



# Sparse principal component regression with adaptive loading



Shuichi Kawano<sup>a,d,\*</sup>, Hironori Fujisawa<sup>b,d</sup>, Toyoyuki Takada<sup>c,d</sup>,  
Toshihiko Shiroishi<sup>c,d</sup>

<sup>a</sup> Graduate School of Information Systems, University of Electro-Communications, 1-5-1 Chofugaoka, Chofu-shi, Tokyo 182-8585, Japan

<sup>b</sup> The Institute of Statistical Mathematics, 10-3 Midori-cho, Tachikawa, Tokyo 190-8562, Japan

<sup>c</sup> Mammalian Genetics Laboratory, National Institute of Genetics, Mishima, Shizuoka 411-8540, Japan

<sup>d</sup> Transdisciplinary Research Integration Center, Research Organization of Information and Systems, Minato-ku, Tokyo 105-0001, Japan

## ARTICLE INFO

### Article history:

Received 24 July 2014

Received in revised form 7 February 2015

Accepted 25 March 2015

Available online 7 April 2015

### Keywords:

Dimension reduction

Identifiability

Principal component regression

Regularization

Sparsity

## ABSTRACT

Principal component regression (PCR) is a two-stage procedure that selects some principal components and then constructs a regression model regarding them as new explanatory variables. Note that the principal components are obtained from only explanatory variables and not considered with the response variable. To address this problem, we propose the sparse principal component regression (SPCR) that is a one-stage procedure for PCR. SPCR enables us to adaptively obtain sparse principal component loadings that are related to the response variable and select the number of principal components simultaneously. SPCR can be obtained by the convex optimization problem for each parameter with the coordinate descent algorithm. Monte Carlo simulations and real data analyses are performed to illustrate the effectiveness of SPCR.

© 2015 Elsevier B.V. All rights reserved.

## 1. Introduction

Principal component analysis (PCA) (Pearson, 1901) is a fundamental statistical tool for dimensionality reduction, data processing, and visualization of multivariate data, with various applications in biology, engineering, and social science. In regression analysis, it can be useful to replace many original explanatory variables with a few principal components, which is called the principal component regression (PCR) (Massy, 1965; Jolliffe, 1982). PCR is widely used in various fields of research and many extensions of PCR have been proposed (see, e.g., Hartnett et al., 1998; Rosipal et al., 2001; Reiss and Ogden, 2007; Wang and Abbott, 2008 and Han and Liu, 2013). Whereas PCR is a useful tool for analyzing multivariate data, this method may not have enough prediction accuracy if the response variable depends on the principal components with small eigenvalues. The problem arises from the two-stage procedure for PCR; a few principal components are selected with large eigenvalues, but without any relation to response variable, and then the regression model is constructed using them as new explanatory variables.

In this paper, we deal with PCA and regression analysis simultaneously, and propose a one-stage procedure for PCR to address this problem. The procedure combines two loss functions; one is the ordinary regression analysis loss and the other is PCA loss with some devices proposed by Zou et al. (2006). In addition, in order to easily interpret estimated principal

\* Corresponding author at: Graduate School of Information Systems, University of Electro-Communications, 1-5-1 Chofugaoka, Chofu-shi, Tokyo 182-8585, Japan. Tel.: +81 42 443 5620.

E-mail addresses: [skawano@ai.is.uec.ac.jp](mailto:skawano@ai.is.uec.ac.jp) (S. Kawano), [fujisawa@ism.ac.jp](mailto:fujisawa@ism.ac.jp) (H. Fujisawa), [ttakada@nig.ac.jp](mailto:ttakada@nig.ac.jp) (T. Takada), [tshirois@nig.ac.jp](mailto:tshirois@nig.ac.jp) (T. Shiroishi).

component loadings and select the number of principal components automatically, we impose the  $L_1$  type regularization on the parameters. This one-stage procedure is called the sparse principal component regression (SPCR) in this paper. SPCR gives sparse principal component loadings that are related to the response variable and selects the number of principal components simultaneously. We also establish a monotonically decreasing estimation procedure for the loss function using the coordinate descent algorithm (Friedman et al., 2010), because SPCR can be obtained via the convex optimization problem for each of parameters.

The partial least squares regression (PLS) (Wold, 1975; Frank and Friedman, 1993) is a dimension reduction technique, which incorporates information between the explanatory variables and the response variable. Recently, Chun and Keleş (2010) have proposed the sparse partial least squares regression (SPLS) that imposes sparsity in the dimension reduction step of PLS, and then constructed a regression model regarding some SPLS components as new explanatory variables, although it is a two-stage procedure. Besides PLS and SPLS, several methods have been proposed for performing dimension reduction and regression analysis simultaneously. Bair et al. (2006) proposed the supervised principal component analysis, which is regression analysis in which the explanatory variables are related to the response variable with respect to correlation. Yu et al. (2006) presented the supervised probabilistic principal component analysis from the Bayesian viewpoint. By imposing the  $L_1$  type regularization into the objective function, Allen et al. (2013) and Chen and Huang (2012) introduced the regularized partial least squares and the sparse reduced-rank regression, respectively. However, none of them integrated the two loss functions for ordinary regression analysis and PCA along with the  $L_1$  type regularization.

Here, we present a characteristic of SPCR with comparison to other sparse methods. Many sparse regression methods are based on the original explanatory variables, but SPCR and SPLS can use new explanatory variables by loadings of the original explanatory variables, which may include more useful explanatory variables, like in PCA. SPCR is obtained from the ordinary regression analysis loss, but SPLS is derived from the covariance loss, which not directly evaluates the ordinary regression analysis loss. This difference is discussed in more detail in Section 3.3.

This paper is organized as follows. In Section 2, we review PCA and the sparse principal component analysis (SPCA) by Zou et al. (2006). We propose SPCR and discuss alternative methods to SPCR in Section 3. Section 4 provides an efficient algorithm for SPCR and a method for selecting tuning parameters in SPCR. Monte Carlo simulations and real data analyses are provided in Section 5. Concluding remarks are given in Section 6. The R language software package `spcr`, which implements SPCR, is available on the Comprehensive R Archive Network (<http://cran.r-project.org>). Supplementary materials can be found in [https://sites.google.com/site/shuichikawanoen/research/suppl\\_spcr.pdf](https://sites.google.com/site/shuichikawanoen/research/suppl_spcr.pdf).

## 2. Preliminaries

### 2.1. Principal component analysis

Let  $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$  be an  $n \times p$  data matrix, where  $n$  and  $p$  denote the sample size and the number of variables, respectively. Without loss of generality, we assume that the column means of the matrix  $X$  are all zero.

PCA is usually implemented by using the singular value decomposition (SVD) of  $X$ . When the SVD of  $X$  is represented by

$$X = UDV^T,$$

the principal components are  $Z = UD$  and the corresponding loadings of the principal components are the columns of  $V$ . Here,  $U$  is an  $n \times n$  orthogonal matrix,  $V = (\mathbf{v}_1, \dots, \mathbf{v}_p)$  is a  $p \times p$  orthogonal matrix, and  $D$  is an  $n \times p$  matrix given by

$$D = \begin{pmatrix} D^* & O_{q,p-q} \\ O_{n-q,q} & O_{n-q,p-q} \end{pmatrix},$$

where  $q = \text{rank}(X)$ ,  $D^* = \text{diag}(d_1, \dots, d_q)$  ( $d_1 \geq \dots \geq d_q > 0$ ), and  $O_{i,j}$  is the  $i \times j$  matrix with all zero elements. Note that the vectors  $V^T \mathbf{x}_1, \dots, V^T \mathbf{x}_n$  are also the principal components, since  $XV = Z$ .

The loading matrix can be obtained by solving the following least squares problem (see, e.g., Hastie et al., 2009);

$$\begin{aligned} \min_B \sum_{i=1}^n \|\mathbf{x}_i - BB^T \mathbf{x}_i\|^2 \\ \text{subject to } B^T B = I_k, \end{aligned} \quad (1)$$

where  $B = (\beta_1, \dots, \beta_k)$  is a  $p \times k$  loading matrix,  $k$  denotes the number of principal components, and  $I_k$  is the  $k \times k$  identity matrix.

The solution is given by

$$\hat{B} = V_k Q^T,$$

where  $V_k = (\mathbf{v}_1, \dots, \mathbf{v}_k)$  and  $Q$  is a  $k \times k$  arbitrary orthogonal matrix.

