

## Correlation pursuit: forward stepwise variable selection for index models

Wenxuan Zhong,

*University of Illinois at Urbana–Champaign, Champaign, USA*

Tingting Zhang,

*University of Virginia, Charlottesville, USA*

Yu Zhu

*Purdue University, West Lafayette, USA*

and Jun S. Liu

*Harvard University, Cambridge, USA*

[Received February 2010. Final revision November 2011]

**Summary.** A stepwise procedure, correlation pursuit (COP), is developed for variable selection under the sufficient dimension reduction framework, in which the response variable  $Y$  is influenced by the predictors  $X_1, X_2, \dots, X_p$  through an unknown function of a few linear combinations of them. Unlike linear stepwise regression, COP does not impose a special form of relationship (such as linear) between the response variable and the predictor variables. The COP procedure selects variables that attain the maximum correlation between the transformed response and the linear combination of the variables. Various asymptotic properties of the COP procedure are established and, in particular, its variable selection performance under a diverging number of predictors and sample size is investigated. The excellent empirical performance of the COP procedure in comparison with existing methods is demonstrated by both extensive simulation studies and a real example in functional genomics.

**Keywords:** Dimension reduction; Projection pursuit regression; Sliced inverse regression; Stepwise regression; Variable selection

### 1. Introduction

Advances in science and technology in the past few decades have led to an explosive growth of high dimensional data across a variety of areas such as genetics, molecular biology, cognitive sciences, environmental sciences, astrophysics, finance and Internet commerce. Compared with their dimensionalities, a large amount of data sets generated from these areas have relatively small sample sizes. Variable (or feature) selection and dimension reduction are more than often key steps in analysing these data. Much progress has been made in the past few decades on variable selection for linear models (see Shao (1998) and Fan and Lv (2010) for a review). In recent years, shrinkage-based procedures for simultaneously estimating regression coefficients and selecting predictors have been particularly attractive to researchers, and many promising

*Address for correspondence:* Wenxuan Zhong, Department of Statistics, University of Illinois at Urbana–Champaign, 725 South Wright Street, Champaign, IL 61820, USA.  
E-mail: wenxuan@illinois.edu

algorithms such as the lasso (Tibshirani, 1996; Zou, 2006; Friedman, 2007), LARS (Efron *et al.*, 2004) and smoothly clipped absolute deviation (SCAD) (Fan and Li, 2001) have been invented.

Let  $Y \in \mathbb{R}$  be a univariate response variable and  $X = (X_1, X_2, \dots, X_p)' \in \mathbb{R}^p$  a vector of  $p$  continuous predictor variables. Throughout this paper, we consider the following sufficient dimension reduction (SDR) model framework as pioneered by Li (1991) and Cook (1994). Let  $\beta_1, \beta_2, \dots, \beta_K$  be  $p$ -dimensional vectors with  $\beta_i = (\beta_{1i}, \beta_{2i}, \dots, \beta_{pi})'$  for  $1 \leq i \leq K$ . The SDR model assumes that  $Y$  and  $X$  are mutually independent conditional on  $\beta_1'X, \beta_2'X, \dots, \beta_K'X$ , i.e.

$$Y \perp X | B'X, \quad (1)$$

where ' $\perp$ ' means 'independent of' and  $B = (\beta_1, \beta_2, \dots, \beta_K)$ . Expression (1) implies that all the information  $X$  contains about  $Y$  is contained in the  $K$  projections  $\beta_1'X, \dots, \beta_K'X$ . A predictor variable  $X_j$  ( $1 \leq j \leq p$ ) is said to be relevant if there is at least one  $i$  ( $1 \leq i \leq K$ ) such that  $\beta_{ji} \neq 0$ . Let  $L$  be the number of relevant predictor variables. When there are a large number of predictors (i.e.  $p$  is large), it is usually safe to impose the *sparsity assumption*, which states that only a small subset of the predictors influences  $Y$  and the others are irrelevant. In the SDR model, this assumption means that both  $K$  and  $L$  are small relative to  $p$ .

In his seminal paper on dimension reduction, Li (1991) proposed a seemingly different model of the form

$$Y = f(\beta_1'X, \beta_2'X, \dots, \beta_K'X, \varepsilon), \quad (2)$$

where  $f$  is an unknown  $(K+1)$ -variate link function and  $\varepsilon$  is a stochastic error independent of  $X$ . It has been shown that the two models (1) and (2) are in fact equivalent (Zeng and Zhu, 2010). We henceforth always refer to  $\beta_1, \beta_2, \dots, \beta_K$  as the SDR directions and the space spanned by these directions as an SDR subspace. In general, SDR subspaces are not unique. To resolve this ambiguity, Cook (1994) introduced the concept of a *central subspace*, which is the intersection of all possible SDR subspaces and is an SDR subspace itself, and showed that the central space is well defined and unique under some general conditions. We denote the central subspace by  $S(B)$  and assume its existence throughout this paper.

Various methods have been developed for estimating  $\beta_1, \dots, \beta_K$  in the literature on SDR. One particular family of methods utilizes inverse regression, which is to regress  $X$  against  $Y$ . The sliced inversion regression (SIR) method that was proposed by Li (1991) is the forerunner of this family of methods. Recognizing that estimation of the SDR directions does not automatically lead to variable selection, Cook (2004) derived various  $\chi^2$ -tests for assessing the contribution of predictor variables to the SDR directions. On the basis of these tests, Li *et al.* (2005) proposed a backward subset selection method for selecting significant predictors. Following the recent trend of using the  $L_1$ - or  $L_2$ -penalty for variable selection, Zhong *et al.* (2005) proposed to regularize the sample covariance matrix of the predictor variables in SIR and developed a procedure called regularized SIR for variable selection. Li (2007) proposed sparse SIR (SSIR) to obtain shrinkage estimates of the SDR directions. Bondell and Li (2009) further adopted the non-negative garrotte method for estimating the SDR directions and showed that the resulting method is consistent in variable selection.

The majority of the aforementioned methods take a two-step approach to variable selection under the SDR model. The first step is to perform dimension reduction, i.e. to estimate the SDR directions; and the second step is to select the relevant variables by using statistical testing or shrinkage methods. Because these methods need to estimate the covariance and conditional covariance matrices of  $X$ , both of which are of dimensions  $p \times p$ , the effectiveness and robustness of the two-step approach are questionable when  $p$  is large relative to  $n$ . Zhu *et al.* (2006) have

shown that the accuracy of estimation of SDR directions deteriorates as  $p$  increases. In other words, the more irrelevant variables there are, the more likely a method fails to estimate the SDR directions accurately, and the less likely the method identifies the true relevant predictor variables.

In this paper, we propose correlation pursuit (COP), which is a stepwise procedure for simultaneous dimension reduction and variable selection under the SDR model. Similar to projection pursuit (Friedman and Tukey, 1974; Huber, 1985), COP defines a projection function to measure the correlation between the transformed response and the projections of  $X$  and pursues a subset of explanatory variables that maximize the projection function. It starts with a randomly selected subset and iterates between finding an explanatory variable (predictor) that significantly improves the current projection function to add to the subset and finding an insignificant predictor to remove from the subset. During each iteration step, COP needs only to consider the predictors that are currently in the subset and one more predictor outside the subset. Therefore, COP can avoid the estimation and inversion of  $p \times p$  covariance and conditional covariance matrices of  $X$  and mitigate the curse of dimensionality. Furthermore, COP performs dimension reduction and variable selection simultaneously. Therefore, dimension reduction and variable selection can be mutually enhanced. Our theoretical investigations as well as simulation studies show that COP is a promising tool for dimension reduction and variable selection in high dimensional data analysis.

The rest of the paper is organized as follows. In Section 2, we give a brief introduction to SIR, following a correlation interpretation of SIR that was provided by Chen and Li (1998). This interpretation was also used in Fung *et al.* (2002) and Zhou and He (2008) for dimension reduction via canonical correlation. In the same section, we describe the COP procedure and derive various test statistics that are used by the procedure. The asymptotic behaviour of the COP procedure is discussed in Section 3. Several implementation issues of the procedure are discussed in Section 4. Simulation and real data examples are reported in Sections 5 and 6 respectively. Additional remarks in Section 7 conclude the paper. An abbreviated version of the proofs of the theorems is provided in Appendix A.

## 2. Correlation pursuit for variable selection

### 2.1. Profile correlation and sliced inverse regression

Let  $\eta$  be an arbitrary direction in  $\mathbb{R}^p$ . We define the *profile correlation* between  $Y$  and  $\eta'X$ , which is denoted by  $P(\eta)$ , as

$$P(\eta) = \max_T [\text{corr}\{T(Y), \eta'X\}], \quad (3)$$

where the maximization is taken over all possible transformations of  $Y$  including non-monotone transformations. The profile correlation  $P(\eta)$  reflects the largest possible correlation between a transformed response  $T(Y)$  and the projection  $\eta'X$ . Let  $\eta_1$  be the direction that maximizes  $P(\eta)$  subject to  $\eta'\Sigma\eta = 1$ , i.e.  $\eta_1 = \arg \max_{\eta'\Sigma\eta=1} \{P(\eta)\}$ . We refer to  $\eta_1$  as the first *principal direction* for the profile correlation between  $Y$  and  $X$  and call  $P(\eta_1)$  the first profile correlation. Direction  $\eta_1$ , or its projection  $\eta_1'X$ , may not entirely characterize the dependence between  $Y$  and  $X$ . Using  $P(\eta)$  as the projection function again, we can look for a second direction, which is denoted by  $\eta_2$ , which is uncorrelated with  $\eta_1'X$ , i.e.  $\eta_2'\Sigma\eta_1 = 0$ , and maximizes  $P(\eta)$ , i.e.  $\eta_2 = \arg \max_{\eta_2'\Sigma\eta_1=0} \{P(\eta)\}$ . We refer to  $\eta_2$  as the second principal direction and  $P(\eta_2)$  as the second profile correlation. This procedure can be continued until no more directions can be found that are orthogonal to the directions obtained and have a non-zero profile correlation with  $Y$ . Suppose that  $\tilde{K}$  principal directions exist between  $Y$  and  $X$ , which are  $\eta_1, \eta_2, \dots$ ,

and  $\eta_{\tilde{K}}$ , with the corresponding profile correlations  $P(\eta_1) \geq P(\eta_2) \geq \dots P(\eta_{\tilde{K}}) > 0$ . We need to impose the following condition to establish the connection between the principal directions and the SDR directions under the SDR model.

*Condition 1* (linearity condition). For any  $\eta$  in  $\mathbb{R}^p$ ,  $E(\eta'X|B'X)$  is linear in  $B'X$ , where  $B$  is as defined in equation (1).

*Proposition 1.* Under the SDR model and the linearity condition, the principal directions  $\eta_1, \eta_2, \dots, \eta_{\tilde{K}}$  are in the central space  $\mathcal{S}(B)$ .

To make this paper self-sufficient, we have included the proof of proposition 1 in Appendix A. Based on the proposition, the principal directions are indeed SDR directions. In general,  $\tilde{K} < K$ . When the link function  $f$  is symmetric along a direction, using correlation alone may fail to recover this direction. For example, if  $Y = X_1^2 + \varepsilon$ , the profile correlation between  $Y$  and  $X_1$  will always be 0. To exclude this possibility, we follow the convention in the SDR literature to impose the following condition.

*Condition 2* (coverage condition). The number of principal directions of profile correlation is equal to the dimensionality of the central subspace, i.e.  $\tilde{K} = K$ .

Under both the linearity and the coverage conditions, the principal directions  $\eta_1, \eta_2, \dots, \eta_K$  form a special basis of the central subspace  $\mathcal{S}(B)$ , i.e.  $\mathcal{S}(B) = \text{span}(\eta_1, \eta_2, \dots, \eta_K)$ . This basis is uniquely defined and is the estimation target of SIR. In the rest of the paper, for ease of discussion, we use  $\beta_1, \beta_2, \dots, \beta_K$  and  $\eta_1, \eta_2, \dots, \eta_K$ , interchangeably.

