Received: 4 April 2014,

Revised: 15 October 2014,

(wileyonlinelibrary.com) DOI: 10.1002/cem.2691

High-throughput data dimension reduction via seeded canonical correlation analysis

Accented: 24 October 2014

Yunju Im, HeyIn Gang and Jae Keun Yoo*

Canonical correlation analysis (CCA) is one of popular statistical methodologies in multivariate analysis, especially, in studying relation of two sets of variables. However, if sample sizes are smaller than the maximum of the dimensions of two sets of variables, it is not plausible to construct canonical coefficient matrices due to failure of inverting sample covariance matrices. In this article, we develop a two step procedure of CCA implemented in such situation. For this, seeded dimension reduction is adapted into CCA. Numerical studies confirm the approach, and two real data analyses are presented. Copyright © 2014 John Wiley & Sons, Ltd.

Keywords: canonical correlation analysis; large p small n; multivariate analysis; seeded dimension reduction

1. INTRODUCTION

Canonical correlation analysis (CCA) is a method to measure an association between two sets of variables and focuses on seeking pairs of linear combinations from two sets of variables by maximizing the Pearson correlation between two sets of variables. The pairs of linear combination and their correlations are called canonical variates and canonical correlations, respectively. A few pairs of canonical variates are believed to be able to represent the relation between the original sets of variables and their variabilities. We briefly introduce the major concepts and statistics of CCA.

Suppose that we are interested in measures of association between two sets of variables of $\mathbf{X} \in \mathbb{R}^{p}$ and $\mathbf{Y} \in \mathbb{R}^{r}$. Define that $\operatorname{cov}(\mathbf{X}) = \Sigma_{X} > 0$, $\operatorname{cov}(\mathbf{Y}) = \Sigma_{y} > 0$, $\operatorname{cov}(\mathbf{X}, \mathbf{Y}) = \Sigma_{xy}$, and $\operatorname{cov}(\mathbf{Y}, \mathbf{X}) = \Sigma_{yx}$. For two linear combinations of \mathbf{X} and \mathbf{Y} , saying $U = \mathbf{a}^{\mathsf{T}}\mathbf{X}$ and $V = \mathbf{b}^{\mathsf{T}}\mathbf{Y}$, it is obtained that $\operatorname{var}(U) = \mathbf{a}^{\mathsf{T}}\Sigma_{X}\mathbf{a}$, $\operatorname{var}(V) = \mathbf{b}^{\mathsf{T}}\Sigma_{y}\mathbf{b}$ and $\operatorname{cov}(U, V) = \mathbf{a}^{\mathsf{T}}\Sigma_{xy}\mathbf{b}$, where $\mathbf{a} \in \mathbb{R}^{p \times 1}$ and $\mathbf{b} \in \mathbb{R}^{r \times 1}$.

We seek to find **a** and **b** to maximize Pearson-correlation between *U* and *V*:

$$\operatorname{cor}(U, V) = \frac{\mathbf{a}^{\mathsf{T}} \boldsymbol{\Sigma}_{xy} \mathbf{b}}{\sqrt{\mathbf{a}^{\mathsf{T}} \boldsymbol{\Sigma}_{x} \mathbf{a}} \sqrt{\mathbf{b}^{\mathsf{T}} \boldsymbol{\Sigma}_{y} \mathbf{b}}}$$
(1)

In CCA, such **a** and **b** are constructed based on the following criteria:

- (1) The first canonical variate pair $(U_1 = \mathbf{a}_1^T \mathbf{X}, V_1^T = \mathbf{b}_1^T \mathbf{Y})$ is constructed from the maximization of (1)
- (2) At the $k \ge 2$ step, the *k*th canonical variate pair ($U_k = \mathbf{a}_k^T \mathbf{X}, V_k = \mathbf{b}_k^T \mathbf{Y}$) is constructed from the maximization of (1) with restriction that $var(U_k) = var(V_k) = 1$ and (U_k, V_k) are uncorrelated with the previous (k 1) canonical variate pairs.
- (3) Repeat Steps 1 and 2 until $k = \min(p, r)$.
- (4) Select the first *d* pairs of (U_k, V_k) to represent the relationship between **X** and **Y**.