## 2.2. Sparse principal component analysis

Zou et al. (2006) proposed an alternative least squares problem given by

$$\min_{A,B} \left\{ \sum_{i=1}^n \|\mathbf{x}_i - AB^T \mathbf{x}_i\|^2 + \lambda \sum_{j=1}^k \|\boldsymbol{\beta}_j\|^2 \right\} \quad (2)$$

subject to  $A^T A = I_k$ ,

where  $A = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_k)$  is a  $p \times k$  matrix and  $\lambda (> 0)$  is a regularization parameter. The minimizer of  $B$  is given by

$$\hat{B} = V_k C Q^T, \quad (3)$$

where  $C = \text{diag}(c_1, \dots, c_k)$ ,  $c_i$  ( $i = 1, \dots, k$ ) is a positive constant, and  $Q$  is an arbitrary orthogonal matrix. The case  $\lambda = 0$  yields the same solution as (1). Formula (2) is a quadratic programming problem with respect to each parameter matrix  $A$  and  $B$ , but Formula (1) is not.

In addition, Zou et al. (2006) proposed to add a sparse regularization term for  $B$  to easily interpret the estimate  $\hat{B}$ , which is called SPCA;

$$\min_{A,B} \left\{ \sum_{i=1}^n \|\mathbf{x}_i - AB^T \mathbf{x}_i\|^2 + \lambda \sum_{j=1}^k \|\boldsymbol{\beta}_j\|^2 + \sum_{j=1}^k \lambda_{1,j} \|\boldsymbol{\beta}_j\|_1 \right\} \quad (4)$$

subject to  $A^T A = I_k$ ,

where  $\lambda_{1,j}$ 's ( $j = 1, \dots, k$ ) are regularization parameters with positive value and  $\|\cdot\|_1$  is the  $L_1$  norm of  $\boldsymbol{\beta}$ . Note that the minimization problem (4) is also the quadratic programming problem with respect to each parameter matrix  $A$  and  $B$ . After simple calculation, the problem (4) becomes

$$\min_{A,B} \sum_{j=1}^k \left\{ \|X\boldsymbol{\alpha}_j - X\boldsymbol{\beta}_j\|^2 + \lambda \|\boldsymbol{\beta}_j\|^2 + \lambda_{1,j} \|\boldsymbol{\beta}_j\|_1 \right\}$$

subject to  $A^T A = I_k$ .

This optimization problem is analogous to the elastic net problem in Zou and Hastie (2005), and hence Zou et al. (2006) proposed an alternating algorithm to estimate  $A$  and  $B$  iteratively. In particular, the LARS algorithm (Efron et al., 2004) is employed to obtain the estimate of  $B$  numerically.

Another approach to obtain a sparse loading matrix is SCoTLASS (Jolliffe et al., 2003). However, Zou et al. (2006) pointed out that the loadings obtained by SCoTLASS are not sparse enough. Also, Lee et al. (2010) and Lee and Huang (2013) developed SPCA for binary data.

## 3. Sparse principal component regression

### 3.1. Sparse principal component regression with adaptive loading

Suppose that we have data for response variables  $y_1, \dots, y_n$  in addition to data  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . We consider regression analysis in the situation that the response variable is explained by variables aggregated by PCA of  $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ . A naive approach is to construct a regression model with a few principal components corresponding to large eigenvalues, which are previously constructed. This approach is called PCR. However, PCR might fail to predict the response if the response is associated with principal components corresponding to small eigenvalues.

To overcome this drawback, we propose SPCR using the principal components  $B^T \mathbf{x}$  as follows:

$$\min_{A,B,\gamma_0,\boldsymbol{\gamma}} \left\{ (1-w) \sum_{i=1}^n (y_i - \gamma_0 - \boldsymbol{\gamma}^T B^T \mathbf{x}_i)^2 + w \sum_{i=1}^n \|\mathbf{x}_i - AB^T \mathbf{x}_i\|^2 \right. \\ \left. + \lambda_\beta (1-\zeta) \sum_{j=1}^k \|\boldsymbol{\beta}_j\|_1 + \lambda_\beta \zeta \sum_{j=1}^k \|\boldsymbol{\beta}_j\|^2 + \lambda_\gamma \|\boldsymbol{\gamma}\|_1 \right\} \quad (5)$$

subject to  $A^T A = I_k$ ,

where  $\gamma_0$  is an intercept,  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_k)^T$  is a coefficient vector,  $\lambda_\beta$  and  $\lambda_\gamma$  are regularization parameters with positive value, and  $w$  and  $\zeta$  are tuning parameters whose values are between zero and one.

The first term in Formula (5) means the least squares loss between the response and the principal components  $B^T \mathbf{x}$ . The second term induces PCA loss of data  $X$ . The tuning parameter  $w$  controls the trade-off between the first and second terms,

and then the value of  $w$  can be determined by a user for any purpose. For example, a smaller value for  $w$  is used when we aim to obtain better prediction accuracies, while a larger value for  $w$  is used when we aim to obtain the exact formulation of the principal component loadings. The third and fifth terms encourage sparsity on  $B$  and  $\gamma$ , respectively. The sparsity on  $B$  enables us to easily interpret the loadings of the principal components. Meanwhile, the sparsity on  $\gamma$  induces automatic selection of the number of principal components. The tuning parameter  $\zeta$  controls the trade-off between the  $L_1$  and  $L_2$  norms for the parameter  $B$ , which was introduced in [Zou and Hastie \(2005\)](#). For detailed roles of this parameter and the  $L_2$  norm, see [Zou and Hastie \(2005\)](#).

We see that (5) is a quadratic programming problem with respect to each parameter, because the problem only combines a regression loss with PCA loss. The optimization problem appears to be simple. However, it is not easy to numerically obtain the estimates of the parameters if we do not introduce the  $L_1$  regularization terms for  $B$  and  $\gamma$ , because there exists an identification problem for  $B$  and  $\gamma$ . For an arbitrary orthogonal matrix  $P$ , we have

$$\gamma^T B^T = \gamma^T P^T P B^T = \gamma^{\dagger T} B^{\dagger T},$$

where  $\gamma^{\dagger} = P\gamma$  and  $B^{\dagger} = B P^T$ . This causes non-unique estimators for  $B$  and  $\gamma$ . However, we incorporate the  $L_1$ -penalties on (5) and then we can expect to obtain the minimizer, because the parameter exists on a hypersphere due to orthogonal invariance in (3) and the  $L_1$ -penalty implies a hypersquare region. For more details, see, e.g., [Tibshirani \(1996\)](#), [Jennrich \(2006\)](#), [Choi et al. \(2011\)](#), and [Hirose and Yamamoto \(in press\)](#). The  $L_1$ -penalties on  $B$  and  $\gamma$  play two types of roles on sparsity and identification problem.

### 3.2. Adaptive sparse principal component regression

In the numerical study in Section 5, we observe that SPCR does not produce enough sparse solution for the loading matrix  $B$ . We, therefore, assign different weights to different parameters in the loading matrix  $B$ . This idea was adopted in the adaptive lasso ([Zou, 2006](#)). Let us consider the weighted sparse principal component regression, given by

$$\begin{aligned} \min_{A, B, \gamma_0, \gamma} \left\{ (1-w) \sum_{i=1}^n (y_i - \gamma_0 - \gamma^T B^T \mathbf{x}_i)^2 + w \sum_{i=1}^n \|\mathbf{x}_i - A B^T \mathbf{x}_i\|^2 \right. \\ \left. + \lambda_{\beta} (1-\zeta) \sum_{j=1}^k \sum_{l=1}^p \omega_{lj} |\beta_{lj}| + \lambda_{\beta} \zeta \sum_{j=1}^k \|\beta_j\|^2 + \lambda_{\gamma} \|\gamma\|_1 \right\} \\ \text{subject to } A^T A = I_k, \end{aligned}$$

where  $\omega_{lj} (> 0)$  is an incorporated weight for the parameter  $\beta_{lj}$ . We call this procedure the adaptive sparse principal component regression (aSPCR). In this paper, we define the weight as  $\omega_{lj} = 1/|\hat{\beta}_{lj}(\text{SPCR})|$ , where  $\hat{\beta}_{lj}(\text{SPCR})$  is an estimate of the parameter  $\beta_{lj}$  obtained from SPCR. In the adaptive lasso, the weight is constructed using the least squares estimators, but it is not applicable due to the identification problem, as described in Section 3.1.

Since aSPCR is a quadratic programming problem with respect to each parameter, we can estimate the parameters according to an efficient estimation algorithm for SPCR. In addition, aSPCR enjoys properties similar to SPCR as described in Section 3.1.

