalents, to find more than one direction.

Further research is envisaged. For example, we show that estimating K can be problematic in the sense that some directions may be “hidden” to some methods, resulting in an underestimation of K . However, different dimension reduction methods may find different directions and this possibility could be potentially be exploited to obtain a combined or meta estimate of K . Another possibility is to explore transformations for linear discriminant analysis (LDA). When Y is discrete, SIR is equivalent to LDA (see, e.g., [Li 2000](#); [Cook & Yin 2001](#)). If g is the number of sub-populations (distinct values in the domain of Y) then LDA can find at most $g - 1$ directions. Consequently predictor transformations may be useful for recovering additional directions. Finally, an exploration of sensitivity to distributional assumptions could also be considered. It is possible that, although normality of the predictor may fail to hold, thus affecting PHD estimation, the estimated PHD direction may still be

adequate to the extent that OLS (which has less restrictive conditions) can benefit following transformation.

A. Proof of Theorem 1

Throughout let $E(\mathbf{x}) = \boldsymbol{\mu}$, $\text{var}(\mathbf{x}) = \boldsymbol{\Sigma}$, $\mathbf{B} = [\beta_1, \dots, \beta_K]$ and recall that $\mathbf{v} \in \mathcal{S}$. It can be shown (see, e.g., [Prendergast 2005](#)) that Condition 1 is equivalent to

$$E(\mathbf{x}|\mathbf{B}^\top \mathbf{x}) = \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{B} (\mathbf{B}^\top \boldsymbol{\Sigma} \mathbf{B})^{-1} \mathbf{B}^\top (E(\mathbf{x}|\mathbf{B}^\top \mathbf{x}) - \boldsymbol{\mu}). \quad (11)$$

Since $\mathbf{v} \in \mathcal{S}$ so that $E(\text{sign}(\mathbf{v}^\top \mathbf{x} - \mathbf{v}^\top \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})) = E\{\text{sign}(\mathbf{v}^\top \mathbf{x} - \mathbf{v}^\top \boldsymbol{\mu})E(\mathbf{x} - \boldsymbol{\mu}|\mathbf{B}^\top \mathbf{x})\}$, it follows from (11) that $\boldsymbol{\Sigma}^{-1}E(t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v}))$ is equal to

$$\boldsymbol{\Sigma}^{-1}E(\text{sign}(\mathbf{v}^\top \mathbf{x} - \mathbf{v}^\top \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})) = \mathbf{B} (\mathbf{B}^\top \boldsymbol{\Sigma} \mathbf{B})^{-1} \mathbf{B}^\top E(t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v})) \in \mathcal{S}. \quad (12)$$

Similarly, by noting that

$$E(\text{sign}(\mathbf{v}^\top \mathbf{x} - \mathbf{v}^\top \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})Y) = E(\text{sign}(\mathbf{v}^\top \mathbf{x} - \mathbf{v}^\top \boldsymbol{\mu})YE(\mathbf{x} - \boldsymbol{\mu}|\mathbf{B}^\top \mathbf{x}))$$

since, from (2), Y is a function of $\mathbf{B}^\top \mathbf{x}$ and ε where ε is independent of \mathbf{x} , we can also show that

$$\boldsymbol{\Sigma}^{-1} \text{cov}(\text{sign}(\mathbf{v}^\top \mathbf{x} - \mathbf{v}^\top \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu}), Y) \in \mathcal{S}. \quad (13)$$

This shows that (6) holds.

For simplicity in what follows, let $\mathbf{e}_t = E(\text{sign}(\mathbf{v}^\top \mathbf{x} - \mathbf{v}^\top \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu}))$. Now, using (12),

$$\text{var}(t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v})) = \boldsymbol{\Sigma} - \mathbf{e}_t \mathbf{e}_t^\top$$

since $E(t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v})t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v})^\top) = E((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top) = \boldsymbol{\Sigma}$. Therefore (e.g., use the Small Rank Adjustment Lemma, [Horn & Johnson 1985](#), page 19)

$$(\text{var}(t_x(\mathbf{x} - \boldsymbol{\mu}; \mathbf{v})))^{-1} = \boldsymbol{\Sigma}^{-1} + \frac{1}{1 - \mathbf{e}_t^\top \boldsymbol{\Sigma}^{-1} \mathbf{e}_t} \boldsymbol{\Sigma}^{-1} \mathbf{e}_t \mathbf{e}_t^\top \boldsymbol{\Sigma}^{-1}$$

In conjunction with (13), this shows that (7) holds, completing the proof since $\boldsymbol{\Sigma}^{-1} \mathbf{e}_t \in \mathcal{S}$.