Then the pairs $(\mathbf{a}_i, \mathbf{b}_i)$ are acquired as follows: $\mathbf{a}_i = \Sigma_x^{-1/2} \psi_i$ and $\mathbf{b}_i = \Sigma_y^{-1/2} \phi_i$ for $i = 1, \dots, q$, where (ψ_1, \dots, ψ_q) and (ϕ_1, \dots, ϕ_q) are the eigenvectors of $\Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1} \Sigma_{yx} \Sigma_y^{-1/2}$ and $\Sigma_y^{-1/2} \Sigma_{yx} \Sigma_x^{-1} \Sigma_{xy} \Sigma_x^{-1/2}$ with the corresponding common nonzero ordered-eigenvalues of $\rho_1^{*2} \ge \dots \ge \rho_q^{*2} \ge 0$, respectively. The matrices of $\mathbf{M}_x = (\mathbf{a}_1, \dots, \mathbf{a}_q)$ and $\mathbf{M}_y = (\mathbf{b}_1, \dots, \mathbf{b}_q)$ are called canonical coefficient matrices.

Hereafter, the CCA by the decompositions of $\Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1} \Sigma_{yx} \Sigma_y^{-1/2}$ and $\Sigma_y^{-1/2} \Sigma_{yx} \Sigma_x^{-1} \Sigma_{xy} \Sigma_x^{-1/2}$ will be called the *standard CCA*. For more details regarding the standard CCA, readers may refer to Johnson and Wichern [1].

One major problem arises with a sample size n less than or equal to max(p, r), because the sample covariance matrices of Σ_X and Σ_{γ} are not invertible. To overcome this, various penalized approaches were proposed. Parkhomenko et al. [2] suggested a CCA to be carried out by penalizing the left and right singular vectors in the covariance matrix of Σ_{xy} . Waaijenborg *et al.* [3] converted the standard CCA to regression forms and adopted a penalized method called elastic net (Zou and Hastie [4]). And, a sparse CCA proposed by Le Cao et al. [5] adapted sparse partial least squares into CCA. Witten et al. [6] applied a penalized matrix decomposition into CCA. One common thing in the various versions of the sparse CCAs is to overcome the matrix inversion problem and to make parts of canonical coefficient matrices zeros. For this, tuning parameters are required, which are typically determined by cross-validation procedures. Therefore, when $\min(p, r)$ is large, the sparse CCA methods often turn out to computationally intensive in practice.

In this paper, we propose an approach of CCA applicable with $n \leq \max(p, r)$ by adapting a seeded dimension reduction

Department of Statistics, Ewha Womans University

^{*} Correspondence to: Jae Keun Yoo, Department of Statistics, Ewha Womans University, Seoul 120-750, Korea E-mail: peter.yoo@ewha.ac.kr

(Cook *et al.* [7]). Different from the sparse CCAs, the proposed method does not make parts of canonical coefficient matrices zeros intentionally. Instead, seed matrices, constructed from Σ_{xy} and Σ_{yx} , are iteratively projected into the marginal covariance matrices of Σ_x and Σ_y . By this, the proposed method can be enjoyed relatively less intensiveness in computation.

This paper is organized as follows. We briefly discuss seeded dimension reduction in Section 2. Section 3 is devoted to associate CCA with seeded dimension reduction for developing seeded canonical correlation analysis. Numerical studies and real data analysis are presented in Section 4. In Section 5, we summarize our works.

We will define notations frequently used. For a matrix $\mathbf{B} \in \mathbb{R}^{p \times q}$ and a subspace S of \mathbb{R}^p , $\mathbf{B}S$ denotes the set of { $\mathbf{B}x : x \in S$ }. A subspace $S(\mathbf{b})$ represents a subspace spanned by the columns of \mathbf{b} . For a symmetric and positive definite matrix Σ , a Σ inner-product in \mathbb{R}^p is defined as $\langle \mathbf{a}, \mathbf{b} \rangle_{\Sigma} = \mathbf{a}^T \Sigma \mathbf{b}$. An orthogonal projection operator onto $S(\mathbf{b})$ relative to $\langle \mathbf{a}, \mathbf{b} \rangle_{\Sigma}$ will be defined as $\mathbf{b}(\mathbf{b}^T \Sigma \mathbf{b})^{\dagger} \mathbf{b}^T \Sigma$, where \dagger stands for the Moore-Penrose inverse (Searle [8]). And, the true rank of \mathbf{M}_x and \mathbf{M}_y will be denoted as d.