Chen and Li (1998) showed that, at the population level, there is an explicit solution for the principal directions. In the proof of their theorem 3.1, Chen and Li (1998) derived that

$$P^2(\eta) = \frac{\eta' \text{var}\{E(X|Y)\}\eta}{\eta' \Sigma \eta} \equiv \frac{\eta' M \eta}{\eta' \Sigma \eta}, \tag{4}$$

where  $M = \Delta \text{var}\{E(X|Y)\}$  is the covariance matrix of the expectation of  $X$  given  $Y$ . Furthermore, the principal directions of profile correlation are the solutions of the eigenvalue decomposition problem

$$M v_i = \lambda_i \Sigma v_i, \quad v_i' \Sigma v_i = 1, \quad \text{for } i = 1, 2, \dots, K; \tag{5}$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_K > 0. \tag{6}$$

The principal directions  $\eta_1, \eta_2, \dots$ , and  $\eta_K$  are the first  $K$  eigenvectors of  $\Sigma^{-1}M$ , and their corresponding eigenvalues are exactly the squared profile correlations, i.e.  $P^2(\eta_i) = \lambda_i$  for  $i = 1, 2, \dots, K$ .

Given independent observations  $\{(\mathbf{x}_i, y_i)\}_{i=1, \dots, n}$  of  $(X, Y)$ , where  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})'$ ,  $\Sigma$  can be estimated by the sample covariance matrix,

$$\hat{\Sigma} = \frac{\sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})'}{n - 1}, \tag{7}$$

where  $\bar{\mathbf{x}}$  is the sample mean of  $\{\mathbf{x}_i\}$ . Li (1991) proposed the following SIR procedure to estimate  $M$ . First, the range of  $\{y_i\}_{i=1}^n$  is divided into  $H$  disjoint intervals, which are denoted as  $S_1, \dots, S_H$ . For  $h = 1, \dots, H$ , the mean vector  $\bar{\mathbf{x}}_h = n_h^{-1} \sum_{y_i \in S_h} \mathbf{x}_i$  is calculated, where  $n_h$  is the number of  $y_i$ s in  $S_h$ . Then,  $M$  is estimated by

$$\hat{M} = \frac{\sum_{h=1}^H n_h (\bar{\mathbf{x}}_h - \bar{\mathbf{x}})(\bar{\mathbf{x}}_h - \bar{\mathbf{x}})'}{n}, \tag{8}$$

and the matrix  $\Sigma^{-1}M$  is estimated by  $\hat{\Sigma}^{-1}\hat{M}$ . The first  $K$  eigenvectors of  $\hat{\Sigma}^{-1}\hat{M}$ , which are denoted by  $\hat{\eta}_1, \hat{\eta}_2, \dots, \hat{\eta}_K$ , are used to estimate the first  $K$  eigenvectors of  $\Sigma^{-1}M$  or, equivalently, the principal directions  $\eta_1, \eta_2, \dots, \eta_K$  respectively. The first  $K$  eigenvalues of  $\hat{\Sigma}^{-1}\hat{M}$ , which are denoted by  $\hat{\lambda}_1, \hat{\lambda}_2, \dots, \hat{\lambda}_K$ , are used to estimate the eigenvalues of  $\Sigma^{-1}M$  or, equivalently, the squared profile correlations  $\lambda_1, \lambda_2, \dots, \lambda_K$  respectively.

### 2.2. Correlation pursuit

The SIR method needs to estimate the two  $p \times p$  covariance matrices  $\Sigma$  and  $M$ , and to obtain the eigenvalue decomposition of  $\hat{\Sigma}^{-1}\hat{M}$ . When a large number of irrelevant variables are present and the sample size  $n$  is relatively small,  $\hat{\Sigma}$  and  $\hat{M}$  become unstable, which leads to very inaccurate estimates of principal directions  $\hat{\eta}_1, \hat{\eta}_2, \dots, \hat{\eta}_K$  (Zhu *et al.*, 2006). As a consequence, those shrinkage-based variable selection methods that rely on  $\hat{\eta}_1, \hat{\eta}_2, \dots, \hat{\eta}_K$  often perform poorly for the SDR model when  $p$  is large. We here propose a stepwise SIR-based procedure for simultaneous dimension reduction (i.e. estimating the principal directions) and variable selection (i.e. identifying true predictors). Our procedure starts with a collection of randomly selected predictors and iterates between an addition step, which selects and adds a predictor to the collection, and a deletion step, which selects and deletes a predictor from the collection. The procedure terminates when no new addition or deletion occurs.

#### 2.2.1. Addition step

Let  $\mathcal{A}$  denote the collection of the indices of the selected predictors and  $X_{\mathcal{A}}$  the collection of the selected variables. Applying SIR to the data involving only the predictors in  $X_{\mathcal{A}}$ , we obtain the estimated squared profile correlations  $\hat{\lambda}_1^{\mathcal{A}}, \hat{\lambda}_2^{\mathcal{A}}, \dots, \hat{\lambda}_K^{\mathcal{A}}$ . Superscript  $\mathcal{A}$  indicates that the estimated squared profile correlations depend on the current subset of selected predictors. Let  $X_t$  be an arbitrary predictor outside  $\mathcal{A}$  and  $\mathcal{A} + t = \mathcal{A} \cup \{t\}$ . Applying SIR to the data involving the predictors in  $\mathcal{A} + t$ , we obtain the estimated squared profile correlations  $\hat{\lambda}_1^{\mathcal{A}+t}, \hat{\lambda}_2^{\mathcal{A}+t}, \dots, \hat{\lambda}_K^{\mathcal{A}+t}$ . Because  $\mathcal{A} \subset \mathcal{A} + t$ , it is easy to see that  $\hat{\lambda}_1^{\mathcal{A}} \leq \hat{\lambda}_1^{\mathcal{A}+t}$ . The difference  $\hat{\lambda}_1^{\mathcal{A}+t} - \hat{\lambda}_1^{\mathcal{A}}$  reflects the amount of improvement in the first profile correlation due to the incorporation of  $X_t$ . We standardize this difference and use the resulting test statistic

$$\text{COP}_1^{\mathcal{A}+t} = \frac{n(\hat{\lambda}_1^{\mathcal{A}+t} - \hat{\lambda}_1^{\mathcal{A}})}{1 - \hat{\lambda}_1^{\mathcal{A}+t}}, \tag{9}$$

to assess the significance of adding  $X_t$  to  $\mathcal{A}$  in improving the first profile correlation. Similarly, the contributions of adding  $X_t$  to the other profile correlations can be assessed by

$$\text{COP}_i^{\mathcal{A}+t} = \frac{n(\hat{\lambda}_i^{\mathcal{A}+t} - \hat{\lambda}_i^{\mathcal{A}})}{1 - \hat{\lambda}_i^{\mathcal{A}+t}}, \tag{10}$$

for  $2 \leq i \leq K$ . The overall contribution of adding  $X_t$  to the improvement in all the  $K$  profile correlations can be assessed by combining the statistics  $\text{COP}_i^{\mathcal{A}+t}$  into one single test statistic

$$\text{COP}_{1:K}^{\mathcal{A}+t} = \sum_{i=1}^K \text{COP}_i^{\mathcal{A}+t}. \tag{11}$$

We further define that

$$\overline{\text{COP}}_{1:K}^A = \max_{t \in \mathcal{A}^c} (\text{COP}_{1:K}^{A+t}). \tag{12}$$

Let  $X_{\bar{t}}$  be a predictor that attains  $\overline{\text{COP}}_{1:K}^A$ , i.e.  $\overline{\text{COP}}_{1:K}^A = \text{COP}_{1:K}^{A+\bar{t}}$ , and let  $c_e$  be a prespecified threshold (details about its choice are deferred to the next two sections). Then, if  $\overline{\text{COP}}_{1:K}^A > c_e$ , we add  $\bar{t}$  to  $\mathcal{A}$ ; otherwise, we do not add any variable.

2.2.2. *Deletion step*

Let  $X_t$  be an arbitrary predictor in  $\mathcal{A}$  and define  $\mathcal{A} - t = \mathcal{A} - \{t\}$ . Let  $\hat{\lambda}_1^{A-t}, \hat{\lambda}_2^{A-t}, \dots, \hat{\lambda}_K^{A-t}$  be the estimated squared profile correlations based on the data involving the predictors in  $\mathcal{A} - t$  only. The effect of deleting  $X_t$  from  $\mathcal{A}$  on the  $i$ th squared profile correlation can be measured by

$$\text{COP}_i^{A-t} = \frac{n(\hat{\lambda}_i^A - \hat{\lambda}_i^{A-t})}{1 - \hat{\lambda}_i^A}, \tag{13}$$

for  $1 \leq i \leq K$ . The overall effect of deleting  $X_t$  is measured by

$$\text{COP}_{1:K}^{A-t} = \sum_{i=1}^K \text{COP}_i^{A-t}, \tag{14}$$

and the least effect from deleting one predictor from  $\mathcal{A}$  is then defined to be

$$\underline{\text{COP}}_{1:K}^A = \min_{t \in \mathcal{A}} (\text{COP}_{1:K}^{A-t}). \tag{15}$$

Let  $X_{\underline{t}}$  be a predictor that achieves  $\underline{\text{COP}}_{1:K}^A$ , and let  $c_d$  be a prespecified threshold for deletion. If  $\underline{\text{COP}}_{1:K}^A < c_d$ , we delete  $X_{\underline{t}}$  from  $\mathcal{A}$ ; otherwise, no deletion happens.

The asymptotic distributions of the proposed statistics and the selection of the thresholds will be discussed in the next two sections. Because the procedure described aims to find predictors that can most significantly improve the profile correlations between  $Y$  and  $X$ , we call it the COP procedure. Below we summarize the COP algorithm.

- Step 1:* set the number of principal directions  $K$  and the threshold values  $c_e$  and  $c_d$ .
- Step 2:* randomly select  $K + 1$  variables as the initial collection of selected variables  $\mathcal{A}$ .
- Step 3:* iterate until no more addition or deletion of predictors can be performed; in the addition step,
  - (a) find  $\bar{t}$  such that  $\text{COP}_{1:K}^{A+\bar{t}} = \overline{\text{COP}}_{1:K}^A$  and
  - (b) if  $\overline{\text{COP}}_{1:K}^A > c_e$ , add  $\bar{t}$  to  $\mathcal{A}$ , i.e. let  $\mathcal{A} = \mathcal{A} + \bar{t}$ ;
 in the deletion step,
  - (a) find  $\underline{t}$  such that  $\text{COP}_{1:K}^{A-\underline{t}} = \underline{\text{COP}}_{1:K}^A$  and
  - (b) if  $\underline{\text{COP}}_{1:K}^A < c_d$ , delete  $\underline{t}$  from  $\mathcal{A}$ , i.e. let  $\mathcal{A} = \mathcal{A} - \underline{t}$ .
- Step 4:* output  $\mathcal{A}$ .

3. Theoretical properties

3.1. *Asymptotic distributions of test statistics in correlation pursuit*

Let us first consider an addition step. We assume that SIR uses a fixed slicing scheme relative to the number of observations  $n$ , i.e. the slices  $S_1, S_2, \dots, S_H$  are fixed (defined by the range of the

response variable) but the number of observations in each slice goes to  $\infty$ . Let  $X_t$  be an arbitrary predictor in  $\mathcal{A}^c$ . Under the null hypothesis  $H_0$  that all the predictors in  $\mathcal{A}^c$  are irrelevant, we have  $\eta_{t1} = \eta_{t2} = \dots = \eta_{tK} = 0$ . Recall that the statistics we propose to measure the contributions of  $X_t$  to the  $K$  profile correlations are  $(\text{COP}_1^{A+t}, \text{COP}_2^{A+t}, \dots, \text{COP}_K^{A+t})'$ , and to measure the overall contribution of  $X_t$  by  $\text{COP}_{1:K}^{A+t}$ . To establish the asymptotic distributions of these statistics, we need to impose a condition on the conditional expectation of  $X_t$  given  $X_{\mathcal{A}}$ .

*Condition 3* (regression condition).  $E(X_t|X_{\mathcal{A}})$  is linear in  $X_{\mathcal{A}}$ .