References

- BREIMAN, L. & FRIEDMAN, J.H. (1985). Estimating optimal transformations for multiple regression and correlation. *J. Amer. Statist. Assoc.* **80**, 580–619. With discussion and with a reply by the authors.
- BRILLINGER, D.R. (1977). The identification of a particular nonlinear time series system. *Biometrika* **64**, 509–515.
- BRILLINGER, D.R. (1983). A generalized linear model with “Gaussian” regressor variables. In *A Festschrift for Erich L. Lehmann*. Wadsworth Statist./Probab. Ser., Belmont, CA: Wadsworth, pp. 97–114.
- COOK, R.D. (1998a). Principal Hessian directions revisited. *J. Amer. Statist. Assoc.* **93**, 84–100. With comments by Ker-Chau Li and a rejoinder by the author.
- COOK, R.D. (1998b). *Regression graphics*. New York: John Wiley & Sons Inc. Ideas for studying regressions through graphics.
- COOK, R.D. & WEISBERG, S. (1991). Discussion of “Sliced inverse regression for dimension reduction”. *J. Amer. Statist. Assoc.* **86**, 328–332.
- COOK, R.D. & YIN, X. (2001). Dimension reduction and visualization in discriminant analysis (with discussion). *Aust. NZ. J. Stat.* **43**, 147–199.
- FARAWAY, J.J. (2006). *Extending the linear model with R: Generalized linear, mixed effects and nonparametric regression models*. Boca Raton, FL: Chapman & Hall/CRC press.
- FOX, J. & WEISBERG, S. (2011). *An R Companion to Applied Regression*. Thousand Oaks, CA: SAGE Publications.
- GARNHAM, A.L. & PRENDERGAST, L.A. (2013). A note on least squares sensitivity in single-index model estimation and the benefits of response transformations. *Electron. J. Stat.* **7**, 1983–2004.
- HALL, P. & LI, K.C. (1993). On almost linearity of low-dimensional projections from high-dimensional data. *Ann. Statist.* **21**, 867–889.
- HORN, R.A. & JOHNSON, C.A. (1985). *Matrix Analysis*. New York: Cambridge University Press.
- HUBER, P.J. (1964). Robust estimation of a location parameter. *Ann. Math. Stat.* **35**, 73–101.
- HUBER, P.J. (1973). Robust regression: asymptotics, conjectures and Monte Carlo. *Ann. Statist.* **1**, 799–821.
- LEEB, H. (2013). On the conditional distributions of low-dimensional projections from high-dimensional data. *Ann. Stat.* **41**, 464–483.
- LI, K.C. (1991). Sliced inverse regression for dimension reduction. *J. Amer. Statist. Assoc.* **86**, 316–342. With discussion and a rejoinder by the author.
- LI, K.C. (1992). On principal Hessian directions for data visualization and dimension reduction: another application of Stein’s lemma. *J. Amer. Statist. Assoc.* **87**, 1025–1039.
- LI, K.C. (2000). High dimensional data analysis via sir/phd approach. <http://www.stat.ucla.edu/~keli/sir-PHD.pdf>.
- LI, K.C. & DUAN, N. (1989). Regression analysis under link violation. *Ann. Statist.* **17**, 1009–1052.
- LIQUET, B. & SARACCO, J. (2012). A graphical tool for selecting the number of slices and the dimension of the model in SIR and SAVE approaches. *Comput. Statist.* **27**, 103–125.
- LUE, H.H. (2001). A study of sensitivity analysis on the method of principal hessian directions. *Comput. Stat.* **16**, 109–130.
- PRENDERGAST, L.A. (2005). Influence functions for sliced inverse regression. *Scand. J. Statist.* **32**, 385–404.
- PRENDERGAST, L.A. (2008). Trimming influential observations for improved single-index model estimated sufficient summary plots. *Comput. Statist. Data Anal.* **52**, 5319–5327.
- PRENDERGAST, L.A. & SHEATHER, S. (2013). On sensitivity of inverse response plot estimation and the benefits of a robust estimation approach. *Scand. J. Stat.* **40**, 219–237.
- PRENDERGAST, L.A. & SMITH, J.A. (2010). Influence functions for dimension reduction methods: an example influence study of principal Hessian direction analysis. *Scand. J. Stat.* **37**, 588–611.
- R CORE TEAM (2013). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org>.