2. SEEDED DIMENSION REDUCTION

The primary purpose of sufficient dimension reduction (SDR) in regression of $\mathbf{Y} \in \mathbb{R}^r | \mathbf{X} \in \mathbb{R}^p$ is to replace the original *p*-dimensional predictors \mathbf{X} with a lower dimensional linear projection without loss of information on the conditional distribution of $\mathbf{Y} | \mathbf{X}$. In other words, SDR seeks to search $\alpha \in \mathbb{R}^{p \times q}$ such that

$$\mathbf{Y} \perp \!\!\!\perp \mathbf{X} | \boldsymbol{\alpha}^{\mathsf{T}} \mathbf{X}, \tag{2}$$

where $\bot\!\!\!\!\perp$ represents independence and $q \leq p$.

Subspaces spanned by the columns of α satisfying (2) are called dimension reduction subspaces, and the minimal subspace among them is called the central subspace $S_{\mathbf{Y}|\mathbf{X}}$. Naturally, the estimation of $S_{\mathbf{Y}|\mathbf{X}}$ is the main stream in the context of SDR.

In general, many SDR methods require the inversion of the covariance matrix Σ_X . However, with sample sizes $n \leq p$, the inversion of the sample version of Σ_X is not plausible, and hence the SDR methods may not be applicable.

In order to overcome the problematic issue, Cook *et al.* [7] proposed a method without matrix inversion. For this to be carried out, a seed matrix $v \in \mathbb{R}^{p \times q}$ should be defined so that $S(v) \subseteq \Sigma_X S_{Y|X}$. In order to avoid complications, it is assumed that $S(v) = \Sigma_X S_{Y|X}$.

Suppose that there is a known subspace $\mathcal{M}_{\mathbf{Y}|\mathbf{X}}$ of \mathbb{R}^p containing $\mathcal{S}_{\mathbf{Y}|\mathbf{X}}$. This indicates that $\Sigma_X^{-1}\mathcal{S}(\nu) \subseteq \mathcal{M}_{\mathbf{Y}|\mathbf{X}}$. Let $\mathbf{P}_{\mathcal{M}_{\mathbf{Y}|\mathbf{X}}}(\Sigma_X) = \mathbf{R}(\mathbf{R}^T \Sigma_X \mathbf{R})^{-1} \mathbf{R}^T \Sigma_X$ be an orthogonal projection operator onto $\mathcal{M}_{\mathbf{Y}|\mathbf{X}}$ relative to < **a**, **b** > Σ_X , where **R** is a $p \times q$ matrix such that $\mathcal{S}(\mathbf{R}) = \mathcal{M}_{\mathbf{Y}|\mathbf{X}}$.

The previous discussion directly implies the following equivalences:

$$\Sigma_{X}^{-1} \boldsymbol{\nu} = \mathbf{P}_{\mathcal{M}_{\mathbf{Y}|\mathbf{X}}(\boldsymbol{\Sigma}_{X})} \Sigma_{X}^{-1} \boldsymbol{\nu} = \mathbf{R} (\mathbf{R}^{\mathsf{T}} \boldsymbol{\Sigma}_{X} \mathbf{R})^{-1} \mathbf{R}^{\mathsf{T}} \boldsymbol{\Sigma}_{X} \boldsymbol{\Sigma}_{X}^{-1} \boldsymbol{\nu}$$

= $\mathbf{R} (\mathbf{R}^{\mathsf{T}} \boldsymbol{\Sigma}_{X} \mathbf{R})^{-1} \mathbf{R}^{\mathsf{T}} \boldsymbol{\nu}.$ (3)

According to the last equivalence of (3), the columns of $\mathbf{R}(\mathbf{R}^{\mathsf{T}}\boldsymbol{\Sigma}_{X}\mathbf{R})^{-1}\mathbf{R}^{\mathsf{T}}$ span $\mathcal{S}_{\mathbf{Y}|\mathbf{X}}$, but the inversion of $\boldsymbol{\Sigma}_{X}$ is not required. If $\mathbf{R}^{\mathsf{T}}\boldsymbol{\Sigma}_{X}\mathbf{R}$ is not invertible, $(\mathbf{R}^{\mathsf{T}}\boldsymbol{\Sigma}_{X}\mathbf{R})^{\dagger}$ is applied instead. To estimate $S_{\mathbf{Y}|\mathbf{X}}$ through $\mathbf{R}(\mathbf{R}^T \Sigma_X \mathbf{R})^{-1} \mathbf{R}^T \nu$, we need to identify the matrix **R**, whose column subspace is large enough to enclose $S_{\mathbf{Y}|\mathbf{X}}$ and small enough to estimate $S_{\mathbf{Y}|\mathbf{X}}$ from available data. In order to find the matrix **R**, iterative projections of ν onto Σ_X were suggested by Cook *et al.* [7]:

$$\mathbf{R}_{u} \equiv (\mathbf{v}, \boldsymbol{\Sigma}_{X} \mathbf{v}, \boldsymbol{\Sigma}_{X}^{2} \mathbf{v}, \dots, \boldsymbol{\Sigma}_{X}^{u-1} \mathbf{v}), \ u = 1, 2, \dots, u^{*}.$$
(4)