*Theorem 1.* Assume that conditions 1 and 2 hold, condition 3 holds for  $(X_{\mathcal{A}}, X_t)$  for any  $X_t \in X_{\mathcal{A}^c}$ , and the squared profile correlations  $\lambda_1, \lambda_2, \dots, \lambda_K$  are positive and different from each other. Then, for any given fixed slicing scheme, under the null hypothesis  $H_0$  that all the predictors in  $\mathcal{A}^c$  are irrelevant, we have that

$$(\text{COP}_1^{A+t}, \text{COP}_2^{A+t}, \dots, \text{COP}_K^{A+t}) \rightarrow (Z_{1t}^2, Z_{2t}^2, \dots, Z_{Kt}^2) \tag{16}$$

in distribution and

$$\text{COP}_{1:K}^{A+t} \rightarrow \sum_{l=1}^K Z_{lt}^2 \tag{17}$$

in distribution as  $n \rightarrow \infty$ . Here,  $(Z_{1t}, Z_{2t}, \dots, Z_{Kt})$  follows the multivariate normal distribution with mean 0 and covariance matrix  $W_{Kt}$ . The explicit expression of  $W_{Kt}$  is given in Appendix A.

The asymptotic distributions in theorem 1 can be much simplified if we impose the following condition on the variance of the conditional expectation of  $X_t$  given  $X_{\mathcal{A}}$ .

*Condition 4* (constant variance condition).  $E[\{X_t - E(X_t|X_{\mathcal{A}})\}^2|X_{\mathcal{A}}]$  is a constant.

*Corollary 1.* Assume that conditions 1 and 2 hold, conditions 3 and 4 hold for  $(X_{\mathcal{A}}, X_t)$  for  $X_t \in X_{\mathcal{A}^c}$  and the squared profile correlations  $\lambda_1, \lambda_2, \dots, \lambda_K$  are positive and different from each other. Then, for any given fixed slicing scheme, under the null hypothesis  $H_0$  that all the predictors in  $\mathcal{A}^c$  are irrelevant, we have that  $\text{COP}_1^{A+t}, \text{COP}_2^{A+t}, \dots, \text{COP}_K^{A+t}$  are asymptotically independent and identically distributed as  $\chi^2(1)$ , and  $\text{COP}_{1:K}^{A+t}$  is asymptotically  $\chi^2(K)$ .

Theorem 1 and corollary 1 characterize the asymptotic behaviours of the test statistics for an arbitrary  $X_t$  in  $\mathcal{A}^c$ . In the COP procedure, however, the predictor that attains the maximum value of  $\text{COP}_{1:K}^{A+t}$  among  $t \in \mathcal{A}^c$ , which is  $\overline{\text{COP}}_{1:K}^{\mathcal{A}}$ , is considered a candidate predictor to enter  $\mathcal{A}$ . Our next theorem characterizes the joint asymptotic behaviour of  $\{\text{COP}_{1:K}^{A+t}\}_{t \in \mathcal{A}^c}$  as well as that of  $\overline{\text{COP}}_{1:K}^{\mathcal{A}}$ .

The linearity, regression and constant variance conditions together are more general than the normality assumption on  $X$  because they only need to hold for the basis of the central subspace (e.g.  $B$  or  $\eta_1, \dots, \eta_K$ ) and a given subset of predictors (e.g.  $\mathcal{A}$ ). If we require that the conditions hold for any projection and any given subset of the predictors, however, then it is equivalent to requiring that  $X$  follows a multivariate normal distribution. To understand the joint behaviour of all the COP statistics, in what follows we impose the normality assumption on  $X$ .

Let  $\mathcal{A} = \{t_j\}_{j=1}^d$  and  $\mathcal{A}^c = \{t_j\}_{j=d+1}^p$  denote the collection of currently selected predictors and its complement respectively. Let  $\Sigma_{\mathcal{A}} = \text{cov}(X_{\mathcal{A}})$ ,  $\Sigma_{\mathcal{A}^c} = \text{cov}(X_{\mathcal{A}^c})$ ,  $\Sigma_{\mathcal{A}\mathcal{A}^c} = \text{cov}(X_{\mathcal{A}}, X_{\mathcal{A}^c})$  and  $\tilde{\Sigma}_{\mathcal{A}^c} = \Sigma_{\mathcal{A}^c} - \Sigma_{\mathcal{A}^c\mathcal{A}}\Sigma_{\mathcal{A}}^{-1}\Sigma_{\mathcal{A}\mathcal{A}^c}$ . Note that  $\Sigma_{\mathcal{A}^c\mathcal{A}} = \Sigma'_{\mathcal{A}\mathcal{A}^c}$ . Let  $\tilde{a} = (\tilde{a}_1, \tilde{a}_2, \dots, \tilde{a}_{p-d})'$  be the vector of the diagonal elements of  $\tilde{\Sigma}_{\mathcal{A}^c}$ . Define  $D_{\mathcal{A}^c} = \text{diag}(\tilde{a}_1, \tilde{a}_2, \dots, \tilde{a}_{p-d})$ , and define  $U_{\mathcal{A}^c} = D_{\mathcal{A}^c}^{-1/2}\tilde{\Sigma}_{\mathcal{A}^c}D_{\mathcal{A}^c}^{-1/2}$ .

*Theorem 2.* Assume that

- (a)  $X$  follows a multivariate normal distribution,
- (b) the coverage condition holds and
- (c) the squared profile correlations  $\lambda_1, \lambda_2, \dots, \lambda_K$  are non-zero and different from each other.

Then, for any fixed slicing scheme, under the null hypothesis  $H_0$  that all the predictors in  $\mathcal{A}^c$  are irrelevant, we have

$$(\text{COP}_{1:K}^{A+t_{d+1}}, \text{COP}_{1:K}^{A+t_{d+2}}, \dots, \text{COP}_{1:K}^{A+t_p}) \xrightarrow{D} \left( \sum_{k=1}^K z_{k,d+1}^2, \dots, \sum_{k=1}^K z_{k,p}^2 \right), \tag{18}$$

and

$$\overline{\text{COP}}_{1:K}^A \xrightarrow{D} \max_{t \in \mathcal{A}^c} \left( \sum_{k=1}^K z_{k,t}^2 \right) \tag{19}$$

as  $n \rightarrow \infty$ . Here  $z_k = (z_{k,d+1}, \dots, z_{k,p})$  for  $k = 1, \dots, K$  are mutually independent and each  $z_k$  follows a multivariate normal distribution with mean 0 and covariance matrix  $U_{\mathcal{A}^c}$ .

We now consider deletion steps of the COP procedure. We let  $\mathcal{A}$  denote the current collection of selected predictors before a deletion step, and we let  $X_t$  be an arbitrary predictor in  $\mathcal{A}$ . Note that  $\text{COP}_k^{\tilde{\mathcal{A}}-t} = \text{COP}_k^{\tilde{\mathcal{A}}+t}$ , where  $\tilde{\mathcal{A}} = \mathcal{A} - t$  for  $1 \leq k \leq K$ . Therefore, results similar to those stated in theorem 1 and corollary 1 can be obtained for  $(\text{COP}_1^{\tilde{\mathcal{A}}-t}, \text{COP}_2^{\tilde{\mathcal{A}}-t}, \dots, \text{COP}_K^{\tilde{\mathcal{A}}-t})$  and  $\text{COP}_{1:K}^{\tilde{\mathcal{A}}-t}$  after some modifications described below. First, our current ‘null hypothesis’, which is denoted as  $H_{0t}$ , is that  $X_t$  and the predictors in  $\mathcal{A}^c$  are irrelevant. Second, the regression and constant variance conditions need to be imposed on the conditional expectation of  $X_t$  given  $X_{\tilde{\mathcal{A}}}$  instead. The asymptotic distribution of  $\overline{\text{COP}}_{1:K}^{\tilde{\mathcal{A}}}$ , however, turns out to be fairly complicated if not entirely elusive, because there is not a common null hypothesis for all  $X_t \in \mathcal{A}$ . In what follows, we shall establish two strong results that have implications for properly selecting the thresholds  $c_e$  and  $c_d$ , as well as for the consistency of the COP procedure in selecting true predictors.

### 3.2. Selection consistency of correlation pursuit

Let  $\mathcal{T}$  be the collection of the true predictors under the SDR model. The principal profile correlation directions are  $\eta_1, \eta_2, \dots, \eta_K$ , which form a basis of the central subspace. Assume that  $S_1, \dots, S_H$  is a fixed slicing scheme that is used by SIR. Let  $p_h = P(y \in S_h)$ ,  $\mathbf{v}_K = (\eta'_1 X - E(\eta'_1 X), \dots, \eta'_K X - E(\eta'_K X))'$  and

$$M_{H,K} = \sum_{h=1}^H p_h \mathbf{L}_{h,K} \mathbf{L}'_{h,K}, \tag{20}$$

where  $\mathbf{L}_{h,K} = E(\mathbf{v}_K | Y \in S_h)$ . A few more conditions are needed for the results that we state in the next two theorems.

*Condition 5.*  $X$  follows a multivariate normal distribution with covariance matrix  $\Sigma$  such that  $\tau_{\min} \leq \lambda_{\min}(\Sigma_p) \leq \lambda_{\max}(\Sigma_p) \leq \tau_{\max}$ , where  $\tau_{\min}$  and  $\tau_{\max}$  are two positive constants, and  $\lambda_{\min}(\cdot)$  and  $\lambda_{\max}(\cdot)$  are the minimum and maximum eigenvalues of a matrix respectively.

*Condition 6.* There is a constant  $\omega_H > 0$  such that  $\lambda_{\min}(M_{H,K}) > \omega_H$ .

*Condition 7.* There are constants  $\sigma_0^2$  and  $v > 0$  such that, for any slice  $S_h$  and any two predictors  $X_i$  and  $X_j$ ,  $\text{var}(X_j | Y \in S_h) \leq \sigma_0^2$  and  $\text{var}(X_i X_j | Y \in S_h) \leq \sigma_0^2$  for all  $i, j = 1, \dots, p$ , and  $h = 1, \dots, H$ . In addition,

$$E(|X_j|^l | Y \in S_h) \leq \frac{l!}{2} \text{var}(X_j | Y \in S_h) v^{l-2}$$



and

$$E(|X_i X_j|^l | Y \in S_h) \leq \frac{l!}{2} \text{var}(X_i X_j | Y \in S_h) v^{l-2}, \quad \text{for } l \geq 2.$$

*Condition 8.* Let  $\eta^j = (\eta_{j1}, \eta_{j2}, \dots, \eta_{jK})'$  in which  $\eta_{jk}$  is the coefficient of  $X_j$  in the  $k$ th principal correlation direction  $\eta_k$ . There is a positive constant  $\varpi > 0$  and a non-negative constant  $\xi_0$ , such that  $\|\eta^j\|^2 > \varpi n^{-\xi_0}$  for  $j \in \mathcal{T}$ , where  $\|\cdot\|$  denotes the standard  $L_2$ -norm.

*Condition 9.*  $\lim_{n \rightarrow \infty} (p) = \infty$  and  $p = o(n^{\varrho_0})$  with  $\varrho_0 \geq 0$  and  $2\varrho_0 + 2\xi_0 < 1$ .

Condition 5 ensures that the variances of the predictors are on a comparable scale and that they are not strongly correlated. Condition 6 assumes a lower bound for the eigenvalues of  $M_{H,K}$ , which is slightly stronger than the coverage condition that ensures SIR to recover all the SDR directions. Condition 7 imposes conditions on the moments of the conditional expectations of  $X$  given  $Y \in S_h$  so that the Bernstein inequalities hold for the conditional sample means. Condition 8 assumes that the coefficients of any true predictors do not decrease to 0 too fast as both  $n$  and  $p$  increase; otherwise, such predictors will not be identifiable asymptotically. Condition 9 allows  $p$  to increase as  $n$  increases, but their rates are constrained. Similar conditions have been used by others for establishing variable selection results for stepwise procedures in linear regression (Wang, 2009; Fan and Lv, 2008).

*Theorem 3.* Let  $\mathcal{A}$  be the set of currently selected predictors and let  $\mathcal{T}$  be the set of true predictors. Let  $\vartheta = \varpi \omega_H \tau_{\min}^2 / 2\tau_{\max}$ . Assume that conditions 5–9 hold. Then, we have

$$P\left\{ \min_{\mathcal{A}: \mathcal{A}^c \cap \mathcal{T} \neq \emptyset} \max_{t \in \mathcal{A}^c \cap \mathcal{T}} (\text{COP}_{1:K}^{A+t}) \geq \vartheta n^{1-\xi_0} \right\} \rightarrow 1, \tag{21}$$

for any fixed slicing scheme as  $n \rightarrow \infty$ .