- SHAKER, A.J. & PRENDERGAST, L.A. (2011). Iterative application of dimension reduction methods. *Electron. J. Stat.* **5**, 1471–1494.
- VENABLES, W.N. & RIPLEY, B.D. (2002). *Modern Applied Statistics with S*. New York: Springer, 4th edn. URL <http://www.stats.ox.ac.uk/pub/MASS4>. ISBN 0-387-95457-0.
- WEISBERG, S. (2002). Dimension reduction regression in R. *J. Stat. Softw.* **7**, 1–22.
- ZHU, L.P., ZHU, L.X. & FENG, Z.H. (2010). Dimension reduction in regressions through cumulative slicing estimation. *J. Amer. Statist. Assoc.* **105**, 1455–1466.

Author Manuscript

Figure 1. Plots of the y_i s versus the $\beta^\top x_i$ s (True Views) and the y_i s versus the $\hat{b}^\top x_i$ s (Estimated Views - ESSPs) for 100 observations generated for Model 1 (Plots A and B) and Model 2 (Plots C and D). OLS was used to estimate the direction of β .

Figure 2. Under the assumption of zero error in Model 2 and choosing $v = \beta$, Plot A provides the plot of the transformed Y versus $\beta^\top x$ and Plot B provides the plot of Y versus the transformed $\beta^\top t_x(x - \mu; v)_s$.

Figure 3. Plots of (A) the ESSP found by RR, (B) an ESSP created using a second direction found using RR following transformation Method 1, (C) residuals versus fitted values from a least squares fit to the dimension reduced predictors from Plots (A) and (B) and the square of these dimension reduced predictors and (D) the corresponding normal quantile-quantile plot.

Author Manuscript

TABLE 1

Average $\text{cor}^2(\mathbf{X}\beta, \mathbf{X}\hat{\mathbf{b}})$ across 10,000 simulated runs for Model 2 with different choices of n and p , where $\hat{\mathbf{b}}$ is the estimate from one of five methods: OLS, PHD, Method 1 (M1) and Method 2 (M2). The designations M1-it and M2-it refer to the iterative estimation scheme for Methods 1 and 2. Standard deviations are in italics.

Method	$p = 10$				$p = 20$			
	$n = 50$	$n = 100$	$n = 200$	$n = 500$	$n = 50$	$n = 100$	$n = 200$	$n = 500$
OLS	0.006 <i>0.171</i>	0.150 <i>0.167</i>	0.147 <i>0.165</i>	0.150 <i>0.165</i>	0.084 <i>0.105</i>	0.081 <i>0.101</i>	0.08 <i>0.099</i>	0.078 <i>0.095</i>
PHD	0.716 <i>0.206</i>	0.897 <i>0.068</i>	0.958 <i>0.023</i>	0.985 <i>0.008</i>	0.258 <i>0.212</i>	0.686 <i>0.170</i>	0.893 <i>0.046</i>	0.965 <i>0.013</i>
M1	0.813 <i>0.234</i>	0.968 <i>0.038</i>	0.989 <i>0.007</i>	0.996 <i>0.002</i>	0.31 <i>0.271</i>	0.833 <i>0.177</i>	0.97 <i>0.017</i>	0.991 <i>0.003</i>
M1-it	0.871 <i>0.211</i>	0.977 <i>0.024</i>	0.990 <i>0.005</i>	0.996 <i>0.002</i>	0.373 <i>0.306</i>	0.905 <i>0.147</i>	0.978 <i>0.010</i>	0.992 <i>0.003</i>
M2	0.833 <i>0.217</i>	0.967 <i>0.041</i>	0.989 <i>0.007</i>	0.996 <i>0.002</i>	0.332 <i>0.279</i>	0.831 <i>0.178</i>	0.969 <i>0.019</i>	0.991 <i>0.004</i>
M2-it	0.875 <i>0.199</i>	0.977 <i>0.029</i>	0.990 <i>0.005</i>	0.996 <i>0.002</i>	0.372 <i>0.295</i>	0.887 <i>0.168</i>	0.978 <i>0.011</i>	0.992 <i>0.003</i>