We call the letter u in (4) a termination index of the projections. It is noted that $S(\mathbf{R}_{u-1}) \subseteq S(\mathbf{R}_u)$ for any $u \ge 2$. Because $S(\mathbf{R}_u)$ forms a nondecreasing sequence, it is important to select a proper termination index u, large enough to guarantee $S(\mathbf{R}_u) = \mathcal{M}_{\mathbf{Y}|\mathbf{X}}$ and small enough to capture $S_{\mathbf{Y}|\mathbf{X}}$. Recently, Yoo [9] suggests bootstrap coefficients of variations to determine the termination index, which does not require any asymptotics and is implemented in a simple way.

3. SEEDED CANONICAL CORRELATION ANALYSIS

3.1. Development

To begin with, we need to see why the seeded dimension reduction method can be adapted to CCA, by investigating the relation between ordinary least squares and canonical coefficient matrices.

Recall the definitions of the canonical coefficient matrices M_x and M_y from the introduction. According to Lee and Yoo [10], it has been shown that

$$\mathcal{S}(\mathbf{M}_{x}) = \mathcal{S}(\boldsymbol{\Sigma}_{x}^{-1}\boldsymbol{\Sigma}_{xy}) \text{ and } \mathcal{S}(\mathbf{M}_{y}) = \mathcal{S}(\boldsymbol{\Sigma}_{y}^{-1}\boldsymbol{\Sigma}_{yx}),$$
(5)

where $\Sigma_x^{-1} \Sigma_{xy}$ and $\Sigma_y^{-1} \Sigma_{yx}$ are the ordinary least squares coefficient matrices of **Y**|**X** and **X**|**Y**, respectively.

The relation in (5) directly indicates that the information on \mathbf{M}_x and \mathbf{M}_y is exhaustively restored through $\Sigma_x^{-1} \Sigma_{xy}$ and $\Sigma_y^{-1} \Sigma_{yx}$. Another interpretation of (5) is that one of possible basis matrices for $S(\Sigma_x^{-1} \Sigma_{xy})$ can become \mathbf{M}_x . Therefore, once any basis matrices of $S(\Sigma_x^{-1} \Sigma_{xy})$ and $S(\Sigma_y^{-1} \Sigma_{xy})$ are known, \mathbf{M}_x and \mathbf{M}_y can be restored through some orthogonalization procedures.

If $\max(p, r) \ge n$, \mathbf{M}_x and \mathbf{M}_y cannot be estimated through the standard CCA due to failure of inverting the sample versions of Σ_x and Σ_y . Meanwhile, the basis matrices of $\Sigma_x^{-1} \Sigma_{xy}$ and $\Sigma_y^{-1} \Sigma_{yx}$ are estimable through seeded dimension reduction with seed matrices of Σ_{xy} and Σ_{yx} . Therefore, although the direct estimation of \mathbf{M}_x and \mathbf{M}_y is not plausible through the standard CCA with $\max(p, r) \ge n$, they can be indirectly restored by estimating basis matrices of $\Sigma_x^{-1} \Sigma_{xy}$ and $\Sigma_y^{-1} \Sigma_{yx}$.

Let $\mathbf{M}_{x,0}$ and $\mathbf{M}_{y,0}$ be matrices resulted from the seeded dimension reduction application with seed matrices of Σ_{xy} and Σ_{yx} . Then the relation in (5) directly indicates that

$$\mathcal{S}(\mathbf{M}_{x}) = \mathcal{S}(\mathbf{M}_{x,0}) \text{ and } \mathcal{S}(\mathbf{M}_{y}) = \mathcal{S}(\mathbf{M}_{y,0}).$$
 (6)