### 3.3. Related work

PLS (see, e.g., [Wold, 1975](#) and [Frank and Friedman, 1993](#)) seeks directions that relate  $X$  to  $\mathbf{y}$  and capture the most variable directions in the  $X$ -space, which is, in general, formulated by

$$\begin{aligned} \mathbf{w}_k = \arg \max_{\mathbf{w}} [\text{Corr}^2(\mathbf{y}, X\mathbf{w}) \text{Var}(X\mathbf{w})] \\ \text{subject to } \mathbf{w}^T \mathbf{w} = 1, \quad \mathbf{w}^T \Sigma_{XX} \mathbf{w}_j = 0, \quad j = 1, \dots, k-1 \end{aligned} \quad (6)$$

for  $k = 1, \dots, p$ , where  $\mathbf{y} = (y_1, \dots, y_n)^T$  and  $\Sigma_{XX}$  is the covariance matrix of  $X$ . The solutions in the problem (6) are derived from NIPALS ([Wold, 1975](#)) or SIMPLS ([de Jong, 1993](#)).

To incorporate sparsity into PLS, SPLS was introduced by [Chun and Keleş \(2010\)](#). The first SPLS direction vector  $\mathbf{c}$  is obtained by

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{c}} \{ -\kappa \mathbf{w}^T M \mathbf{w} + (1-\kappa)(\mathbf{c} - \mathbf{w})^T M (\mathbf{c} - \mathbf{w}) + \lambda_{1, \text{SPLS}} \|\mathbf{c}\|_1 + \lambda_{2, \text{SPLS}} \|\mathbf{c}\|^2 \} \\ \text{subject to } \mathbf{w}^T \mathbf{w} = 1, \end{aligned} \quad (7)$$

where  $M = X^T \mathbf{y} \mathbf{y}^T X$ , and  $\kappa, \lambda_{1, \text{SPLS}}, \lambda_{2, \text{SPLS}}$  are tuning parameters with positive value. Note that the problem (7) becomes the original maximum eigenvalue problem of PLS when  $\kappa = 1, \lambda_{1, \text{SPLS}} = 0$ , and  $\lambda_{2, \text{SPLS}} = 0$ . This SPLS problem is solved by alternately estimating the parameters  $\mathbf{w}$  and  $\mathbf{c}$ . The idea is similar to that used in SPCA. [Chun and Keleş \(2010\)](#) furthermore

introduced the SPLS–NIPALS and SPLS–SIMPLS algorithm for deriving the rest of the direction vectors, and then predicted the response variable by a linear model with SPLS loading vectors as new explanatory variables; it is a two-stage procedure.

To emphasize a difference between our proposed method and the related work described above, we consider an example as follows. Suppose that

$$y = a_1x_1 + a_2x_2 + \varepsilon, \quad x_j \sim N(0, \tau_j^2), \quad \varepsilon \sim N(0, \sigma^2).$$

This model has another expression in the form

$$y = a_1^*z_1 + a_2^*z_2 + \varepsilon, \quad z_j \sim N(0, 1), \quad a_j^* = a_j\tau_j.$$

The covariance structures are given by

$$\text{Cov}(y, x_j) = a_j\tau_j^2, \quad \text{Cov}(y, z_j) = a_j^* = a_j\tau_j.$$

Let us select the explanatory variable that maximizes the covariance:

$$\max_x \text{Cov}(y, x) \quad \text{or} \quad \max_z \text{Cov}(y, z) = \max_z \text{Corr}(y, z).$$

Consider the case  $(a_1, a_2, \tau_1, \tau_2) = (8, 1, 1, 3)$ . It follows that  $a_1^* = 8$ ,  $a_2^* = 3$ ,  $a_1\tau_1^2 = 8$  and  $a_2\tau_2^2 = 9$ . In this case, it is clear that the first variable  $(x_1, z_1)$  has a larger effect in  $y$  than the second variable  $(x_2, z_2)$ . Remember that PLS and SPLS are based on the maximization of covariance, so that they will firstly select the variable  $z_1$  on the second maximization, whereas they will firstly select the variable  $x_2$  on the first maximization. Therefore, on the first maximization, PLS and SPLS fail to select the explanatory variable largely associated with the response. Meanwhile, SPCR will select the first variable  $(x_1, z_1)$  on both maximizations, because the prediction error remains unchanged after normalization.

## 4. Implementation

### 4.1. Computational algorithm

For estimating the parameter  $A$ , we utilize the same algorithm given by Zou et al. (2006). The parameters  $B$  and  $\gamma$  are estimated by the coordinate descent algorithm (Friedman et al., 2010), because the optimization problems include the  $L_1$  regularization terms, respectively.

The optimization problem in aSPCR is rewritten as follows:

$$\begin{aligned} \min_{A, B, \gamma_0, \gamma} & \left[ (1-w) \sum_{i=1}^n \left\{ y_i - \gamma_0 - \sum_{j=1}^k \gamma_j \left( \sum_{l=1}^p \beta_{lj} x_{il} \right) \right\}^2 + w \sum_{j=1}^k \sum_{i=1}^n \left( y_{ji}^* - \sum_{l=1}^p \beta_{lj} x_{il} \right)^2 \right. \\ & \left. + \lambda_\beta (1-\zeta) \sum_{j=1}^k \sum_{l=1}^p \omega_{lj} |\beta_{lj}| + \lambda_\beta \zeta \sum_{j=1}^k \sum_{l=1}^p \beta_{lj}^2 + \lambda_\gamma \sum_{j=1}^k |\gamma_j| \right] \end{aligned}$$

subject to  $A^T A = I_k$ ,

where  $y_{ji}^*$  is the  $i$ th element of the vector  $X\alpha_j$ . SPCR is a special case of aSPCR with  $\omega_{lj} = 1$ . The detailed algorithm is given as follows.

$\beta_{lj}$  given  $\gamma_0, \gamma_j$  and  $A$ : The coordinate-wise update for  $\beta_{lj}$  has the form:

$$\hat{\beta}_{l'j'} \leftarrow \frac{S \left( \sum_{i=1}^n x_{il'} \left\{ (1-w) Y_i \gamma_{j'} + Y_{ji}^* w \right\}, \frac{\lambda_\beta \omega_{l'j'} (1-\zeta)}{2} \right)}{\left\{ (1-w) \gamma_{j'}^2 + w \right\} \sum_{i=1}^n x_{il'}^2 + \lambda_\beta \zeta}, \quad (l' = 1, \dots, p; j' = 1, \dots, k), \quad (8)$$

where

$$\begin{aligned} Y_i &= y_i - \gamma_0 - \sum_{j=1}^k \sum_{l \neq l'} \gamma_j \beta_{lj} x_{il} - \sum_{j \neq j'} \gamma_j \beta_{l'j} x_{il'}, \\ Y_{ji}^* &= y_{ji}^* - \sum_{l \neq l'} \beta_{lj'} x_{il}, \end{aligned}$$

and  $S(z, \eta)$  is the soft-thresholding operator with value

$$\text{sign}(z)(|z| - \eta)_+ = \begin{cases} z - \eta & (z > 0 \text{ and } \eta < |z|) \\ z + \eta & (z < 0 \text{ and } \eta < |z|) \\ 0 & (\eta \geq |z|). \end{cases}$$

$\gamma_j$  given  $\gamma_0$ ,  $\beta_{ij}$  and  $A$ : The update expression for  $\gamma_j$  is given by

$$\hat{\gamma}_{j'} \leftarrow \frac{S \left( (1-w) \sum_{i=1}^n y_i^{**} x_{ij'}^*, \frac{\lambda_{j'}}{2} \right)}{(1-w) \sum_{i=1}^n x_{ij'}^{*2}}, \quad (j' = 1, \dots, k), \quad (9)$$

where

$$x_{ij}^* = \beta_j^T \mathbf{x}_i, \\ y_i^{**} = y_i - \gamma_0 - \sum_{j \neq j'} \gamma_j x_{ij}^*.$$

$A$  given  $\gamma_0$ ,  $\beta_{ij}$  and  $\gamma_j$ : The estimate of  $A$  is obtained by

$$\hat{A} = UV^T,$$

where  $(X^T X)B = UDV^T$ .

$\gamma_0$  given  $\beta_{ij}$ ,  $\gamma_j$  and  $A$ : The estimate of  $\gamma_0$  is derived from

$$\hat{\gamma}_0 = \frac{1}{n} \sum_{i=1}^n \left\{ y_i - \sum_{j=1}^k \hat{\gamma}_j \left( \sum_{l=1}^p \hat{\beta}_{jl} x_{il} \right) \right\}.$$

This procedure is iterated until a convergence condition is satisfied.

#### 4.2. More efficient algorithm

To speed up our algorithm, we apply the covariance updates, which was proposed by [Friedman et al. \(2010\)](#), into the parameter updates.