The probability statement (21) is not just about one given collection of predictors. It considers all the possible collections that do not include all the true predictors yet, i.e.  $\{\mathcal{A}: \mathcal{A}^c \cap \mathcal{T} \neq \emptyset\}$ . In other words, it considers all the possible scenarios where the null hypothesis  $H_0$  is not true. Further note that  $\max_{t \in \mathcal{A}^c \cap \mathcal{T}} (\text{COP}_{1:K}^{A+t}) \neq \overline{\text{COP}}_{1:K}^A$ . Because  $\max_{t \in \mathcal{A}^c} (\text{COP}_{1:K}^{A+t}) \geq \max_{t \in \mathcal{A}^c \cap \mathcal{T}} (\text{COP}_{1:K}^{A+t})$ , from equation (21), we have

$$P\left\{ \min_{\mathcal{A}: \mathcal{A}^c \cap \mathcal{T} \neq \emptyset} (\overline{\text{COP}}_{1:K}^A) \geq \vartheta n^{1-\xi_0} \right\} \rightarrow 1. \tag{22}$$

This result implies that by setting  $c_e$  to  $\vartheta n^{1-\xi_0}$  or smaller, if the COP procedure has not collected all the true predictors yet, then with probability going to 1 (as  $n \rightarrow \infty$ ) it will continue to select a predictor to the current collection. Thus, the addition step of COP will not stop until all the true predictors have been selected. Another way to interpret expression (22) is that the selection power of the COP procedure converges to 1 asymptotically.

*Theorem 4.* Assume that conditions 5–9 hold. Then we have

$$P\left\{ \max_{\mathcal{A}: \mathcal{A}^c \cap \mathcal{T} = \emptyset} \max_{t \in \mathcal{A}^c} (\text{COP}_{1:K}^{A+t}) < C n^{\varrho} \right\} \rightarrow 1, \tag{23}$$

for  $\varrho > \frac{1}{2} + \varrho_0$ , and any positive constant  $C$ , under any fixed slicing scheme with  $n \rightarrow \infty$ .

Theorem 4 has two implications. The first regards the addition step of COP. Once all the true predictors have been selected, i.e.  $\mathcal{A}^c \cap \mathcal{T} = \emptyset$ , the probability that it will select a false predictor from  $\mathcal{A}^c$  converges to 0. The second implication concerns the deletion step. Consider one collection of selected predictors  $\tilde{\mathcal{A}}$  and assume that  $\tilde{\mathcal{A}}$  contains all the true predictors and also

some irrelevant ones, i.e.  $\tilde{\mathcal{A}} \supset \mathcal{T}$ . Clearly,

$$\underline{\text{COP}}_{1:K}^{\tilde{\mathcal{A}}} \leq \min_{t \in \tilde{\mathcal{A}} - \mathcal{T}} (\text{COP}_{1:K}^{\tilde{\mathcal{A}}-t}) \leq \max_{\mathcal{A}:\mathcal{A}^c \cap \mathcal{T} = \emptyset} \max_{t \in \mathcal{A}^c} \left( \sum_{k=1}^K \text{COP}_k^{\mathcal{A}+t} \right). \tag{24}$$

Therefore,

$$P(\underline{\text{COP}}_{1:K}^{\tilde{\mathcal{A}}} < Cn^\varrho) \rightarrow 1. \tag{25}$$

In other words, with probability going to 1, the COP procedure will delete an irrelevant predictor from the current collection.

One possible choice of the thresholds is  $\chi_c^2 = \vartheta n^{1-\xi_0}$  and  $\chi_d^2 = \vartheta n^{1-\xi_0}/2$ . From theorem 3, asymptotically, the COP algorithm will not stop selecting variables until all the true predictors have been included. Moreover, once all the true predictors have been included, according to theorem 4, all the redundant variables will be removed from the selected variables.

### 4. Implementation issues

When implementing the COP algorithm, we need to specify the number of profile correlation directions  $K$ , the thresholds for the addition and deletion steps  $c_e$  and  $c_d$ , and the slicing scheme, particularly the number of slices  $H$ . A proper specification of these tuning parameters is critical for the success of the COP algorithm.

#### 4.1. Slicing schemes and the choice of $H$

Li (1991) suggested that, in terms of estimation, the performance of SIR is robust to the number of slices in general. The COP algorithm uses SIR to derive test statistics for selecting variables. It is of interest to understand the effect of a slicing scheme on the testing procedures involved. Again, we consider an addition step in the COP procedure. Let  $\mathcal{A}$  be the current collection of selected predictors. Let  $X_t$  be an arbitrary predictor in  $\mathcal{A}^c$ .

*Theorem 5.* Assume that  $X$  follows a multivariate normal distribution. Then, for any given fixed slicing scheme, we have

$$P\left(\frac{\text{COP}_{1:K}^{\mathcal{A}}}{n} \geq C_{H,\mathcal{A}+t}\right) \rightarrow 1, \quad \text{as } n \rightarrow \infty, \tag{26}$$

where

$$C_{H,\mathcal{A}+t} = (\tilde{\eta}_{t,\mathcal{A}})' M_{H,K} \tilde{\eta}_{t,\mathcal{A}} / \sigma_{t,\mathcal{A}}^2, \tag{27}$$

$\sigma_{t,\mathcal{A}}^2 = \text{var}(X_t | X_{\mathcal{A}})$ ,  $\tilde{\eta}_{t,\mathcal{A}} = \text{cov}(X_t, \mathbf{v}_K | X_{\mathcal{A}})$  and  $M_{H,K}$  is defined in equation (20).

The difference between theorem 1 and theorem 5 is that the latter does not assume that  $X_t$  is an irrelevant predictor. When  $X_t$  is indeed a true predictor, then  $\eta^t$  is not a zero vector and  $\max_{t \in \mathcal{A}^c \cap \mathcal{T}} (C_{H,\mathcal{A}+t})$  is greater than 0. The larger  $C_{H,\mathcal{A}+t}$  is, the more likely  $X_t$  will be added to  $\mathcal{A}$ . The next result shows that a finer slicing scheme leads to higher power for the addition step by COP. For any two different slicing schemes  $S = (S_1, \dots, S_{H_1})$  and  $S' = (S'_1, \dots, S'_{H_2})$ , we say that  $S'$  is a refinement of  $S$ , which is denoted by  $S' \leq S$ , if, for any  $S'_h \in S'$ , there is an  $S_h \in S$  such that  $S'_h \subseteq S_h$ .

*Proposition 2.* Suppose that  $S$  and  $S'$  are two slicing schemes such that  $S' \preceq S$ . Then, for any  $\eta \in \mathbb{R}^K$ , we have

$$\eta' M_{H_2, K} \eta \geq \eta' M_{H_1, K} \eta, \tag{28}$$

where  $M_{H_2, K}$  and  $M_{H_1, K}$  are defined as in equation (20) under the slicing schemes  $S'$  and  $S$  respectively.

Proposition 2 implies that the constant  $C_{H, \mathcal{A}}$  in theorem 5 becomes larger when a finer slicing scheme is used. This further suggests that the power of the COP procedure in selecting true predictors tends to increase if a slicing scheme uses a larger number of slices. However, when a slicing scheme uses a larger number of slices, the number of observations in each slice will decrease, which makes the estimate of  $E(X|y \in S_h)$  less accurate and further makes the estimates of  $M = \text{cov}\{E(X|Y)\}$  and its eigenvalues  $\lambda_1, \dots, \lambda_K$  less stable. The success of the COP procedure hinges on a good balance between the number of slices and the number of observations in each slice. We observed from intensive simulation studies that, with a reasonable number of observations in each slice (say 20 or more), a large number of slices is preferred.

#### 4.2. Choice of $c_e$ and $c_d$

Section 3 has characterized the asymptotic distributions or behaviours of the test statistics that are involved in the COP procedure. In theory, these results (theorems 4 and 5) can be used for choosing the thresholds  $c_e$  and  $c_d$ . In practice, however, these thresholds should be used with much caution because of the following concerns. First, the distributions that were obtained in Section 3 are for a single addition or deletion step and under various assumptions. Second, the distributions are valid only in an asymptotic sense. In what follows, we propose to use a cross-validation (CV) procedure for selecting  $c_e$  and  $c_d$ .

Let  $\{\alpha_i\}_{1 \leq i \leq d}$  be a prespecified grid on a subinterval in  $(0, 1)$  and  $\{\chi_{\alpha_i, K}^2\}_{1 \leq i \leq d}$  be the collection of the  $100\alpha_i$ th percentile of  $\chi_K^2$ . For convenience, we consider only the  $m$  pairs of  $c_e = \chi_{\alpha_i, K}^2$  and  $c_d = \chi_{\alpha_i - 0.05, K}^2$  for  $1 \leq i \leq m$ . Note that  $c_d < c_e$  and that there is only one tuning parameter that we need to determine. We follow the general fivefold CV scheme to select the best pair of  $c_e$  and  $c_d$ . We randomly divide the original data into five equal-sized subsets and then apply the COP procedure to any four subsets to generate the estimation and variable selection results. The remaining subset of the data is used to test the model and to generate a performance measurement. The performance measurements are averaged and the result is used as the CV score. We choose the pair of  $c_e$  and  $c_d$  that maximizes the CV score.

We define the performance measure that is used in the CV procedure as follows. Suppose that  $\mathcal{A}$  is the collection of selected predictors and  $\eta_{1, \mathcal{A}}, \dots, \eta_{K, \mathcal{A}}$  are the estimates of the principal profile correlation directions produced by applying the COP procedure to the training data set. We consider the first principal profile correlation direction first. Recall that  $\eta_{1, \mathcal{A}}$  is the direction that achieves the maximum correlation of a linear projection of  $X$  and the transformed response  $Y$ , and the optimal transformation is  $T_1(Y) = E(\eta'_{1, \mathcal{A}} X | Y)$  (theorem 3.1 in Chen and Li (1998)). With  $\eta_{1, \mathcal{A}}$  estimated by  $\hat{\eta}_{1, \mathcal{A}}$  by using the training data, we apply LOESS proposed by Cleveland (1979) to fit  $T_1(Y)$  using the training data and we denote the fitted transformation as  $\hat{T}_1(\cdot)$ . Let  $\tilde{X}$  and  $\tilde{Y}$  be the data matrix and the response vector of the testing data set. Then, the squared profile correlation between  $\tilde{X}$  and  $\tilde{Y}$  based on the direction  $\hat{\eta}_{1, \mathcal{A}}$  and transformation  $\hat{T}_1(\cdot)$  is computed as  $\text{corr}^2\{\hat{T}_1(\tilde{Y}), \hat{\eta}'_{1, \mathcal{A}} \tilde{X}\}$ . Similarly, the squared profile correlations between  $\tilde{X}$  and  $\tilde{Y}$  along  $\hat{\eta}_{2, \mathcal{A}}, \dots, \hat{\eta}_{K, \mathcal{A}}$  can be calculated. The overall performance measure is defined to be

$$\text{PC} = \sum_{k=1}^K \text{corr}^2\{\hat{T}_k(\tilde{Y}), \hat{\eta}'_{k, \mathcal{A}} \tilde{X}\}. \tag{29}$$

The CV score for any pair  $(c_e, c_d)$  is defined to be the average PC over the five possible partitions of the training–test data sets.