TABLE 2

Average $\text{cor}^2(\mathbf{X}\beta, \mathbf{X}\hat{\mathbf{b}})$ across 10,000 simulated data sets for Model 3 with different choices of n , where $\hat{\mathbf{b}}$ is the estimate from various methods. The designations RR, RM1 and RM2 refer to Methods OLS, M1 and M2 but where robust regression M -estimation has been used in the regression step with the Huber weight function. The designations M1-trim and M2-trim refer to Methods M1 and M2 but where 10% of observations with the largest Cook's distance were trimmed. Standard deviations are in italics.

Method	$n = 100$	$n = 200$	$n = 500$	Method	$n = 100$	$n = 200$	$n = 500$
OLS	0.007 <i>0.013</i>	0.004 <i>0.009</i>	0.002 <i>0.005</i>	M3	0.164 <i>0.124</i>	0.235 <i>0.137</i>	0.268 <i>0.146</i>
RR	0.034 <i>0.051</i>	0.035 <i>0.050</i>	0.035 <i>0.047</i>	RM1	0.718 <i>0.257</i>	0.955 <i>0.063</i>	0.987 <i>0.018</i>
PHD	0.662 <i>0.229</i>	0.920 <i>0.035</i>	0.978 <i>0.008</i>	RM2	0.752 <i>0.244</i>	0.961 <i>0.015</i>	0.982 <i>0.006</i>
M1	0.271 <i>0.243</i>	0.518 <i>0.311</i>	0.705 <i>0.317</i>	M1-trim	0.655 <i>0.264</i>	0.924 <i>0.106</i>	0.970 <i>0.044</i>
M2	0.437 <i>0.226</i>	0.582 <i>0.212</i>	0.662 <i>0.197</i>	M2-trim	0.613 <i>0.241</i>	0.875 <i>0.088</i>	0.927 <i>0.049</i>

TABLE 3

Average $cor^2(\mathbf{X}\beta, \mathbf{X}\hat{\mathbf{b}})$ across 10,000 simulated data sets for Model 3 with two different choices of $n, p = 20$, where $\hat{\mathbf{b}}$ is the estimate from various methods. The designation M2 refers to transformation Method 2 and M2-it refers to this transformation with iterative estimation. Standard deviations are in italics.

Method	$n = 100$	$n = 200$	Method	$n = 100$	$n = 200$
SIR	0.057 <i>0.086</i>	0.059 <i>0.088</i>	SAVE SIR (M2)	0.121 <i>0.175</i>	0.494 <i>0.365</i>
CUME	0.058 <i>0.078</i>	0.059 <i>0.078</i>	SAVE SIR (M2-it)	0.162 <i>0.243</i>	0.686 <i>0.429</i>
SAVE	0.110 <i>0.134</i>	0.327 <i>0.241</i>	CUVE CUME (M2)	0.905 <i>0.124</i>	0.983 <i>0.007</i>
CUVE	0.808 <i>0.130</i>	0.937 <i>0.026</i>	CUVE CUME (M2-it)	0.940 <i>0.114</i>	0.988 <i>0.005</i>