We can rewrite the update of the parameter  $B$  in (8) in the form

$$\sum_{i=1}^n x_{il'} \left\{ (1-w) Y_i \gamma_{j'} + Y_{ji}^* w \right\} = (1-w) \gamma_{j'} \sum_{i=1}^n x_{il'} r_i + w \sum_{i=1}^n x_{il'} r_{ji}^* + \tilde{\beta}_{lj'} \sum_{i=1}^n x_{il'}^2 \left\{ (1-w) \gamma_{j'}^2 + w \right\},$$

where  $\tilde{\beta}_{lj'}$  is the current estimate of  $\beta_{lj'}$ ,  $r_i = y_i - \gamma_0 - \sum_{j=1}^k \sum_{l=1}^p \gamma_j \tilde{\beta}_{lj} x_{il}$  and  $r_{ji}^* = y_{ji}^* - \sum_{l=1}^p \tilde{\beta}_{lj} x_{il}$ . After simple calculation, the first term on the right-hand side (up to  $(1-w)\gamma_{j'}$ ) becomes

$$\sum_{i=1}^n x_{il'} r_i = \sum_{i=1}^n x_{il'} y_i - \gamma_0 \sum_{i=1}^n x_{il'} - \sum_{j: |\tilde{\beta}_{lj}| > 0} \gamma_j \tilde{\beta}_{lj} \mathbf{x}_i^T \mathbf{x}_l, \quad (10)$$

and the second term on the right-hand side (up to  $w$ ) is

$$\sum_{i=1}^n x_{il'} r_{ji}^* = \sum_{i=1}^n x_{il'} y_{ji}^* - \sum_{l: |\tilde{\beta}_{lj'}| > 0} \tilde{\beta}_{lj'} \mathbf{x}_i^T \mathbf{x}_l. \quad (11)$$

These formulas largely reduces computational task, because we update only the last term on (10) and (11) when the estimate of  $\beta_{lj'}$  is non-zero, while we do not update (10) and (11) when the estimate of  $\beta_{lj'}$  is zero.

Similarly, the update of the parameter  $\gamma$  in (9) is written as

$$\sum_{i=1}^n y_i^{**} x_{ij'}^* = \sum_{i=1}^n s_i x_{ij'}^* + \tilde{\gamma}_{j'} \sum_{i=1}^n x_{ij'}^{*2}, \quad (12)$$

where  $\tilde{\gamma}_{j'}$  is the current estimate of  $\gamma_{j'}$ . The first term on the right becomes

$$\sum_{i=1}^n s_i x_{ij'}^* = \sum_{i=1}^n y_i x_{ij'}^* - \gamma_0 \sum_{i=1}^n x_{ij'}^* - \sum_{j: |\tilde{\gamma}_j| > 0} \tilde{\gamma}_j \mathbf{x}_j^{*T} \mathbf{x}_j^*. \quad (13)$$

Therefore we update only the last term on (13) when the estimate of  $\gamma_{j'}$  is non-zero, while we do not update (13) when the estimate of  $\gamma_{j'}$  is zero.

### 4.3. Selection of tuning parameters

SPCR and aSPCR depend on four tuning parameters  $(w, \zeta, \lambda_\beta, \lambda_\gamma)$ . To avoid this hard computational task, we propose that the values of  $w$  and  $\zeta$  are set by a user, and then optimize only two tuning parameters  $\lambda_\beta$  and  $\lambda_\gamma$ .

The tuning parameter  $\zeta$  takes the role in the trade-off between the  $L_1$  and  $L_2$  penalties on  $B$ . For example, in elastic net, the value of  $\zeta$  is usually determined by a user, according to the purpose (Zou and Hastie, 2005). The tuning parameter  $w$  plays an important role in prediction accuracy. It seems that a small value of  $w$  is good for prediction. We tried many simulations with various values of  $w$  for various situations. A very small value of  $w$  often presented an unstable result and slow convergence. It would be because if the penalty for  $B$  was not enough, then the parameter  $B$  would be too flexible. Then, we concluded to set  $w = 0.1$  in this study. The value of  $w$  can also be determined by a user, according to the purpose. If a user wants to have explanatory variables close to principal components, a larger value of  $w$  would be favorable.

The tuning parameters  $\lambda_\beta$  and  $\lambda_\gamma$  are optimized using  $K$ -fold cross-validation. When the original dataset is divided into  $K$  datasets  $(\mathbf{y}^{(1)}, X^{(1)}), \dots, (\mathbf{y}^{(K)}, X^{(K)})$ , the CV criterion is given by

$$CV = \frac{1}{K} \sum_{k=1}^K \|\mathbf{y}^{(k)} - \hat{\gamma}_0^{(-k)} \mathbf{1}_{(k)} - X^{(k)} \hat{B}^{(-k)} \hat{\gamma}^{(-k)}\|^2,$$

where  $\mathbf{1}_{(k)}$  is a vector of which the elements are all one, and  $\hat{\gamma}_0^{(-k)}, \hat{B}^{(-k)}, \hat{\gamma}^{(-k)}$  are the estimates computed with the data removing the  $k$ -th part. We employed  $K = 5$  in our numerical study. The tuning parameters  $\lambda_\beta$  and  $\lambda_\gamma$  were, respectively, selected from 10 equally-spaced values on  $[\lambda_{\min}, \lambda_{\max}]$ , where  $\lambda_{\min}$  and  $\lambda_{\max}$  were determined according to the function `glmnet` in R.

## 5. Numerical study

### 5.1. Monte Carlo simulations

Monte Carlo simulations were conducted to investigate the performances of our proposed method. Three models were examined in this study.

In the first model, we considered the 10-dimensional covariate vector  $\mathbf{x} = (x_1, \dots, x_{10})^T$  according to a multivariate normal distribution with mean zero vector and variance-covariance matrix  $\Sigma_1$ , and generated the response  $y$  from the linear regression model given by

$$y_i = \xi_1^* x_{i1} + \xi_2^* x_{i2} + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad i = 1, \dots, n.$$

We used  $\xi_1^* = 2, \xi_2^* = 1, \Sigma_1 = I_{10}$  (Case 1(a)), where  $I_{10}$  is the  $10 \times 10$  identity matrix, and  $\xi_1^* = 8, \xi_2^* = 1, \Sigma_1 = \text{diag}(1, 3^2, \dots, 1)$  (Case 1(b)). Case 1(a) is a simple situation. Case 1(b) corresponds to the situation discussed in Section 3.3.

In the second model, we considered the 20-dimensional covariate vector  $\mathbf{x} = (x_1, \dots, x_{20})^T$  according to a multivariate normal distribution  $N(\mathbf{0}_{20}, \Sigma_2)$ , and generated the response  $y$  by

$$y_i = 4\mathbf{x}_i^T \xi^* + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad i = 1, \dots, n.$$

We used  $\Sigma_2 = \text{block diag}(\Sigma_2^*, I_{11})$  and  $\xi^* = (\mathbf{v}_1^*, 0, \dots, 0)^T$ , where  $(\Sigma_2^*)_{ij} = 0.9^{|i-j|}$  ( $i, j = 1, \dots, 9$ ) and  $\mathbf{v}_1^* = (-1, 0, 1, 1, 0, -1, -1, 0, 1)$  is a sparse approximation of the fourth eigenvector of  $\Sigma_2^*$  (Case 2). This case deals with the situation where the response is associated with the principal component loading with small eigenvalue. Note that even if each explanatory variable  $\mathbf{x}$  is normalized, the principal component  $\mathbf{x}^T \xi$  does not have unit variance in general.

In the third model, we assumed the 30-dimensional covariate vector  $\mathbf{x} = (x_1, \dots, x_{30})^T$  according to a multivariate normal distribution  $N(\mathbf{0}_{30}, \Sigma_3)$ , and generated the response  $y$  by

$$y_i = 4\mathbf{x}_i^T \xi_1^* + 4\mathbf{x}_i^T \xi_2^* + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad i = 1, \dots, n.$$

We used  $\Sigma_3 = \text{block diag}(\Sigma_2^*, \Sigma_3^*, I_{15})$  with  $(\Sigma_3^*)_{ij} = 0.9^{|i-j|}$  ( $i, j = 1, \dots, 6$ ), and  $\xi_1^* = (\mathbf{v}_1^*, 0, \dots, 0)^T$ . Two cases were considered for  $\xi_2^* = (0, \dots, 0, \mathbf{v}_2^*, 0, \dots, 0)^T$ , where the first nine and last 15 values are zero. First, we used  $\mathbf{v}_2^* = \underbrace{(1, \dots, 1)}_6$

that is an approximation of the first eigenvector of  $\Sigma_3^*$  (Case 3(a)). Second, we used  $\mathbf{v}_2^* = (1, 0, -1, -1, 0, 1)$  that is a sparse approximation of the third eigenvector of  $\Sigma_3^*$  (Case 3(b)). Case 3 is a more complex situation.