### 4.3. Selection of the number of directions $K$

To determine  $K$ , the number of principal profile correlation directions, we adopt a Bayesian information criterion type of criterion proposed by Zhu *et al.* (2006). For any given  $K$  between 1 and  $J$ , where  $J \leq \max(n, p)$  is a reasonable upper bound chosen by the user, we apply the COP procedure with  $K = k$ . Suppose that the resulting collection of the selected predictors is  $\mathcal{A}_k$  and the cardinality of  $\mathcal{A}_k$  is  $p_k$ . Using the data involving only the selected predictors, we can estimate  $M = \text{cov}\{E(X_{\mathcal{A}_k} | Y)\}$  as before and denote the result as  $\hat{M}$ . Let  $\hat{\theta}_1 \geq \hat{\theta}_2 \geq \dots \geq \hat{\theta}_p$  be the eigenvalues of  $\hat{M} + I_{p_k}$ , where  $I_{p_k}$  is the  $p_k \times p_k$  identity matrix, and let  $\tau$  be the number of  $\hat{\theta}_i$ s that are greater than 1. Define

$$G(k) = -\log\{L(k)\} + \frac{\log(n)}{2}k(2p_k - k + 1), \tag{30}$$

where  $\log\{L(k)\} = \sum_{i=\min(\tau, k)+1}^p \{\log(\hat{\theta}_i) + 1 - \hat{\theta}_i\}$ . We choose  $K = \arg \min_{1 \leq k \leq J} \{G(k)\}$ . In the original criterion that was proposed by Zhu *et al.* (2006), they showed that the criterion produces a consistent estimate of  $K$  for fixed  $p_k$ . Our simulation study shows that the modified criterion leads to the correct specification of  $K$  for the COP procedure and can be generally used in practice.

## 5. Simulation study

We have performed extensive simulation studies to compare the COP algorithm with a few existing variable selection methods and we shall present three examples in this section. When implementing the COP algorithm in these examples, we use the CV procedure and the  $G$  information criterion that was discussed in the previous section to select the thresholds  $c_e$  and  $c_d$  and the dimensionality  $K$  respectively. The grid that was used for selecting  $c_e$  is  $\{\chi_{0.90, K}^2, \chi_{0.95, K}^2, \chi_{0.99, K}^2, \chi_{0.999, K}^2, \chi_{0.9999, K}^2\}$ , and the associated grid for selecting  $c_d$  is  $\{\chi_{0.85, K}^2, \chi_{0.90, K}^2, \chi_{0.94, K}^2, \chi_{0.949, K}^2, \chi_{0.9499, K}^2\}$ . The range that was used for selecting  $K$  is from 1 to 4 (i.e.  $J = 4$ ). For SSIR, we used the grid  $\{0, 0.1, \dots, 0.9, 1\} \times \{0, 0.1, \dots, 0.9, 1\}$  to select the pair of tuning parameters that leads to its best performance. Both COP and SSIR involve slicing the range of the response variable, for which we use the same scheme to facilitate fair comparison.

### 5.1. Linear models

In this example, we consider the linear model

$$Y = X\beta + \sigma\varepsilon, \tag{31}$$

where  $X = (X_1, X_2, \dots, X_p)'$  follows a  $p$ -variate normal distribution with mean 0 and covariances  $\text{cov}(X_i, X_j) = \rho^{|i-j|}$  for  $1 \leq i, j \leq p$ , and  $\varepsilon$  is independent of  $X$  and follows  $N(0, 1)$ . The variable selection methods that we compare the COP procedure with include the lasso, SCAD (Fan and Li, 2001), MARS and SSIR (Li, 2007). The R packages `SIS`, `lasso` and `mda` are used to run SCAD, the lasso and MARS respectively. The tuning parameters that are involved in SCAD and the lasso are selected by CV. We use the code that was provided by the original authors to run SSIR. In this example, we consider two specifications of the linear model given below: scenario 1,

$$p = 8, \quad \beta = (3, 1.5, 2, 0, 0, 0, 0, 0)', \quad \sigma = 3, \quad \rho = 0.5;$$

scenario 2,

$$p = 1000, \quad \beta = (3, 1.5, 1, 1, 2, 1, 0.9, 1, 1, 1, 0, \dots, 0)', \quad \sigma = 1, \quad \rho = 0.5.$$

Under scenario 1, model (31) involves three true predictors and five irrelevant variables, and was originally used in Tibshirani (1996) and Fan and Li (2001) to demonstrate the empirical performances of the lasso and SCAD. We randomly generated 100 data sets from scenario 1, each with 40 data points (i.e.  $n = 40$ ), and applied the aforementioned methods to the data sets. Two quantities were used to measure the variable selection performance of each method, which are the average number of irrelevant predictors falsely selected as true predictors (which is denoted by FP) and the average number of true predictors falsely excluded as irrelevant predictors (which is denoted by FN). Under scenario 1.1, the FPs and FNs range from 0 to 5 and from 0 to 3 respectively, with small values indicating good performances in variable selection. The FP- and FN-values of the methods tested are reported in Table 1.

Under scenario 2, model (31) involves 10 true predictors and 990 irrelevant predictors and is clearly more challenging than scenario 1. We randomly generated 100 data sets each with 200 data points (i.e.  $n = 200$ ) from scenario 2. In each data set,  $n < p$ . Similarly to scenario 1, we applied the methods mentioned above to the data sets and report the FP- and FN-values of these methods in Table 1. The tuning parameters in all these methods are determined by CV.

From the left-hand panel of Table 1, under scenario 1, SSIR has the lowest FP-value (FP = 0.19), i.e. the average number of irrelevant variables selected by SSIR is 0.19, and COP has the third lowest FP-values (0.71). The other methods tend to have more false positive results than SSIR and COP. In terms of FNs, the order of the methods ranked from the lowest to the highest is MARS, SCAD, the lasso, COP and SSIR. The relative sub-par performance of COP and SSIR is because these two methods are developed for variable selection under models that are more general than the linear model.

From the right-hand panel of Table 1, under scenario 2, COP has the lowest FP-value (FP = 2.28). In terms of FN, the lasso and MARS have the lowest value with COP following modestly behind. Compared with MARS, COP has a much lower FP-value and a slightly higher FN-value. SSIR breaks down under scenario 2 because the variance-covariance matrix of  $X$  is no longer invertible. In terms of both FP and FN, COP outperformed SCAD under this scenario. One explanation for this comparison result is that SCAD involves non-convex optimization and can be unstable in implementation.

**Table 1.** Performance comparison under linear models†

Method	Results for $p=8, n=40,$ $\sigma=3, \rho=0.5$		Results for $p=1000, n=200,$ $\sigma=1, \rho=0.5$	
	FP (0, 5)	FN (0, 3)	FP (0, 990)	FN (0, 10)
LASSO	0.77 (0.093)	0.16 (0.037)	8.87 (0.586)	0.00 (0.000)
SCAD	0.67 (0.094)	0.10 (0.030)	6.05 (0.926)	1.16 (0.150)
MARS	4.00 (0.059)	0.04 (0.020)	30.64 (0.165)	0.00 (0.000)
SSIR	0.19 (0.051)	0.96 (0.068)	‡	‡
COP	0.71 (0.080)	0.56 (0.066)	2.28 (0.203)	0.75 (0.095)

†FP is the average number of irrelevant variables that are falsely selected by the method, and FN is the average number of true variables that are falsely excluded by the method; the numbers in parentheses are the standard error of FP or FN.

‡The algorithm broke down.

5.2. Non-linear multiple-index models

In this example, we consider the multiple-index model

$$Y = \frac{X_1 + X_2 + \dots + X_d}{0.5 + (1.5 + X_2 + X_3 + X_4)^2} + \sigma\varepsilon, \tag{32}$$

where  $X_1, \dots, X_p$  are independent identically distributed  $N(0, 1)$  random variables,  $\varepsilon$  is  $N(0, 1)$  and independent of  $X$ , and  $d$  and  $\sigma$  are parameters that need to be further specified. This model was originally used in Li (1991) for demonstrating the performance of SIR. It is not difficult to see that, given the two projections  $X_1 + X_2 + \dots + X_d$  and  $X_2 + X_3 + X_4$ ,  $Y$  and  $X$  are independent of each other. The dimensionality of the central subspace of model (32) is 2, and the collection of true predictors is  $\{X_1, \dots, X_d\} \cup \{X_2, X_3, X_4\}$ . Because model (32) is non-linear, methods that were designed specifically for linear models such as the lasso and SCAD are clearly at a disadvantage. Therefore, in this example, we compare the performances of MARS, SSIR and COP only.

By specifying  $p$ ,  $d$  and  $\sigma$  at different values, we have the following three scenarios: scenario 3,

$$p = 30, \quad d = 3, \quad \sigma = 0.1;$$

scenario 4,

$$p = 30, \quad d = 3, \quad \sigma = 2;$$

scenario 5,

$$p = 400, \quad d = 8, \quad \sigma = 0.1.$$

For each scenario, we generated 100 data sets each with 200 observations (i.e.  $n = 200$ ) and applied MARS, SSIR and COP to each data set. The resulting FP- and FN-values are reported in Table 2.

For scenario 3, MARS achieved the lowest FN-value (0.03), but its FP-value was unacceptably high (16.55); SSIR had the lowest FP-values, but its FN-value was the highest among the three. The FP- and FN-values of COP were between the extremes. It appears that the performances of SSIR and COP are similar under scenario 3. For scenario 4, COP outperformed SSIR in terms of both FP- and FN-values. MARS again achieved the lowest FN-value (0.32) at the expense of

**Table 2.** Performance comparison under the multiple-index model†

Method	Results for $\sigma = 0.1$ , $p = 30, d = 3$		Results for $\sigma = 2$ , $p = 30, d = 3$		Results for $\sigma = 0.1$ , $p = 400, d = 8$	
	FP (0, 26)	FN (0, 4)	FP (0, 26)	FN (0, 4)	FP (0, 292)	FN (0, 8)
MARS	16.55 (0.174)	0.03 (0.017)	17.18 (0.186)	0.32 (0.053)	‡	‡
SSIR	0.12 (0.033)	0.91 (0.029)	4.14 (0.288)	1.76 (0.115)	‡	‡
COP	1.88 (0.149)	0.83 (0.038)	3.26 (0.210)	1.71 (0.104)	8.93 (0.576)	0.18 (0.081)

†FP is the average number of irrelevant variables that are falsely selected by the method, and FN is the average number of true variables that are falsely excluded by the method; the numbers in parentheses is the standard error of FP or FN.

‡The algorithm broke down.

an unacceptable FP-value (17.18). Scenario 5 is the most challenging among the three scenarios, in which the number of predictors exceeds the number of observations. Both MARS and SSIR broke down under this scenario. However, COP still demonstrated an excellent performance with its FP- and FN-values reasonably low.

5.3. Heteroscedastic models

In the previous examples, the true predictors affect only the mean response. In this example, we consider the heteroscedastic model

$$Y = \frac{0.2\varepsilon}{1.5 + \sum_{j=1}^p \beta_{j,1} X_j}, \tag{33}$$

where  $X = (X_1, X_2, \dots, X_p)'$  follows a  $p$ -variate normal distribution with mean 0 and covariances  $\text{cov}(X_i, X_j) = \rho^{|i-j|}$  for  $1 \leq i, j \leq p$ ,  $\varepsilon$  is independent of  $X$  and follows  $N(0, 1)$ , and  $\beta_{j,1}$  equals 1 for  $1 \leq j \leq 8$  and equals 0 for  $j \geq 9$ . Note that the central subspace is spanned by  $\beta_1 = (\beta_{1,1}, \beta_{2,1}, \dots, \beta_{p,1})'$  and the number of true predictors is 8. We further specify  $\rho$  and  $p$  in equation (33) and consider the following three scenarios: scenario 6,

$$\rho = 0, \quad p = 500;$$

scenario 7,

$$\rho = 0, \quad p = 1000;$$

scenario 8,

$$\rho = 0.3, \quad p = 1500.$$

For each scenario, we generated 100 data sets each with  $n = 1000$  observations and applied MARS, SSIR and COP to the data sets. The FP- and FN-values of the three methods are listed in Table 3.

Under scenario 6, both SSIR and COP outperformed MARS. The FN-value of SSIR (0.99) is less than that of COP (1.21), but the FP-value (52.54) is much larger than that of COP (5.71). Under both scenarios 7 and 8, in which  $p$  is much larger than  $n$ , SSIR broke down, but COP still demonstrated excellent performances. The performances of MARS under these two scenarios were fairly poor.