Next, the original two sets of **X** and **Y** are replaced with $\mathbf{M}_{x,0}^{\mathsf{I}}\mathbf{X}$ and $\mathbf{M}_{y,0}^{\mathsf{T}}\mathbf{Y}$. Because the ranks of \mathbf{M}_x and \mathbf{M}_y , denoted as *d*, are less than min(*p*, *r*), these replacements result in dimension reductions of **X** and **Y**, and it will be called *initialized CCA*, hereafter. It should be noted that the initialized CCA is guaranteed to be equally informative to the standard CCA application by (6). Although the ranks of $\mathbf{M}_{x,0}^{\mathsf{T}} \mathbf{X}$ and $\mathbf{M}_{y,0}^{\mathsf{T}} \mathbf{Y}$ are theoretically less than min(p, r), $\mathbf{M}_{x,0}$ and $\mathbf{M}_{y,0}$ are still $p \times r$ and $r \times p$ matrices, respectively. Also, $\mathbf{M}_{x,0}$ and $\mathbf{M}_{y,0}$ may not be satisfied with the orthonormality required in the standard CCA. Therefore, the CCA dimension reductions are completed by the standard CCA application with $\mathbf{M}_{x,0}^{\mathsf{T}} \mathbf{X}$ and $\mathbf{M}_{y,0}^{\mathsf{T}} \mathbf{Y}$, which yields the final canonical coefficient matrices of \mathbf{M}_x and \mathbf{M}_y . This procedure will be called *finalized CCA*. To perform the finalized CCA in practice, we need to have that d is relatively smaller than n. It is not heavy in practice, because this is the main reason for dimension reduction.

The proposed two-step CCA procedure done by the initialized and finalized CCAs will be called *seeded canonical correlation analysis*. Two implementations of the seeded CCA will be presented in next section.

3.2. Implementation

Recall that $\mathbf{M}_{x,0}$ and $\mathbf{M}_{y,0}$ are resulted from the initialized CCA application. The replacements of **X** and **Y** by $\mathbf{M}_{x,0}^{\mathsf{T}} \mathbf{X}$ and $\mathbf{M}_{y,0}^{\mathsf{T}} \mathbf{Y}$ cause initial dimension reductions, which are theoretically equivalent to the standard CCA. In implementation perspective, however, the dimensions of Σ_{xy} and Σ_{yx} used in the initialized CCA can affect the accuracies in estimation of $\mathbf{M}_{x,0}$ and $\mathbf{M}_{y,0}$ and numerical stability in computation. Hence, if the dimensions of Σ_{xy} and Σ_{yx} are large, Σ_{xy} and Σ_{yx} should be adjusted small enough to yield $\mathbf{M}_{x,0}$ and $\mathbf{M}_{y,0}$ equally informative to the standard CCA. Depending on the size of min(*p*, *r*), two versions of seeded CCA are proposed.

3.2.1. Case 1: min(p, r) fairly small

The first case is that min(p, r) is fairly smaller than n. For simplicity and without loss of generality, we assume that r < p, and the maximum number of pairs of canonical variates is equal to r. Then the initialized CCA is carried out with $\hat{\Sigma}_{xy}$ and $\hat{\Sigma}_{yx}$ as the seed matrices, and **X** alone is replaced with $\hat{\mathbf{M}}_{x,0}^{\mathsf{T}} \mathbf{X}$. Next, by the finalized CCA with $\hat{\mathbf{M}}_{x,0}^{\mathsf{T}} \mathbf{X}$ and the original **Y**, the CCA reduction is completed.

3.2.2. Case 2: min(p, r) large compared with n

The second case is that min(*p*, *r*) is not fairly small compared with *n*, or even larger than *n*. Then, as the seed matrices in the initialized CCA, the sets of the eigenvectors corresponding to *m* largest eigenvalues of $\hat{\Sigma}_{xy}$ and $\hat{\Sigma}_{yx}$ are used. The verification of these replacements is summarized in the next lemma.

Lemma 1

Assume that $d < \min(p, r)$, where d stand for the true ranks of rank(\mathbf{M}_X) and rank(\mathbf{M}_y). Let $v_X \in \mathbb{R}^{p \times d}$ and $v_y \in \mathbb{R}^{r \times d}$ represent the eigenvectors corresponding to the d largest eigenvalues of Σ_{xy} and Σ_{yx} , respectively. Then, $\mathcal{S}(\mathbf{M}_X) = \mathcal{S}(\Sigma_X^{-1}v_X)$ and $\mathcal{S}(\mathbf{M}_y) = \mathcal{S}(\Sigma_y^{-1}v_y)$.