The sample size was set to  $n = 50, 200$ . The standard error  $\sigma$  was set to 0.1 or 1. Our proposed methods, SPCR and aSPCR, were fitted to the simulated data with one or 10 components ( $k = 1, 10$ ) for Case 1, one or five components ( $k = 1, 5$ ) for Case 2, and 10 components ( $k = 10$ ) for Case 3. Our proposed methods were compared with SPLS, PLS, and PCR. SPLS was computed by the package `spls` in R, and PLS and PCR by the package `pls` in R. The number of components and the values of tuning parameters in SPLS, PLS, and PCR were selected by 10-fold cross-validation. The performance was evaluated by  $MSE = E[(y - \hat{y})^2]$ . The simulation was conducted 100 times and the MSE was estimated by 1000 random samples. To estimate the parameter for SPCR and aSPCR, we used the algorithm proposed in Section 4 with the convergence condition



**Table 1**Mean (standard deviation) of MSE for  $\sigma = 0.1$ . The bold values correspond to the smallest mean.

Case	$k$	$n$	aSPCR	SPCR	SPLS	PLS	PCR
1(a)	1	50	<b><math>1.095 \times 10^{-2}</math></b> ( $9.906 \times 10^{-4}$ )	$1.654 \times 10^{-1}$ ( $8.799 \times 10^{-1}$ )	$2.952 \times 10^{-1}$ ( $3.919 \times 10^{-1}$ )	$8.877 \times 10^{-1}$ ( $3.885 \times 10^{-1}$ )	4.643 ( $6.325 \times 10^{-1}$ )
		200	<b><math>1.019 \times 10^{-2}</math></b> ( $5.088 \times 10^{-4}$ )	$5.735 \times 10^{-2}$ ( $4.702 \times 10^{-1}$ )	$3.167 \times 10^{-2}$ ( $3.095 \times 10^{-2}$ )	$2.249 \times 10^{-1}$ ( $9.559 \times 10^{-2}$ )	4.605 ( $5.240 \times 10^{-1}$ )
	10	50	$1.156 \times 10^{-2}$ ( $1.072 \times 10^{-3}$ )	$1.162 \times 10^{-2}$ ( $1.107 \times 10^{-3}$ )	<b><math>1.118 \times 10^{-2}</math></b> ( $1.304 \times 10^{-3}$ )	$1.283 \times 10^{-2}$ ( $1.380 \times 10^{-3}$ )	$1.282 \times 10^{-2}$ ( $1.379 \times 10^{-3}$ )
		200	$1.029 \times 10^{-2}$ ( $5.063 \times 10^{-4}$ )	$1.031 \times 10^{-2}$ ( $5.628 \times 10^{-4}$ )	<b><math>1.021 \times 10^{-2}</math></b> ( $5.120 \times 10^{-4}$ )	$1.054 \times 10^{-2}$ ( $5.216 \times 10^{-4}$ )	$1.054 \times 10^{-2}$ ( $5.218 \times 10^{-4}$ )
	1	50	<b><math>1.250 \times 10^{-2}</math></b> ( $2.220 \times 10^{-3}$ )	$1.465 \times 10^{-2}$ ( $2.778 \times 10^{-3}$ )	$4.043 \times 10^1$ ( $1.869 \times 10^1$ )	$4.595 \times 10^1$ ( $1.148 \times 10^1$ )	$6.650 \times 10^1$ (4.517)
		200	<b><math>1.131 \times 10^{-2}</math></b> ( $7.155 \times 10^{-4}$ )	$1.186 \times 10^{-2}$ ( $7.808 \times 10^{-4}$ )	$3.975 \times 10^1$ ( $1.531 \times 10^1$ )	$4.532 \times 10^1$ (5.048)	$6.457 \times 10^1$ (2.919)
	10	50	<b><math>1.092 \times 10^{-2}</math></b> ( $9.594 \times 10^{-4}$ )	$1.120 \times 10^{-2}$ ( $1.118 \times 10^{-3}$ )	$1.125 \times 10^{-2}$ ( $1.460 \times 10^{-3}$ )	$1.256 \times 10^{-2}$ ( $1.289 \times 10^{-3}$ )	$1.256 \times 10^{-2}$ ( $1.288 \times 10^{-3}$ )
		200	$1.027 \times 10^{-2}$ ( $4.968 \times 10^{-4}$ )	<b><math>1.022 \times 10^{-2}</math></b> ( $5.056 \times 10^{-4}$ )	$1.025 \times 10^{-2}$ ( $5.079 \times 10^{-4}$ )	$1.054 \times 10^{-2}$ ( $5.366 \times 10^{-4}$ )	$1.054 \times 10^{-2}$ ( $5.362 \times 10^{-4}$ )
2	1	50	<b><math>1.241 \times 10^{-2}</math></b> ( $1.738 \times 10^{-3}$ )	$1.614 \times 10^{-2}$ ( $3.601 \times 10^{-3}$ )	$1.978 \times 10^1$ (1.909)	$1.979 \times 10^1$ (1.851)	$2.038 \times 10^1$ (1.272)
		200	<b><math>1.051 \times 10^{-2}</math></b> ( $6.754 \times 10^{-4}$ )	$1.102 \times 10^{-2}$ ( $8.276 \times 10^{-4}$ )	$1.418 \times 10^1$ (4.475)	$1.571 \times 10^1$ (2.938)	$1.967 \times 10^1$ ( $8.374 \times 10^{-1}$ )
	5	50	<b><math>1.313 \times 10^{-2}</math></b> ( $2.207 \times 10^{-3}$ )	$1.548 \times 10^{-2}$ ( $3.708 \times 10^{-3}$ )	$3.946 \times 10^{-1}$ ( $6.452 \times 10^{-1}$ )	1.946 (1.337)	$2.118 \times 10^1$ (1.426)
		200	<b><math>1.077 \times 10^{-2}</math></b> ( $7.140 \times 10^{-4}$ )	$1.091 \times 10^{-2}$ ( $7.768 \times 10^{-4}$ )	$1.667 \times 10^{-2}$ ( $1.274 \times 10^{-2}$ )	$8.268 \times 10^{-2}$ ( $4.039 \times 10^{-2}$ )	$1.978 \times 10^1$ ( $8.926 \times 10^{-1}$ )
	10	50	<b><math>1.831 \times 10^{-2}</math></b> ( $4.842 \times 10^{-3}$ )	$2.191 \times 10^{-2}$ ( $6.641 \times 10^{-3}$ )	$3.438 \times 10^{-1}$ ( $4.319 \times 10^{-1}$ )	$8.493 \times 10^{-1}$ ( $6.014 \times 10^{-1}$ )	$2.839 \times 10^1$ (5.090)
		200	<b><math>1.158 \times 10^{-2}</math></b> ( $8.208 \times 10^{-4}$ )	$1.166 \times 10^{-2}$ ( $8.225 \times 10^{-4}$ )	$1.247 \times 10^{-2}$ ( $1.597 \times 10^{-3}$ )	$2.407 \times 10^{-2}$ ( $7.115 \times 10^{-3}$ )	$2.172 \times 10^1$ ( $1.463 \times 10^{-1}$ )
	10	50	<b><math>1.721 \times 10^{-2}</math></b> ( $5.311 \times 10^{-3}$ )	$2.180 \times 10^{-2}$ ( $6.390 \times 10^{-3}$ )	$4.852 \times 10^{-1}$ ( $6.966 \times 10^{-1}$ )	1.295 ( $9.401 \times 10^{-1}$ )	$3.676 \times 10^1$ (2.676)
		200	$1.185 \times 10^{-2}$ ( $9.778 \times 10^{-4}$ )	<b><math>1.167 \times 10^{-2}</math></b> ( $8.533 \times 10^{-4}$ )	$1.201 \times 10^{-2}$ ( $1.710 \times 10^{-3}$ )	$2.972 \times 10^{-2}$ ( $1.030 \times 10^{-2}$ )	$3.373 \times 10^1$ (1.605)

$\max_i |\hat{\theta}_i^{(r+1)} - \hat{\theta}_i^{(r)}| < 10^{-3}$ , where  $\hat{\theta}_i^{(r)}$  is the  $i$ th estimate of the parameter  $\theta$  for  $r$ -th iteration and  $\theta = (\gamma_0, \gamma^T, \{\text{vec}(B)\}^T)^T$ . The tuning parameter  $\zeta$  was set to be 0.01, because the number of parameters ( $pk + k + 1$ ) is larger than sample sizes ( $n$ ).

Tables 1 and 2 show the means and standard deviations of MSEs for  $\sigma = 0.1, 1$ , and present similar results. PCR was clearly the worst. SPLS was better than PLS, and aSPCR was basically better than SPCR. Therefore, we compare our methods, SPCR and aSPCR, with SPLS in more details.