**Table 3.** Performance comparison under the heteroscedastic model†

Method	Results for $\rho = 0$ , $n = 1000, p = 500$		Results for $\rho = 0$ , $n = 1000, p = 1000$		Results for $\rho = 0.3$ , $n = 1000, p = 1500$	
	FP (0, 492)	FN (0, 8)	FP (0, 992)	FN (0, 8)	FP (0, 1492)	FN (0, 8)
MARS	212.15 (0.428)	4.83 (0.116)	230.33 (0.372)	6.16 (0.129)	236.60 (0.524)	6.84 (0.126)
SSIR	52.54 (1.970)	0.88 (0.149)	‡	‡	‡	‡
COP	5.79 (0.365)	1.21 (0.030)	13.14 (0.734)	1.29 (0.037)	21.36 (0.937)	1.5 (0.039)

†FP is the average number of irrelevant variables that are falsely selected by the method, and FN is the average number of true variables that are falsely excluded by the method; the numbers in parentheses are the standard error of FP or FN.

‡The algorithm broke down.

## 6. Application: predict gene expression from sequences by using next generation sequencing data

Embryonic stem cells (ESCs) maintain self-renewal and pluripotency as they have the ability to differentiate into all cell types. To enhance the understanding of the ESC development, predictive models, such as regression models, can be constructed in which the gene expression is regarded as the response variable and various features that are associated with gene regulating transcription factors (TFs) are taken as the predictors. Examples of such features include motif scores based on position-specific weight matrices of motifs recognized by the TFs (Conlon *et al.*, 2003), and ‘ChIP-chip’ log-ratios.

Recently, the emerging next generation sequencing technologies, in particular, ‘RNA-Seq’ and ‘ChIP-Seq’, have offered researchers an unprecedented opportunity to build predictive models for complex biological processes such as gene regulation. Compared with the traditional hybridization-based methods, such as microarrays, RNA-Seq and ChIP-Seq provide more accurate quantification of gene expression and TF–DNA binding locations respectively (Mortazavi *et al.*, 2008; Wilhelm *et al.*, 2008; Nagalakshmi *et al.*, 2008; Boyer *et al.*, 2005; Johnson *et al.*, 2007).

To quantify gene expression in RNA-Seq data, one may calculate RPKM, the number of reads per kilobase of exon region per million mapped reads, which has been shown to be proportional to the gene expression levels (Cloonan *et al.*, 2008). From ChIP-Seq data, Ouyang *et al.* (2009) proposed a feature named the transcription factor association strength (TFAS), which has been shown to explain a much higher proportion of gene expression variation than traditional predictors in predictive models. In particular, for each TF, the TFAS for each gene is computed as a weighted sum of the corresponding ChIP-Seq signal strengths, where the weights reflect the proximity of the signal to the gene. We here examine whether we can build a better predictive model for gene expressions by combining both TFASs and motif scores of TFs in mouse ESCs.

To achieve this, we compiled a data set consisting of gene expressions, TFASs and motif scores. In this data set, the RPKMs were calculated as gene expression levels from RNA-Seq data in mouse ESCs (Cloonan *et al.*, 2008). The TFASs of 12 TFs were calculated from the ChIP-Seq experiments in mouse ESCs (Chen *et al.*, 2008). In addition, we supplement this data set with motif scores of putative mouse TFs. From the TF database TRANSFAC, we compiled a list of 300 mouse TF binding motifs. For each gene, a matching score was calculated by using the scoring system that was described in Zhong *et al.* (2005) for each TF binding motif. The matching score can be considered intuitively as the expected number of occurrences of a TF binding motif on the gene’s promoter region. To build a predictive model in mouse ESCs, we treat the gene expression as the response variable and the 12 TFASs as well as the 300 TF motif matching scores as predictors. More precisely, the response is a vector with 12408 entries and the data matrix is a  $12408 \times 312$  matrix with  $(i, j)$ th entry representing the TFAS score of the  $i$ th gene’s promoter region for TF  $j$  if  $j \leq 12$ , representing the matching score of the  $i$ th gene’s promoter region for TF  $j$  if  $j > 12$ .

We applied COP to this data set. The procedure identified two principal directions and selected in total 42 predictors. The first squared profile correlation is  $\lambda_1 = 0.67$ , and the second squared profile correlation is  $\lambda_2 = 0.20$ . Among the 12 TFASs calculated from ChIP-Seq, eight were selected by COP. In particular, Oct4 is a well-known master regulator regulating pluripotency, and Klf4 regulates differentiation (Cai *et al.*, 2010). Evidence also suggests that, at these early stages of development, STAT3 activation is required for self-renewal of ESCs (Matsuda *et al.*, 1999). Among the 300 TF motif scores, 34 of them are selected by COP. To understand further what extra information TF motif scores provide, we annotate the functions of the 34 TFs. It is of interest to note that 24 of the 34 selected motifs correspond to TFs that are either regulators



**Table 4.** Motifs identified

Development	COUP-TF, AP2, Sp1, CHOP C/EBpalph, NF-AT Pax, Pax8, GABP, En1, TTF1 PITX2, NKx2-2, HIXA4, ZF5, PPAR direct repeat 1
Cancer	IRF1, EVI1, NF1, GKLf, Whn VDR, POU6F1, Arnt, Cdx2
8 selected TFASs	E2F1, Mycn, ZFx, Klf4 Tcfcp2/1, Oct4, Stat3, Smad1

for development or cancer related; Table 4. Since ESCs are in a developmental phase, it is not surprising to have active TFs regulating general development. Some recent evidence suggests that tumour suppressors that control cancer cell proliferation also regulate stem cell self-renewal (Pardal *et al.*, 2005). Thus, a careful study of these cancer-related TFs could lead to a better understanding of the stem cell regulatory network.

## 7. Discussion

The contribution of the COP procedure to the development of variable selection methodologies for high dimensional regression analysis is twofold. First, it does not impose any assumption on the relationship between the response variable and the predictors, and the SDR framework that the COP procedure relies on includes fully non-parametric models as special cases. Therefore, COP can be considered a model-free variable selection procedure that is applicable in any high dimensional data analysis. Second, as demonstrated by our simulation studies, the COP procedure can effectively handle hundreds of thousands of predictors, which can be extremely challenging to other existing methods for variable selection beyond linear or parametric models. Like linear stepwise regression, the COP procedure may encounter issues that are typical to stepwise procedures as discussed in Miller (1984). Nonetheless, we believe that the COP procedure should become an indispensable member of the repository of variable selection tools and we recommend its broad use. When a parametric model is postulated for the relationship between the response and the predictor variables and model-specific variable selection methods are available, we recommend the use of COP together with these methods as a safeguard against possible model misidentification. We have implemented the COP procedure using programming language R, and the R package can be downloaded from <http://cran.r-project.org/web/packages/COP/> or requested from the authors directly.

As a trade-off, the COP procedure imposes various assumptions on the distribution of the predictors, of which the linearity assumption is the most fundamental and crucial. When the linearity condition is required to hold for any lower dimensional projection, it is equivalent to requiring that the joint distribution of the predictors is elliptically contoured (Eaton, 1986). Hall and Li (1993) established the fact that low dimensional projections from high dimensional data approximately satisfy the linearity condition, which to a certain degree alleviates the concern of the linearity assumption and explains why SIR and the COP procedure worked well under mild violation of the assumption. When the linearity condition is heavily violated, data reweighting schemes such as the Voronoi reweighting scheme (Cook and Nachtsheim, 1994) can be used to correct the violation. We plan to incorporate such schemes into the COP procedure in the future.

When the number of the predictors is extremely large, the performance of the COP procedure can be compromised. This is also so for variable selection methods under the linear

model. Lately, Fan and Lv (2008) have advocated a two-step approach to attack so-called ultra-high dimensionality. The first step is to perform screening to reduce the dimensionality from ultrahigh to high or moderately high, and then, in the second step, variable selection methods are applied to identify the true predictors. The same approach can be used for variable selection under the SDR framework. More precisely, we can apply the forward COP procedure, which is simply the COP procedure with the deletion step removed, to reduce the dimensionality of a problem from ultrahigh to moderately high. The forward COP procedure is much easier to implement and computationally more efficient than the COP procedure. Then, the usual COP procedure is applied to the reduced data to select the true predictors. This approach is currently under investigation and the results will be reported in a future publication.

**Acknowledgements**

We thank Xuming He, Steve Portnoy and John Marden for helpful suggestions. This work was supported by National Institutes of Health grant U01 ES016011, a Department of Energy grant from the Office of Science (Biological and Environmental Research) and National Science Foundation grant DMS 1120256 to WZ, National Institutes of Health grant R01-HG02518-02 and National Science Foundation grant DMS 1007762 to JL and National Science Foundation grant DMS 0707004 to YZ.

**Appendix A**

**A.1. Proof of proposition 1**

Let  $\mathcal{S}^\perp(B)$  denote the space of vectors such that, for any  $\rho \in \mathcal{S}^\perp(B)$  and any  $\beta \in \mathcal{S}(B)$ ,  $\rho' \Sigma \beta = 0$ . Let  $\mathcal{S}^\perp(\tilde{K})$  be the space of vectors such that for any  $\rho \in \mathcal{S}^\perp(\tilde{K})$   $\rho' \Sigma \eta_k = 0$  for  $k = 1, \dots, \tilde{K}$ . We shall show that  $\mathcal{S}^\perp(B) \subseteq \mathcal{S}^\perp(\tilde{K})$ , which means, for any  $\rho \in \mathcal{S}^\perp(B)$ ,  $P(\rho) = 0$ . First, because, for any  $T$ ,  $T(Y) \perp \eta' X | B' X$ , then

$$\text{cov}\{T(Y), \eta' X\} = E\{T(Y)\eta' X\} = E\{E\{T(Y)|B' X\} E(\eta' X|B' X)\}.$$

Because of the linearity condition, for any  $\rho \in \mathcal{S}^\perp(B)$ ,  $E(\rho' X|B' X) = c_1 \beta'_1 X + \dots + c_K \beta'_K X$ , where  $c_1, \dots, c_K$  are linear coefficients. In addition, since  $\text{cov}(\rho' X, \beta'_k X) = 0$  for  $k = 1, \dots, K$ ,  $E(\rho' X|B' X) = 0$ . Consequently,

$$\text{corr}^2\{T(Y), \rho' X\} = \frac{\text{cov}^2\{T(Y), \rho' X\}}{\text{var}\{T(Y)\}\text{var}(\rho' X)} = 0,$$

$P(\rho) = 0$  and  $\mathcal{S}^\perp(B) \subseteq \mathcal{S}^\perp(\tilde{K})$ . Proposition 1 holds.

**A.2. Proof of theorem 1**

Without loss of generality, we let  $\mathcal{A} = \{1, \dots, d\}$  and  $t = d + 1$ . Let  $X^{(j)}$  be the vector of  $n$  independent identically distributed observations of the  $j$ th variable for  $j = 1, \dots, d + 1$ . We assume that the predictors have been centred to have zero sample mean. Denote  $\mathbf{X}_{n \times j} = (X^{(1)}, \dots, X^{(j)})$  for  $j = d, d + 1$ . We let

$$\hat{M}^{(j)} = \sum_{h=1}^H \frac{n_h}{n} \bar{\mathbf{X}}_h^{(j)T} \bar{\mathbf{X}}_h^{(j)} \quad \text{for } j = d, d + 1$$

where  $\bar{\mathbf{X}}_h^{(j)}$  ( $j = d, d + 1$ ) is the average of the first  $j$  variables for those individuals whose responses fall into the  $h$ th slice  $S_h$ ,  $h = 1, \dots, H$ . Let  $n_h$  be the number of observations in the  $h$ th slice,  $h = 1, \dots, H$ . Let  $\hat{\lambda}_i^{(j)}$  be the  $i$ th largest eigenvalue of  $\hat{\Sigma}_j^{-1} \hat{M}^{(j)}$  for  $j = d, d + 1$ , where  $\hat{\Sigma}_j$  is the sample variance-covariance matrix of  $\mathbf{X}_{n \times j}$ . It is difficult to see the asymptotic distribution of  $\hat{\lambda}_i^{(d+1)} - \hat{\lambda}_i^{(d)}$  for  $i = 1, \dots, K$  directly based on  $\hat{\Sigma}_j^{-1} \hat{M}^{(j)}$  for  $j = d, d + 1$ . We did some transformations such that the transformed  $\hat{\Sigma}_j^{-1} \hat{M}^{(d)}$  (with eigenvalues unchanged) is a submatrix of the transformed  $\hat{\Sigma}_j^{-1} \hat{M}^{(d+1)}$ .