TABLE 4

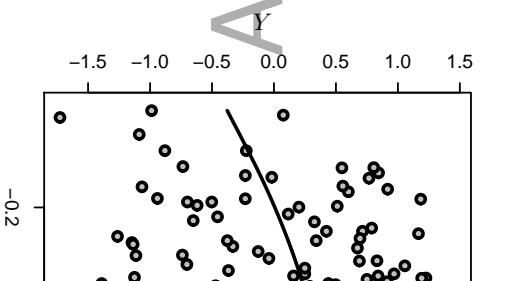
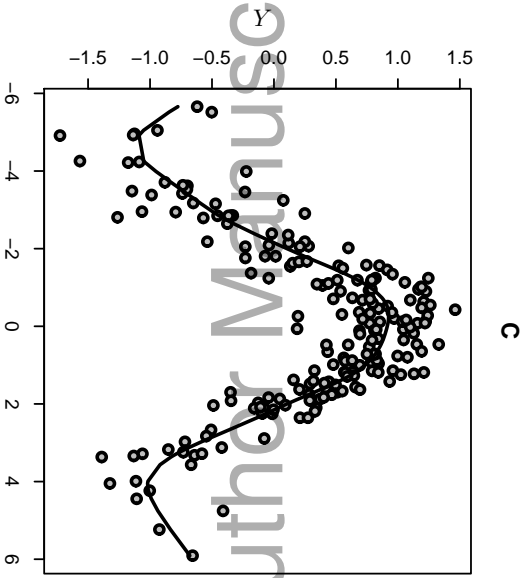
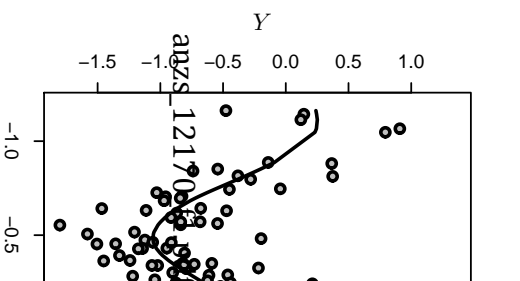
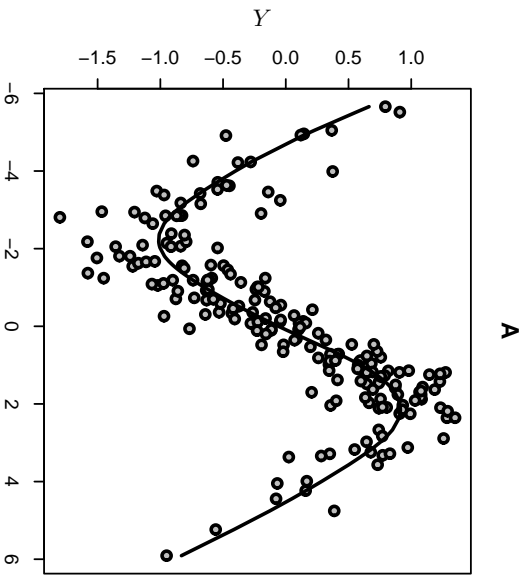
Average first and second canonical correlations (\bar{r}_1, \bar{r}_2) between $\mathbf{X}[\beta_1, \beta_2]$ and $\mathbf{X}\hat{\mathbf{B}}$ across 10,000 simulated runs for data generated from Model 4 with different choices of n and p where $\hat{\mathbf{b}}$ is the estimate from one of five methods: OLS, PHD, Method 1 (M1) and Method 2 (M2). Standard deviations are in italics.

Method	$n = 100$		$n = 200$		$n = 500$		$n = 1000$	
	\bar{r}_1	\bar{r}_2	\bar{r}_1	\bar{r}_2	\bar{r}_1	\bar{r}_2	\bar{r}_1	\bar{r}_2
SIR	0.743 <i>0.154</i>	0.264 <i>0.19</i>	0.844 <i>0.103</i>	0.296 <i>0.212</i>	0.939 <i>0.037</i>	0.393 <i>0.254</i>	0.971 <i>0.016</i>	0.53 <i>0.271</i>
PHD	0.912 <i>0.059</i>	0.386 <i>0.226</i>	0.961 <i>0.024</i>	0.403 <i>0.232</i>	0.985 <i>0.008</i>	0.407 <i>0.235</i>	0.993 <i>0.004</i>	0.409 <i>0.233</i>
PHD OLS	0.914 <i>0.059</i>	0.632 <i>0.219</i>	0.961 <i>0.023</i>	0.805 <i>0.125</i>	0.985 <i>0.008</i>	0.918 <i>0.045</i>	0.993 <i>0.004</i>	0.958 <i>0.022</i>
OLS,M1	0.927 <i>0.056</i>	0.671 <i>0.192</i>	0.972 <i>0.017</i>	0.827 <i>0.098</i>	0.99 <i>0.006</i>	0.923 <i>0.041</i>	0.995 <i>0.003</i>	0.96 <i>0.021</i>
OLS,M2	0.933 <i>0.054</i>	0.671 <i>0.194</i>	0.974 <i>0.016</i>	0.827 <i>0.099</i>	0.991 <i>0.005</i>	0.923 <i>0.042</i>	0.996 <i>0.002</i>	0.96 <i>0.021</i>