In Case 1(a), aSPCR was basically better than SPLS for  $k = 1$  and competitive to SPLS for  $k = 10$ . In Case 1(b), SPCR and aSPCR provided much smaller MSEs than SPLS for  $k = 1$  and were competitive to SPLS for  $k = 10$ . The results for  $k = 1$  correspond to that discussed in Section 3.3. SPCR and aSPCR could appropriately select the loading related to the response. In Case 2, SPCR and aSPCR provided much smaller MSEs than SPLS for  $k = 1$ , like in Case 1(b) for  $k = 1$ , and aSPCR was better than SPLS for  $k = 5$ . In addition, SPCR and aSPCR provided almost the same MSEs for  $k = 1$  as those for  $k = 5$ . This means that SPCR and aSPCR could adaptively select the principal component loading with small eigenvalue. In Case 3, SPCR and aSPCR were better than SPLS. In complex situations for  $n = 50$ , aSPCR outperforms SPLS. We also compared our methods with lasso, adaptive lasso (aLasso), elastic net (EN), adaptive elastic net (aEN), and ordinary least squares (OLS) (see the supplementary material Appendix A). Our methods were better than or competitive with them, like SPLS was better than or competitive with them (Chun and Keleş, 2010).

We also computed the true positive rate (TPR) and the true negative rate (TNR) for aSPCR, SPCR, and SPLS, which are defined by

$$\text{TPR} = \frac{1}{100} \sum_{k=1}^{100} \frac{|\{j : \hat{\xi}_j^{(k)} \neq 0 \wedge \xi_j^* \neq 0\}|}{|\{j : \xi_j^* \neq 0\}|},$$

$$\text{TNR} = \frac{1}{100} \sum_{k=1}^{100} \frac{|\{j : \hat{\xi}_j^{(k)} = 0 \wedge \xi_j^* = 0\}|}{|\{j : \xi_j^* = 0\}|},$$

where  $\hat{\xi}_j^{(k)}$  is the estimated  $j$ th coefficient for the  $k$ -th simulation, and  $|\{*\}|$  is the number of elements included in a set  $\{*\}$ . Tables 3 and 4 show the means and standard deviations of TPR and TNR, and present similar results. In all cases, most of



**Table 2**Mean (standard deviation) of MSE for  $\sigma = 1$ . The bold values correspond to the smallest mean.

Case	$k$	$n$	aSPCR	SPCR	SPLS	PLS	PCR
1(a)	1	50	<b>1.266</b> ( $8.134 \times 10^{-1}$ )	1.638 (1.361)	1.475 ( $4.789 \times 10^{-1}$ )	1.999 ( $4.331 \times 10^{-1}$ )	5.663 ( $6.464 \times 10^{-1}$ )
		200	1.159 ( $8.267 \times 10^{-1}$ )	1.333 (1.169)	<b>1.031</b> ( $5.665 \times 10^{-2}$ )	1.256 ( $1.225 \times 10^{-1}$ )	5.598 ( $5.593 \times 10^{-1}$ )
	10	50	1.123 ( $1.163 \times 10^{-1}$ )	1.194 ( $1.142 \times 10^{-1}$ )	<b>1.122</b> ( $1.357 \times 10^{-1}$ )	1.283 ( $1.388 \times 10^{-1}$ )	1.282 ( $1.377 \times 10^{-2}$ )
		200	1.023 ( $4.983 \times 10^{-2}$ )	1.034 ( $5.214 \times 10^{-2}$ )	<b>1.021</b> ( $5.136 \times 10^{-2}$ )	1.054 ( $5.208 \times 10^{-2}$ )	1.054 ( $5.218 \times 10^{-2}$ )
	1(b)	50	<b>1.191</b> ( $1.260 \times 10^{-1}$ )	1.283 ( $1.383 \times 10^{-1}$ )	$4.144 \times 10^1$ ( $1.871 \times 10^1$ )	$4.711 \times 10^1$ ( $1.137 \times 10^1$ )	$6.748 \times 10^1$ (4.646)
		200	<b>1.030</b> ( $5.226 \times 10^{-2}$ )	1.062 ( $5.493 \times 10^{-2}$ )	$4.050 \times 10^1$ ( $1.565 \times 10^1$ )	$4.629 \times 10^1$ (5.246)	$6.560 \times 10^1$ (3.078)
	10	50	<b>1.088</b> ( $9.089 \times 10^{-2}$ )	1.160 ( $1.158 \times 10^{-1}$ )	1.140 ( $1.470 \times 10^{-1}$ )	1.295 ( $1.552 \times 10^{-1}$ )	1.297 ( $1.552 \times 10^{-1}$ )
		200	<b>1.020</b> ( $5.216 \times 10^{-2}$ )	1.031 ( $5.392 \times 10^{-2}$ )	1.023 ( $5.404 \times 10^{-2}$ )	1.054 ( $5.479 \times 10^{-2}$ )	1.054 ( $5.475 \times 10^{-2}$ )
	2	50	<b>1.284</b> ( $2.522 \times 10^{-1}$ )	1.583 ( $3.245 \times 10^{-1}$ )	$2.079 \times 10^1$ (1.788)	$2.084 \times 10^1$ (2.012)	$2.140 \times 10^1$ (1.295)
		200	<b>1.058</b> ( $5.566 \times 10^{-2}$ )	1.120 ( $6.347 \times 10^{-2}$ )	$1.568 \times 10^1$ (4.475)	$1.695 \times 10^1$ (2.981)	$2.086 \times 10^1$ ( $8.458 \times 10^{-1}$ )
	5	50	<b>1.279</b> ( $2.434 \times 10^{-1}$ )	1.576 ( $3.221 \times 10^{-1}$ )	2.017 (1.048)	3.398 (1.442)	$2.224 \times 10^1$ (1.476)
		200	<b>1.060</b> ( $5.671 \times 10^{-2}$ )	1.119 ( $6.323 \times 10^{-2}$ )	1.075 ( $5.837 \times 10^{-2}$ )	1.175 ( $7.427 \times 10^{-2}$ )	$2.097 \times 10^1$ ( $8.876 \times 10^{-1}$ )
3(a)	10	50	<b>1.607</b> ( $4.250 \times 10^{-1}$ )	2.274 ( $6.044 \times 10^{-1}$ )	2.403 ( $8.958 \times 10^{-1}$ )	2.724 ( $7.205 \times 10^{-1}$ )	$2.961 \times 10^1$ (5.070)
		200	<b>1.088</b> ( $7.104 \times 10^{-2}$ )	1.162 ( $7.882 \times 10^{-2}$ )	1.156 ( $2.621 \times 10^{-1}$ )	1.187 ( $7.714 \times 10^{-2}$ )	$2.277 \times 10^1$ (1.539)
3(b)	10	50	<b>1.482</b> ( $3.094 \times 10^{-1}$ )	2.180 ( $5.990 \times 10^{-1}$ )	2.364 ( $9.068 \times 10^{-1}$ )	3.081 ( $8.959 \times 10^{-1}$ )	$3.793 \times 10^1$ (2.835)
		200	<b>1.085</b> ( $6.686 \times 10^{-2}$ )	1.165 ( $7.719 \times 10^{-2}$ )	1.158 ( $4.742 \times 10^{-1}$ )	1.192 ( $7.631 \times 10^{-2}$ )	$3.482 \times 10^1$ (1.698)

TPRs are very high. For TNR, SPLS provides higher ratios for simple situations (Cases 1(a) and 1(b)), while aSPCR provides higher ratios for complex situations (Cases 2, 3(a), and 3(b)). In particular, in Cases 3(a) and 3(b) for  $n = 50$ , TNRs of aSPCR are much higher than those of SPLS.

## 5.2. Real data analyses

We examined the effectiveness of our proposed method through real data analyses. Five benchmark datasets were used—housing, energy, forest, concrete, and communities. The datasets were obtained from the UCI database (<http://archive.ics.uci.edu/ml/index.html>). The sample size and the numbers of covariates in these datasets are summarized in Table 5. For the ‘energy’ dataset, we used two types of response variables, following the explanation of the webpage. They are called ‘energy1’ and ‘energy2’ in this section. For the ‘communities’ dataset, we removed missing values by pre-processing. The covariates were standardized for each dataset.

First, using the ‘housing’ dataset, we illustrate a behavior of aSPCR. The estimates  $\hat{B}$  and  $\hat{\gamma}$  for aSPCR with  $k = 5$  were given by

$$\begin{aligned}\hat{\beta}_1 &= (0.025, 0.03, 0, 0, 0.058, 0, 0, 0, 0, 0.053, 0, 0, 0)^T, \\ \hat{\beta}_2 &= (0, 0, 0, -0.007, 0, -0.032, 0, 0.036, -0.028, 0, 0.024, -0.009, 0.045)^T, \\ \hat{\gamma}_1 &= -33.87, \quad \hat{\gamma}_2 = -83.30.\end{aligned}$$

We also observed that  $\hat{\gamma}_3 = \hat{\gamma}_4 = \hat{\gamma}_5 = 0$ , and then  $\hat{\beta}_3$ ,  $\hat{\beta}_4$ , and  $\hat{\beta}_5$  are omitted. The  $L_1$ -penalties on  $B$  and  $\gamma$  could produce zero estimates for the parameters, and then caused automatic selection of principal components. In addition, we have

$$\hat{B}\hat{\gamma} = (-0.87, 1.01, 0, 0.65, -1.97, 2.68, 0, -3.06, 2.38, -1.80, -2.03, 0.79, -3.75)^T,$$

which suggests that the third and seven variables (that is, variables **indus** and **age**) are irrelevant with the response variable. This fact was pointed out in some literatures (see, e.g., Shao and Rao, 2000; Khalili, 2010 and Leng, 2010).