Let

$$\hat{\gamma}_{n \times 1} = (\hat{\gamma}_1, \dots, \hat{\gamma}_n)^T = \frac{1}{\hat{\sigma}} \{I - \mathbf{X}_{n \times d} (\mathbf{X}_{n \times d}^T \mathbf{X}_{n \times d})^{-1} \mathbf{X}_{n \times d}^T\} X^{(d+1)},$$

where  $\hat{\sigma}^2$  is the sample variance of  $\{I - \mathbf{X}_{n \times d}(\mathbf{X}_{n \times d}^T \mathbf{X}_{n \times d})^{-1} \mathbf{X}_{n \times d}^T\} X^{(d+1)}$ . Denote  $\tilde{\gamma}_h = n_h^{-1} \sum_{y_i \in S_h} \hat{\gamma}_i$ . Let  $\gamma = X_{d+1} - E(X_{d+1} | X_1, \dots, X_d)$ , and  $\gamma_{n \times 1}$  be the  $n$  regression error terms of the  $n$  observed  $X_{d+1}$  on  $X_1, \dots, X_d$ . Then  $\gamma_{n \times 1}$  are independent and identically distributed with mean 0 and a finite variance. Under the null hypothesis  $H_0: \eta_{d+1,i} = 0, i = 1, \dots, K$ , we have  $E(\gamma | y) = E\{E(\gamma | X_1, \dots, X_d) | y\} = 0$  for any  $y$ . Let  $\bar{\gamma}$  be the mean of  $\gamma_{n \times 1}$ . Then

$$\hat{\gamma}_{n \times 1} = (\hat{\gamma}_1, \dots, \hat{\gamma}_n)^T = \frac{1}{\hat{\sigma}} \{I - \mathbf{X}_{n \times d}(\mathbf{X}_{n \times d}^T \mathbf{X}_{n \times d})^{-1} \mathbf{X}_{n \times d}^T\} (\gamma_{n \times 1} - \bar{\gamma}).$$

With transformations on  $\hat{\Sigma}_j^{-1} \hat{M}^{(d)}$ , we showed that  $\hat{\lambda}_i^{(d+1)} - \hat{\lambda}_i^{(d)}$  for  $i = 1, \dots, K$  equals a squared linear combination of  $\tilde{\gamma}_h$ . Thus, we just need to show that  $(\tilde{\gamma}_1, \dots, \tilde{\gamma}_H)$  converges to a multivariate normal distribution, and we complete the proof. Let  $(z_1, \dots, z_d)' = \Sigma_d^{-1/2}(x_1, \dots, x_d)'$ . Define four matrices,  $A_{H \times H}, B_{H \times d}, E_{d \times d}$  and  $\Gamma_{H \times d}$ , where  $A_{H \times H} = \text{diag}\{\text{var}(\gamma | y \in S_1), \dots, \text{var}(\gamma | y \in S_H)\} / \sigma^2$ , the  $(h, j)$ th entry of  $B_{H \times d}$  is  $\sqrt{p_h} \text{cov}(z_j \gamma, \gamma | y \in S_h) / \sigma^2$ , the  $(j, j')$ th entry of  $E_{d \times d}$  equals  $\text{cov}(z_j \gamma, z_{j'} \gamma) / \sigma^2$ , the  $(h, j)$ th entry of  $\Gamma_{H \times d}$  is  $\sqrt{p_h} E(z_j | y \in S_h)$  and  $\sigma^2 = \lim_{n \rightarrow \infty} (\hat{\sigma}^2) = \text{var}(\gamma)$ . Let  $\Upsilon$  be a  $d \times d$  matrix and

$$\Upsilon = \Gamma_{H \times d}^T A_{H \times H} \Gamma_{H \times d} - \Gamma_{H \times d}^T B_{H \times d} \Gamma_{H \times d}^T \Gamma_{H \times d} - \Gamma_{H \times d}^T \Gamma_{H \times d} B_{H \times d}^T \Gamma_{H \times d} + \Gamma_{H \times d}^T \Gamma_{H \times d} E_{d \times d} \Gamma_{H \times d}^T \Gamma_{H \times d}.$$

Define  $\tilde{Q}$  to be a  $d \times K$  matrix with  $j$ th column  $q_j / \sqrt{\{\lambda_j^{(d)}(1 - \lambda_j^{(d)})\}}$ , where  $q_j$  is the  $j$ th eigenvector of the limiting matrix  $\lim_{n \rightarrow \infty} (\hat{\Sigma}_d^{-1/2} \hat{M}^{(d)} \hat{\Sigma}_d^{-1/2})$ , and  $\lambda_j^{(d)} = \lim_{n \rightarrow \infty} (\hat{\lambda}_j^{(d)})$ . Then  $W_{Kt} = \tilde{Q}^T \Upsilon \tilde{Q}$ .

### A.3. Proof of corollary 1

With an additional condition that  $E(\gamma^2 | X_1, \dots, X_d)$  is constant, we can show that the asymptotic variance matrix of  $(\tilde{\gamma}_1, \dots, \tilde{\gamma}_H)$  adopts a special form, with which the asymptotic standard  $\chi^2$ -distribution can be derived.

### A.4. Proof of theorem 2

Without loss of generality, we let  $\mathcal{A} = \{X_1, \dots, X_d\}$ . Following the notation that was used in theorem 1, let  $\gamma_j = X_j - E(X_j | X_i, i \in \mathcal{A})$  for  $j \in \mathcal{A}^c$ , and

$$\hat{\gamma}^j = (\hat{\gamma}_{j,1}, \dots, \hat{\gamma}_{j,n})' = \frac{1}{\hat{\sigma}_j} \{\mathbf{I}_n - \mathbf{X}_{n \times d}(\mathbf{X}_{n \times d}' \mathbf{X}_{n \times d})^{-1} (\mathbf{X}_{n \times d})'\} X^{(j)}.$$

Let  $\tilde{\gamma}_h^j = n_h^{-1} \sum_{y_i \in S_h} \hat{\gamma}_{j,i}$ . Similarly to the proof of theorem 1, we basically show  $\tilde{\gamma}_h^j$  for  $j = d + 1, \dots, p$  and  $h = 1, \dots, H$  converge to a multivariate normal distribution.

### A.5. Proof of theorem 3

We use the same notation as defined in the proof of theorem 2. Let  $\tilde{\gamma}_h^j = n_h^{-1} \sum_{y_i \in S_h} \hat{\gamma}_{j,i}$ . Let  $\hat{\lambda}_k^{(d)}$  be defined as in the proof of theorem 1. First, for any  $t$ ,  $\text{COP}_{1:K}^{A+t} \geq n (\sum_{k=1}^K \hat{\lambda}_k^{(d+t)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)})$ , and

$$\left| \sum_{k=1}^K \hat{\lambda}_k^{(d+1)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)} \right| \geq \left| \sum_{k=1}^K (\lambda_k^{(d+1)} - \lambda_k^{(d)}) \right| - \left| \sum_{k=1}^K (\lambda_k^{(d+1)} - \hat{\lambda}_k^{(d+1)}) \right| - \left| \sum_{k=1}^K (\lambda_k^{(d)} - \hat{\lambda}_k^{(d)}) \right|.$$

Since  $X$  follows a multivariate normal distribution, from Li (1991),  $\lambda_k^{(d)} = \lambda_k^{(d+1)} = 0$  for  $k > K$ ; then

$$\sum_{k=1}^K (\lambda_k^{(d+1)} - \lambda_k^{(d)}) = \lim_{n \rightarrow \infty} \{\text{tr}(\tilde{\Omega}^{(d+1)}) - \text{tr}(\hat{\Omega}^{(d)})\} = \lim_{n \rightarrow \infty} \left\{ \sum_{h=1}^H \frac{n_h}{n} (\tilde{\gamma}_h^j)^2 \right\} = \sum_{h=1}^H p_h \frac{E^2(\gamma_j | y \in S_h)}{\sigma_j^2}.$$

We need to use the two lemmas 1 and 2 that are stated below. The proofs of the two lemmas are omitted here. From lemma 1,

$$\max_{t \in \mathcal{A}^c \cap \mathcal{T}} \left\{ n \left( \sum_{k=1}^K \hat{\lambda}_k^{(d+1)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)} \right) \right\} \geq \varpi \omega_H n^{1-\xi_0} \frac{\tau_{\min}^2}{\tau_{\max}} - \left| \sum_{k=1}^K n(\lambda_k^{(d+1)} - \hat{\lambda}_k^{(d+1)}) \right| - \left| \sum_{k=1}^K n(\lambda_k^{(d)} - \hat{\lambda}_k^{(d)}) \right|.$$

Then, as long as

$$\max_{\mathcal{A} \subseteq \{1, \dots, p\}} \left\{ n \left| \sum_{k=1}^K (\lambda_k^{(d)} - \hat{\lambda}_k^{(d)}) \right| \right\} \leq \vartheta n^{1-\xi_0} / 2,$$

we have

$$\min_{\mathcal{A}: \mathcal{A}^c \cap \mathcal{T} \neq \emptyset} \max_{t \in \mathcal{A}^c \cap \mathcal{T}} (\text{COP}_{1:K}^{\mathcal{A}+t}) \geq \vartheta n^{1-\xi_0}.$$

From lemma 2,

$$P \left( \max_{\mathcal{A} \subseteq \{1, \dots, p\}} \left| \sum_{k=1}^K \lambda_k^{(d)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)} \right| > \frac{\vartheta n^{-\xi_0}}{2} \right) \leq 2Kp(p+1)C_1 \exp \left( -C_2 n^{1-2\xi_0} \frac{\tau_{\min}^2 \vartheta^2}{256K^2 p^2} \right).$$

Under condition 8, since  $p = o(n^{\vartheta_0})$  with  $2\vartheta_0 + 2\xi_0 < 1$ ,  $P(\max_{\mathcal{A} \subseteq \{1, \dots, p\}} |\sum_{k=1}^K \lambda_k^{(d)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)}| > \vartheta n^{-\xi_0} / 2) \rightarrow 0$ , and  $P\{\min_{\mathcal{A}: \mathcal{A}^c \cap \mathcal{T} \neq \emptyset} \max_{t \in \mathcal{A}^c \cap \mathcal{T}} (\text{COP}_{1:K}^{\mathcal{A}+t}) \geq \vartheta n^{1-\xi_0}\} \rightarrow 1$ .

### A.6. Proof of theorem 4

Since

$$\left| \sum_{k=1}^K \hat{\lambda}_k^{(d+1)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)} \right| \leq \left| \sum_{k=1}^K (\lambda_k^{(d+1)} - \lambda_k^{(d)}) \right| + \left| \sum_{k=1}^K (\lambda_k^{(d+1)} - \hat{\lambda}_k^{(d+1)}) \right| + \left| \sum_{k=1}^K (\lambda_k^{(d)} - \hat{\lambda}_k^{(d)}) \right|,$$

and, with  $\mathcal{T} \subseteq \mathcal{A}$ ,  $|\sum_{k=1}^K (\lambda_k^{(d+1)} - \lambda_k^{(d)})| = 0$ , then, from lemma 2,  $P(\max_{\mathcal{A} \subseteq \{1, \dots, p\}} |\sum_{k=1}^K \lambda_k^{(d)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)}| > \varepsilon) \rightarrow 0$  for  $\varepsilon > Cn^{\vartheta_0-1/2}$  and theorem 4 holds.