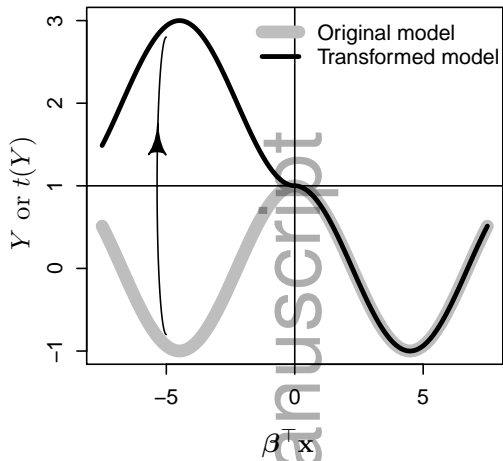
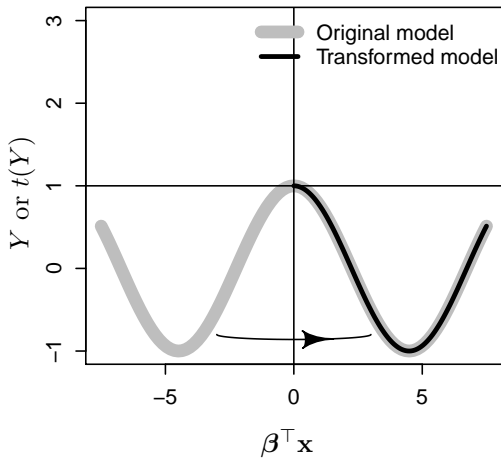
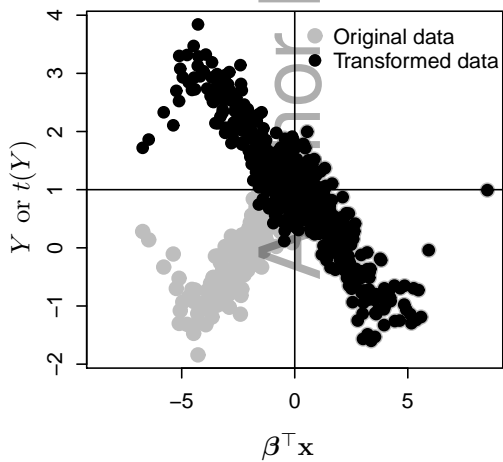
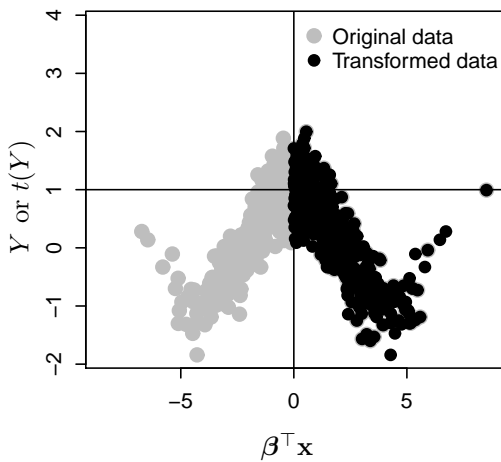
TABLE 5

Proportion of times that $\hat{K} = 0, 1, 2$ or ≥ 3 when using the estimated eigenvalues from the methods SIR, PHD, CUME and CUVE for Models 2-4 and with $n = 200$ and 500 for 1000 simulation runs.

Method	$\hat{K} = 0$	$n = 200$				$n = 500$			
		$\hat{K} = 1$	$\hat{K} = 2$	$\hat{K} \geq 3$	$\hat{K} = 0$	$\hat{K} = 1$	$\hat{K} = 2$	$\hat{K} \geq 3$	
Model 2 (K = 1) SIR	0.950	0.049	0.001	0.000	0.951	0.047	0.002	0.000	
PHD	0.000	0.973	0.026	0.001	0.000	0.952	0.042	0.006	
CUME	-	0.736	0.263	0.001	-	0.764	0.236	0.000	
CUVE	-	0.996	0.004	0.000	-	1.000	0.000	0.000	
Model 3 (K = 1) SIR	0.954	0.046	0.000	0.000	0.960	0.038	0.002	0.000	
PHD	0.785	0.215	0.000	0.000	0.755	0.245	0.000	0.000	
CUME	-	0.452	0.547	0.001	-	0.462	0.535	0.003	
CUVE	-	0.007	0.747	0.246	-	0.983	0.017	0.000	
Model 4 (K = 2) SIR	0.628	0.344	0.025	0.003	0.092	0.828	0.076	0.004	
PHD	0.011	0.894	0.091	0.004	0.000	0.867	0.122	0.011	
CUME	-	0.997	0.003	0.000	-	1.000	0.000	0.000	
CUVE	-	0.526	0.468	0.006	-	0.993	0.007	0.000	