Next, we randomly used 100 observations as training data to estimate the parameters and the remaining observations as test data to estimate the MSE. In addition, we randomly used 50 or 100 observations as training data for the ‘communities’

**Table 3**Mean (standard deviation) of TPR and TNR for  $\sigma = 0.1$ . The bold values correspond to the largest TPR and TNR.

Case	$k$	$n$	TPR			TNR		
			aSPCR	SPCR	SPLS	aSPCR	SPCR	SPLS
1(a)	1	50	<b>1</b> (0)	0.970 (0.171)	0.930 (0.174)	<b>1</b> (0)	0.615 (0.285)	0.982 (0.053)
		200	<b>1</b> (0)	0.990 (0.100)	<b>1</b> (0)	<b>1</b> (0)	0.631 (0.318)	<b>1</b> (0)
	10	50	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	0.693 (0.368)	0.496 (0.287)	<b>0.930</b> (0.130)
		200	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	0.562 (0.316)	0.528 (0.265)	<b>0.911</b> (0.160)
	1(b)	50	<b>1</b> (0)	<b>1</b> (0)	0.870 (0.220)	<b>1</b> (0)	0.061 (0.158)	0.926 (0.158)
		200	<b>1</b> (0)	<b>1</b> (0)	0.905 (0.197)	<b>1</b> (0)	0.070 (0.089)	0.963 (0.087)
2	1	50	<b>1</b> (0)	<b>1</b> (0)	0.870 (0.220)	<b>1</b> (0)	0.061 (0.158)	0.926 (0.158)
		200	<b>1</b> (0)	<b>1</b> (0)	0.905 (0.197)	<b>1</b> (0)	0.070 (0.089)	0.963 (0.087)
	10	50	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	<b>0.900</b> (0.253)	0.631 (0.246)	0.891 (0.192)
		200	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	0.757 (0.295)	0.753 (0.241)	<b>0.913</b> (0.146)
	5	50	<b>1</b> (0)	<b>1</b> (0)	0.995 (0.028)	<b>0.859</b> (0.111)	0.304 (0.196)	0.775 (0.135)
		200	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	0.905 (0.075)	0.387 (0.252)	<b>0.931</b> (0.073)
3(a)	10	50	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	<b>0.862</b> (0.102)	0.289 (0.168)	0.503 (0.146)
		200	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	<b>0.903</b> (0.062)	0.316 (0.216)	0.816 (0.079)
3(b)	10	50	<b>1</b> (0)	<b>1</b> (0)	0.998 (0.014)	<b>0.854</b> (0.092)	0.271 (0.155)	0.516 (0.165)
		200	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	<b>0.916</b> (0.061)	0.294 (0.182)	0.822 (0.083)

dataset to investigate the case where the number of explanatory variables is very large and the sample size is not large. Two cases with the sample sizes 50 and 100 were called ‘communities1’ and ‘communities2’, respectively. To estimate the parameter for SPCR and aSPCR, we used the algorithm proposed in Section 4 with the same convergence condition as in Section 5.1. The procedure was repeated 50 times.

Our proposed methods, SPCR and aSPCR, were compared with eight competing methods used in Section 5.1. For the ‘housing’, ‘energy’, ‘forest’ and ‘concrete’ datasets, the number of principal components or PLS components was set to  $k = 5$  for aSPCR, SPCR, SPLS, PLS, and PCR, while  $k = 10$  was adopted for the ‘communities’ dataset because the number of explanatory variables is very large. The tuning parameters  $\lambda_\beta$ ,  $\lambda_\gamma$ ,  $\zeta$  in SPCR and aSPCR were selected by five-fold cross-validation;  $\lambda_\beta$  and  $\lambda_\gamma$  were selected in similar manners to Section 5.1 and  $\zeta$  was selected from 0.1, 0.3, 0.5, 0.7, and 0.9 (in the ‘communities’ dataset, we fixed  $\zeta = 0.1$  due to large computational burden). The tuning parameters in other methods were selected in similar manners to in Section 5.1.

Table 6 shows the means and standard deviations of MSEs. For the ‘energy’ and ‘communities’ datasets, we could not obtain MSEs for OLS, because the design matrix included high correlation or high dimensionality. PCR was clearly worst except for the ‘forest’ and ‘communities’ datasets. SPCR and aSPCR were better than other methods for the ‘housing’, ‘energy2’ and ‘concrete’ datasets, and better or close to other methods for the ‘energy1’ dataset. The MSEs for the ‘forest’ and ‘communities’ datasets showed a different behavior from those for other datasets. PCR presented a good MSE, although aSPCR provided a smaller MSE than PCR in the ‘forest’ dataset, and SPCR was superior to PCR in the ‘communities’ dataset in terms of MSE. Furthermore, for almost all datasets, aSPCR was superior to SPLS, PLS, and PCR.

## 6. Concluding remarks

We proposed a one-stage procedure for PCR, which is constructed by combining a regression loss with PCA loss along with  $L_1$  type regularization. We called this procedure SPCR. SPCR enabled us to adaptively provide sparse principal component loadings that are associated with the response and to select the number of principal components automatically. The estimation algorithm for SPCR was established via the coordinate descent algorithm. To obtain a more sparse regression model, we also proposed aSPCR, which assigns different weights to different parameters in the loading matrix  $B$  in the estimation procedure. In numerical study, SPCR and aSPCR showed a good behavior in terms of prediction accuracy, TPR, and TNR.

**Table 4**Mean (standard deviation) of TPR and TNR for  $\sigma = 1$ . The bold values correspond to the largest TPR and TNR.

Case	$k$	$n$	TPR			TNR		
			aSPCR	SPCR	SPLS	aSPCR	SPCR	SPLS
1(a)	1	50	<b>0.970</b> (0.171)	0.910 (0.287)	0.910 (0.193)	0.791 (0.247)	0.258 (0.277)	<b>0.953</b> (0.128)
		200	0.970 (0.171)	0.940 (0.238)	<b>1</b> (0)	0.870 (0.183)	0.250 (0.255)	<b>0.998</b> (0.012)
	10	50	<b>1</b> (0)	0.990 (0.100)	<b>1</b> (0)	0.802 (0.334)	0.227 (0.168)	<b>0.931</b> (0.141)
		200	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	0.737 (0.353)	0.318 (0.204)	<b>0.911</b> (0.164)
	1(b)	50	<b>1</b> (0)	<b>1</b> (0)	0.870 (0.220)	0.580 (0.211)	0.012 (0.057)	<b>0.915</b> (0.166)
		200	<b>1</b> (0)	<b>1</b> (0)	0.900 (0.201)	0.730 (0.179)	0.007 (0.029)	<b>0.966</b> (0.083)
	10	50	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	<b>0.933</b> (0.178)	0.552 (0.239)	0.896 (0.210)
		200	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	0.951 (0.213)	0.647 (0.270)	<b>0.955</b> (0.129)
2	1	50	<b>1</b> (0)	<b>1</b> (0)	0.543 (0.313)	<b>0.865</b> (0.182)	0.172 (0.139)	0.726 (0.317)
		200	<b>1</b> (0)	<b>1</b> (0)	0.860 (0.215)	<b>0.930</b> (0.122)	0.202 (0.153)	0.775 (0.253)
	5	50	<b>1</b> (0)	<b>1</b> (0)	0.993 (0.032)	<b>0.872</b> (0.191)	0.176 (0.145)	0.648 (0.200)
		200	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	0.892 (0.190)	0.205 (0.150)	<b>0.896</b> (0.111)
	3(a)	50	0.999 (0.008)	<b>1</b> (0)	0.998 (0.011)	<b>0.885</b> (0.148)	0.142 (0.101)	0.423 (0.220)
		200	<b>1</b> (0)	<b>1</b> (0)	0.999 (0.008)	<b>0.901</b> (0.164)	0.165 (0.122)	0.846 (0.163)
3(b)	10	50	<b>1</b> (0)	<b>1</b> (0)	0.999 (0.010)	<b>0.880</b> (0.130)	0.184 (0.128)	0.430 (0.202)
		200	<b>1</b> (0)	<b>1</b> (0)	0.998 (0.020)	<b>0.875</b> (0.203)	0.223 (0.162)	0.864 (0.148)

**Table 5**

Sample size and the numbers of covariates in real datasets.