*Lemma 1.* Under the same conditions as in theorem 3, for any  $\mathcal{A} \subseteq \{1, \dots, p\}$  and  $\mathcal{A}^c \cap \mathcal{T} \neq \emptyset$ ,

$$\max_{j \in \mathcal{A}^c \cap \mathcal{T}} \left\{ \sum_{h=1}^H p_h E^2(\gamma_j | y \in S_h) / \sigma_j^2 \right\} \geq \tau_{\min}^2 \varpi \omega_H n^{-\xi_0} / \tau_{\max} > 0.$$

*Lemma 2.* Under the same conditions as in lemma 1,

$$P \left( \max_{\mathcal{A} \subseteq \{1, \dots, p\}} \left| \sum_{k=1}^K \lambda_k^{(d)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)} \right| > \varepsilon \right) \leq 2Kp(p+1)C_1 \exp \left( -C_2 n \frac{\tau_{\min}^2 \varepsilon^2}{64K^2 p^2} \right).$$

### A.7. Proof of theorem 5

For coherence, we use the same notation as defined in the proof of theorem 1. Without loss of generality, let  $\mathcal{A} = \{1, \dots, d\}$  and  $t = d + 1$ . Under the assumption that  $\mathbf{X}_{n \times (d+1)}$  has a multivariate normal distribution, we derive the limiting value of  $(\sum_{k=1}^K \hat{\lambda}_k^{(d+1)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)})$  as  $n \rightarrow \infty$  for fixed slices. Let  $\Xi_{K \times K}$  be the variance-covariance matrix of  $\mathbf{v}_K$ .

Because  $\{X_1, \dots, X_{d+1}\}$  follow a multivariate normal distribution, we have  $\gamma = X_{d+1} - \rho_0 + \sum_{i=1}^d \rho_i X_i$  and  $\gamma \sim N(0, \sigma_{d+1}^2)$  where the  $\rho_i$  are the coefficients. Since we assume that the response depends only on  $K$  linear combinations of  $\mathbf{X}_{n \times (d+1)}$ ,  $\tilde{\Omega}^{(d+1)}$  and  $\hat{\Omega}^{(d)}$  have at most  $K$  non-zero eigenvalues, and

$$\begin{aligned} \frac{\sum_{k=1}^K \hat{\lambda}_k^{(d+1)}}{\text{tr}(\tilde{\Omega}^{(d+1)})} &\xrightarrow{P} 1, \\ \frac{\sum_{k=1}^K \hat{\lambda}_k^{(d)}}{\text{tr}(\hat{\Omega}^{(d)})} &\xrightarrow{P} 1, \\ \frac{\sum_{k=1}^K \hat{\lambda}_k^{(d+1)} - \sum_{k=1}^K \hat{\lambda}_k^{(d)}}{\text{tr}(\tilde{\Omega}^{(d+1)}) - \text{tr}(\hat{\Omega}^{(d)})} &\xrightarrow{P} 1. \end{aligned}$$

We have the following three results.

- (a)  $\text{tr}(\hat{\Omega}^{(d+1)}) - \text{tr}(\hat{\Omega}^{(d)}) = \sum_{h=1}^H n_h (\tilde{\gamma}_h)^2 / n.$
- (b)  $\tilde{\gamma}_h \xrightarrow{P} E(\gamma|y \in S_h) / \sigma_{d+1}, h = 1, \dots, H.$
- (c) Since  $E(\gamma|\mathbf{v}_K) = \tilde{\eta}'_{t,\mathcal{A}} \Xi_{K \times K}^{-1} \mathbf{v}_K,$  then

$$E(\gamma|y \in S_h) = E\{E(\gamma|\mathbf{v}_K)|y \in S_h\} = \tilde{\eta}'_{t,\mathcal{A}} \Xi_{K \times K}^{-1} \mathbf{L}_{H,K}.$$

Combining results (a)–(c) we have  $\sum_{h=1}^H n_h (\tilde{\gamma}_h)^2 / n \rightarrow^P \tilde{\eta}'_{t,\mathcal{A}} \Xi_{K \times K}^{-1} \mathbf{M}_{H,K} \Xi_{K \times K}^{-1} \tilde{\eta}_{t,\mathcal{A}}.$  Since  $\Xi_{K \times K}^{-1} = \mathbf{I}_{K \times K},$

$$\sum_{k=1}^K \hat{\lambda}_k^{A+t} - \sum_{k=1}^K \hat{\lambda}_k^A \xrightarrow{P} \frac{1}{\sigma_{d+1}^2} \tilde{\eta}'_{t,\mathcal{A}} \mathbf{M}_{H,K} \tilde{\eta}_{t,\mathcal{A}},$$

and theorem 5 holds.

### A.8. Proof of proposition 2

Note that  $\eta' M_{H,K} \eta = \text{var}\{E(\eta' \mathbf{v}_K | y \in S_h)\}$  and

$$\text{var}\{E(\eta' \mathbf{v}_K | y \in S'_h)\} = \text{var}\{E(\eta' \mathbf{v}_K | y \in S_h)\} + \text{var}[\text{var}\{E(\eta' \mathbf{v}_K | y \in S'_h) | S_h\}].$$

Thus, proposition 2 holds.

## References

Bondell, H. D. and Li, L. (2009) Shrinkage inverse regression estimation for model-free variable selection. *J. R. Statist. Soc. B*, **71**, 287–299.

Boyer, L. A., Lee, T. I., Cole, M. F., Johnstone, S. E., Levine, S. S., Zucker, J. P., Guenther, M. G., Kumar, R. M., Murray, H. L., Jenner, R. G., Gifford, D. K., Melton, D. A., Jaenisch, R. and Young, R. A. (2005) Core transcriptional regulatory circuitry in human embryonic stem cells. *Cell*, **122**, 947–956.

Cai, J., Xie, D., Fan, Z., Chipperfield, H., Marden, J., Wong, W. H. and Zhong, S. (2010) Modeling co-expression across species for complex traits: insights to the difference of human and mouse embryonic stem cells. *PLOS Comput. Biol.*, **6**, article e1000707.

Chen, C.-H. and Li, K.-C. (1998) Can SIR be as popular as multiple linear regression? *Statist. Sin.*, **8**, 289–316.

Chen, X., Xu, H., Yuan, P., Fang, F., Huss, M., Vega, V. B., Wong, E., Orlov, Y. L., Zhang, W., Jiang, J., Loh, Y. H., Yeo, H. C., Yeo, Z. X., Narang, V., Govindarajan, K. R., Leong, B., Shahab, A., Ruan, Y., Bourque, G., Sung, W. K., Clarke, N. D., Wei, C. L. and Ng, H. H. (2008) Integration of external signaling pathways with the core transcriptional network in embryonic stem cells. *Cell*, **133**, 1106–1117.

Cleveland, W. (1979) Robust locally weighted regression and smoothing scatterplots. *J. Am. Statist. Ass.*, **74**, 829–836.

Cloonan, N., Forrest, A. R. R., Kolle, G., Gardiner, B. B. A., Faulkner, G. J., Brown, M. K., Taylor, D. F., Steptoe, A. L., Wani, S., Bethel, G., Robertson, A. J., Perkins, A. C., Bruce, S. J., Lee, C. C., Ranade, S. S., Peckham, H. E., Manning, J. M., McKernan, K. J. and Grimmond, S. M. (2008) Stem cell transcriptome profiling via massive-scale mRNA sequencing. *Nat. Meth.*, **5**, 613–619.

Conlon, E., Liu, X., Lieb, J. and Liu, J. (2003) Integrating regulatory motif discovery and genome-wide expression analysis. *Proc. Natn. Acad. Sci. USA*, **100**, 3339–3344.

Cook, R. D. (1994) *An Introduction to Regression Graphics*. New York: Wiley.

Cook, R. (2004) Testing predictor contributions in sufficient dimension reduction. *Ann. Statist.*, **32**, 1062–1092.

Cook, R. D. and Nachtsheim, C. J. (1994) Reweighting to achieve elliptically contoured covariates in regression. *J. Am. Statist. Ass.*, **89**, 592–599.

Eaton, M. L. (1986) A characterization of spherical distributions. *J. Multiv. Anal.*, **20**, 272–276.

Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R. (2004) Least angle regression. *Ann. Statist.*, **32**, 407–499.

Fan, J. and Li, R. (2001) Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Am. Statist. Ass.*, **96**, 1348–1360.

Fan, J. and Lv, J. (2008) Sure independence screening for ultrahigh dimensional feature space (with discussion). *J. R. Statist. Soc. B*, **70**, 849–911.

Fan, J. and Lv, J. (2010) A selective overview of variable selection in high dimensional feature space. *Statist. Sin.*, **20**, 101–148.

Friedman, J. H. (2007) Pathwise coordinate optimization. *Ann. Appl. Statist.*, **1**, 302–332.

Friedman, J. H. and Tukey, J. W. (1974) A projection pursuit algorithm for explanatory data analysis. *IEEE Trans. Comput. C*, **23**, 881–889.

Fung, W., He, X., Liu, L. and Shi, P. (2002) Dimension reduction based on canonical correlation. *Statist. Sin.*, **12**, 1093–1114.

- Hall, P. and Li, K.-C. (1993) On almost linearity of low dimensional projections from high dimensional data. *Ann. Statist.*, **21**, 867–889.
- Huber, P. J. (1985) Projection pursuit. *Ann. Statist.*, **13**, 435–475.
- Johnson, D. S., Mortazavi, A., Myers, R. M. and Wold, B. (2007) Genome-wide mapping of in vivo protein-DNA interactions. *Science*, **316**, 1497–1502.
- Li, K.-C. (1991) Sliced inverse regression for dimension reduction. *J. Am. Statist. Ass.*, **86**, 316–327.
- Li, L. (2007) Sparse sufficient dimension reduction. *Biometrika*, **94**, 603–613.
- Li, L., Cook, R. D. and Nachtshiem, C. J. (2005) Model-free variable selection. *J. R. Statist. Soc. B*, **67**, 285–299.
- Matsuda, T., Nakamura, T., Nakao, K., Arai, T., Katsuki, M., Heike, T. and Yokota, T. (1999) STAT3 activation is sufficient to maintain an undifferentiated state of mouse embryonic stem cells. *EMBO J.*, **18**, 4261–4269.
- Miller, A. J. (1984) Selection of subsets of regression variables (with discussion). *J. R. Statist. Soc. A*, **147**, 389–425.
- Mortazavi, A., Williams, B. A., McCue, K., Schaeffer, L. and Wold, B. (2008) Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat. Meth.*, **5**, 621–628.
- Nagalakshmi, U., Wang, Z., Waern, K., Shou, C., Raha, D., Gerstein, M. and Snyder, M. (2008) The transcriptional landscape of the yeast genome defined by RNA sequencing. *Science*, **320**, 1344–1349.
- Ouyang, Z., Zhou, Q. and Wong, W. H. (2009) Chip-seq of transcription factors predicts absolute and differential gene expression in embryonic stem cells. *Proc. Natn. Acad. Sci. USA*, **106**, 21521–21526.
- Pardal, R., Molofsky, A. V., He, S. and Morrison, S. J. (2005) Stem cell self-renewal and cancer cell proliferation are regulated by common networks that balance the activation of proto-oncogenes and tumor suppressors. *Cold Spring Harb. Symp. Quant. Biol.*, **70**, 177–185.
- Shao, J. (1998) An asymptotic theory for linear model selection (with discussion). *Statist. Sin.*, **7**, 221–264.
- Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *J. R. Statist. Soc. B*, **58**, 267–288.
- Wang, H. (2009) Forward regression for ultra-high dimensional variable screening. *J. Am. Statist. Ass.*, **104**, 1512–1524.
- Wilhelm, B. T., Marguerat, S., Watt, S., Schubert, F., Wood, V., Goodhead, I., Penkett, C. J., Rogers, J. and Bahler, J. (2008) Dynamic repertoire of a eukaryotic transcriptome surveyed at single-nucleotide resolution. *Nature*, **453**, 1239–1243.
- Zeng, P. and Zhu, Y. (2010) An integral transform method for estimating the central mean and central subspaces. *J. Multiv. Anal.*, **101**, 271–290.
- Zhong, W., Zeng, P., Ma, P., Liu, J. and Zhu, Y. (2005) Regularized sliced inverse regression for motif discovery. *Bioinformatics*, **21**, 4169–4175.
- Zhou, J. and He, X. (2008) Dimension reduction based on constrained canonical correlation and variable filtering. *Ann. Statist.*, **36**, 1649–1668.
- Zhu, L., Miao, B. and Peng, H. (2006) On sliced inverse regression with high-dimensional covariates. *J. Am. Statist. Ass.*, **101**, 630–643.
- Zou, H. (2006) The adaptive lasso and its oracle properties. *J. Am. Statist. Ass.*, **101**, 1418–1429.