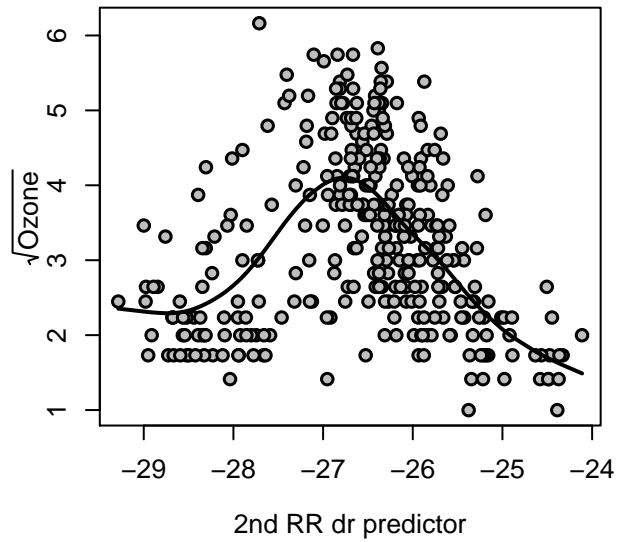
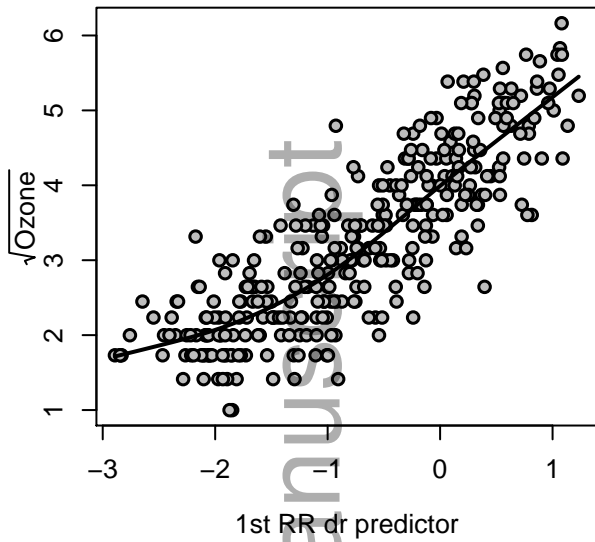


anzs_1217088

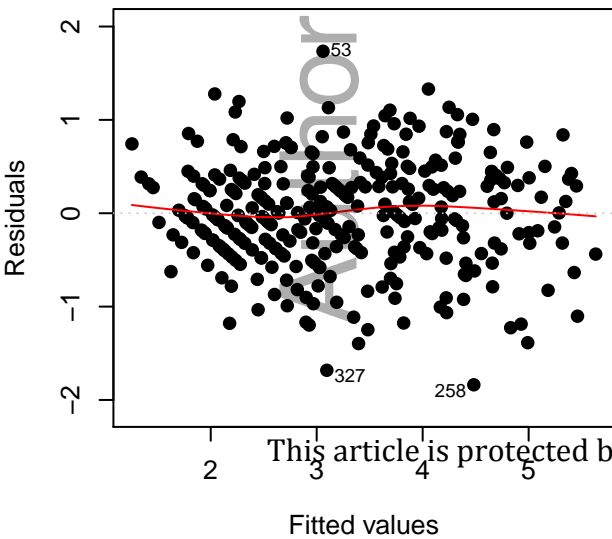
A**B****C****D**

A: RR view 1

B: RR view 2 using Method 1



C: Residuals versus fitted



D: Normal Q-Q plot