	Housing	Energy	Forest	Concrete	Communities
Sample size	506	768	517	1030	1993
# of covariates	13	8	10	8	101

**Table 6**

Mean (standard deviation) of MSE for real datasets. The bold values correspond to the smallest mean.

	Housing	Energy1	Energy2	Forest	Concrete	Communities1	Communities2
SPCR	<b>28.94</b> (4.402)	11.07 (0.612)	<b>9.248</b> (0.468)	4680 (466.6)	<b>121.2</b> (12.13)	<b><math>2.823 \times 10^{-2}</math></b> ( $4.259 \times 10^{-3}$ )	<b><math>2.229 \times 10^{-2}</math></b> ( $1.992 \times 10^{-3}$ )
aSPCR	29.18 (4.628)	11.04 (0.568)	9.275 (0.519)	4569 (757.4)	121.3 (12.25)	$3.998 \times 10^{-2}$ ( $1.789 \times 10^{-2}$ )	$3.589 \times 10^{-2}$ ( $1.115 \times 10^{-2}$ )
SPLS	30.24 (3.845)	11.14 (0.634)	9.405 (0.532)	4579 (652.0)	125.5 (11.74)	$3.384 \times 10^{-2}$ ( $8.897 \times 10^{-3}$ )	$2.409 \times 10^{-2}$ ( $4.654 \times 10^{-3}$ )
PLS	29.78 (4.467)	11.18 (0.627)	9.386 (0.567)	4683 (471.0)	122.0 (11.67)	$6.124 \times 10^{-2}$ ( $3.948 \times 10^{-2}$ )	$3.677 \times 10^{-2}$ ( $8.100 \times 10^{-3}$ )
PCR	30.45 (3.478)	14.90 (0.641)	12.50 (0.523)	4599 (600.7)	159.4 (27.04)	$2.932 \times 10^{-2}$ ( $4.464 \times 10^{-3}$ )	$2.254 \times 10^{-2}$ ( $2.503 \times 10^{-3}$ )
Lasso	29.80 (4.055)	11.21 (0.637)	9.379 (0.529)	4534 (752.4)	122.1 (12.08)	$3.221 \times 10^{-2}$ ( $1.118 \times 10^{-2}$ )	$2.355 \times 10^{-2}$ ( $4.149 \times 10^{-3}$ )
aLasso	30.16 (4.601)	11.06 (0.549)	9.313 (0.521)	4542 (709.6)	122.7 (13.96)	$4.309 \times 10^{-2}$ ( $2.013 \times 10^{-2}$ )	$2.334 \times 10^{-2}$ ( $3.211 \times 10^{-3}$ )
EN	29.56 (3.953)	11.18 (0.624)	9.340 (0.516)	4534 (749.5)	122.7 (12.36)	$3.612 \times 10^{-2}$ ( $3.644 \times 10^{-2}$ )	$2.311 \times 10^{-2}$ ( $3.124 \times 10^{-3}$ )
aEN	30.05 (4.057)	<b>11.03</b> (0.552)	9.286 (0.527)	4572 (628.5)	123.1 (14.32)	$4.590 \times 10^{-2}$ ( $3.827 \times 10^{-2}$ )	$2.322 \times 10^{-2}$ ( $3.112 \times 10^{-3}$ )
OLS	29.66 (6.592)	–	–	<b>4318</b> (951.7)	123.1 (13.65)	–	–

## Acknowledgments

The authors would like to thank the anonymous reviewers for their constructive and helpful comments. This work was supported by the Bio-diversity Research Project of the Transdisciplinary Research Integration Center, Research Organization of Information and Systems.

## Appendix A. Supplementary data

Supplementary material related to this article can be found online at <http://dx.doi.org/10.1016/j.csda.2015.03.016>.

## References

- Allen, G.I., Peterson, C., Vannucci, M., Maletić-Savatić, M., 2013. Regularized partial least squares with an application to NMR spectroscopy. *Stat. Anal. Data Min.* 5, 302–314.
- Bair, E., Hastie, T., Paul, D., Tibshirani, R., 2006. Prediction by supervised principal components. *J. Amer. Statist. Assoc.* 101, 119–137.
- Chen, L., Huang, J.Z., 2012. Sparse reduced-rank regression for simultaneous dimension reduction and variable selection. *J. Amer. Statist. Assoc.* 107, 1533–1545.
- Choi, J., Zou, H., Oehlert, G., 2011. A penalized maximum likelihood approach to sparse factor analysis. *Stat. Interface* 3, 429–436.
- Chun, H., Keleş, S., 2010. Sparse partial least squares regression for simultaneous dimension reduction and variable selection. *J. R. Stat. Soc. Ser. B* 72, 3–25.
- de Jong, S., 1993. SIMPLS: An alternative approach to partial least squares regression. *Chemometr. Intell. Lab. 18*, 251–263.
- Efron, B., Hastie, T., Johnstone, I., Tibshirani, R., 2004. Least angle regression. *Ann. Statist.* 32, 407–499.
- Frank, I., Friedman, J., 1993. A statistical view of some chemometrics regression tools. *Technometrics* 35, 109–135.
- Friedman, J., Hastie, T., Tibshirani, R., 2010. Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* 33, 1–22.
- Han, F., Liu, H., 2013. Robust sparse principal component regression under the high dimensional elliptical model. *Adv. Neural Inf. Process. Syst.* 26, 1941–1949.
- Hartnett, M.K., Lightbody, G., Irwin, G.W., 1998. Dynamic inferential estimation using principal components regression (PCR). *Chemometr. Intell. Lab.* 40, 215–224.
- Hastie, T., Tibshirani, R., Friedman, J., 2009. *The Elements of Statistical Learning*, second ed. Springer, New York.
- Hirose, K., Yamamoto, M., 2014. Sparse estimation via nonconcave penalized likelihood in a factor analysis model. *Stat. Comput.* in press.
- Jennrich, R.I., 2006. Rotation to simple loadings using component loss functions: The oblique case. *Psychometrika* 71, 173–191.
- Jolliffe, I.T., 1982. A note on the use of principal components in regression. *Appl. Stat.* 31, 300–303.
- Jolliffe, I.T., Trendafilov, N.T., Uddin, M., 2003. A modified principal component technique based on the LASSO. *J. Comput. Graph. Statist.* 12, 531–547.
- Khalili, A., 2010. New estimation and feature selection methods in mixture-of-experts models. *Canad. J. Statist.* 38, 519–539.
- Lee, S., Huang, J.Z., 2013. A coordinate descent MM algorithm for fast computation of sparse logistic PCA. *Comput. Statist. Data Anal.* 62, 26–38.
- Lee, S., Huang, J.Z., Hu, J., 2010. Sparse logistic principal components analysis for binary data. *Ann. Appl. Stat.* 4, 1579–1601.
- Leng, C., 2010. Variable selection and coefficient estimation via regularized rank regression. *Statist. Sinica* 20, 167–181.
- Massy, W.F., 1965. Principal components regression in explanatory statistical research. *J. Amer. Statist. Assoc.* 60, 234–256.
- Pearson, K., 1901. On lines and planes of closest fit to systems of points in space. *Phil. Mag.* 2, 559–572.
- Reiss, P.T., Ogden, R.T., 2007. Functional principal component regression and functional partial least squares. *J. Amer. Statist. Assoc.* 102, 984–996.
- Rosipal, R., Girolami, M., Trejo, L.J., Cichocki, A., 2001. Kernel PCA for feature extraction and de-noising in non-linear regression. *Neural Comput. Appl.* 10, 231–243.
- Shao, J., Rao, J.S., 2000. The GIC for model selection: a hypothesis testing approach. *J. Statist. Plann. Inference* 88, 215–231.
- Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B* 58, 267–288.
- Wang, K., Abbott, D., 2008. A principal components regression approach to multilocus genetic association studies. *Genet. Epidemiol.* 32, 108–118.
- Wold, H., 1975. Soft modeling by latent variables: The nonlinear iterative partial least squares approach. In: Gani, J. (Ed.), *Perspectives in Probability and Statistics, Papers in Honor of MS Bartlett*. pp. 520–540.
- Yu, S., Yu, K., Tresp, V., Kriegel, H.-P., Wu, M., 2006. Supervised probabilistic principal component analysis. In: *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 464–473.
- Zou, H., 2006. The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.* 101, 1418–1429.
- Zou, H., Hastie, T., 2005. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B* 67, 301–320.
- Zou, H., Hastie, T., Tibshirani, R., 2006. Sparse principal component analysis. *J. Comput. Graph. Statist.* 15, 265–286.