Proof: We will prove only that $S(\mathbf{M}_x) = S(\boldsymbol{\Sigma}_x^{-1}\boldsymbol{\nu}_x)$, because, by applying the same arguments, $S(\mathbf{M}_y) = S(\boldsymbol{\Sigma}_y^{-1}\boldsymbol{\nu}_y)$ is established. By relation (5), rank $(\boldsymbol{\Sigma}_x^{-1}\boldsymbol{\Sigma}_{xy}) = d$, which directly implies that rank $(\boldsymbol{\Sigma}_{yx}) = d$. Therefore, we have $S(\boldsymbol{\Sigma}_{yx}) = S(\boldsymbol{\nu}_x)$ by the construction of $\boldsymbol{\nu}_x$. Then we have $S(\mathbf{M}_x) = \boldsymbol{\Sigma}_x^{-1}S(\boldsymbol{\Sigma}_{yx}) = \boldsymbol{\Sigma}_x^{-1}S(\boldsymbol{\Sigma}_{yx}) = \boldsymbol{\Sigma}_x^{-1}S(\boldsymbol{\Sigma}_{yx}) = \boldsymbol{\Sigma}_x^{-1}S(\boldsymbol{\Sigma}_{yx}) = \boldsymbol{\Sigma}_x^{-1}S(\boldsymbol{\Sigma}_{yx})$.

Denote the two sets of the eigenvectors corresponding to the *m* largest eigenvalues of $\hat{\Sigma}_{xy}$ and $\hat{\Sigma}_{yx}$ as $\hat{\nu}_x \in \mathbb{R}^{p \times m}$ and $\hat{\nu}_y \in \mathbb{R}^{r \times m}$, respectively. Through the initialized CCA with $\hat{\nu}_x$ and $\hat{\nu}_y$ as seed matrices, $\hat{\mathbf{M}}_{x,0}$ and $\hat{\mathbf{M}}_{y,0}$ are constructed, and $\hat{\mathbf{M}}_{x,0}^{\mathsf{T}}\mathbf{X}$ and $\hat{\mathbf{M}}_{y,0}^{\mathsf{T}}\mathbf{Y}$ replace the original sets of \mathbf{X} and \mathbf{Y} . Again, the CCA reduction is completed by the finalized CCA with $\hat{\mathbf{M}}_{x,0}^{\mathsf{T}}\mathbf{X}$ and $\hat{\mathbf{M}}_{y,0}^{\mathsf{T}}\mathbf{Y}$.

Here, it is important to choose *m* large enough to have $S(\hat{\Sigma}_{yx}) = S(\hat{v}_x)$ and $S(\hat{\Sigma}_{xy}) = S(\hat{v}_y)$ but small enough to estimate \mathbf{M}_x and \mathbf{M}_y more accurately and to have numerical stability. To select a proper value of *m* without avoiding complexity, we adopt two simple and widely accepted ways in practice. One is graphical determination by a scree plot for eigenvalue of $\hat{\Sigma}_{xy}$, and the other is the number of eigenvalues whose sum is to cover 60% or above of the total variations of $\hat{\Sigma}_{xy}$.

4. NUMERICAL STUDIES AND DATA ANALYSIS

4.1. Numerical studies

To confirm the effectiveness of the proposed method, three simulation examples are presented. In all examples, we consider n = 100, and the number of iteration is always 100. The termination index u^* is fixed at 3 for simplicity.

Example 1

The first example is the case that $\min(p, r)$ is fairly small. Define η to have its first and last of 20% of elements equal to $(0.4)^{-1/2}$ and all other elements are equal to 0. The first set **X** of *p* variables was generated from $MN(\mu, \Sigma)$, where *MN* stands for multivariate normal distribution, $\mu = 0$ and $\Sigma = \operatorname{diag}(\lambda)$. In the example, we considered either 10 or 100 for *p* and two choices for λ : $\lambda_1 = (2/3)\operatorname{diag}(2, \ldots, 2, 1, \ldots, 1)^T$ with equal multiplicity between 1 and 2 and $\lambda_2 = (2/p + 1)(p, (p - 1), \ldots, 1)^T$.

The second set **Y** was constructed in three different cases: (1) $Y_1 = \eta^T \mathbf{X} + \varepsilon_1$ and $Y_2 = \eta^T \mathbf{X} + \varepsilon_2$; (2) $Y_1 = \eta^T \mathbf{X} + \varepsilon_1$ and $Y_2 = (\eta^T \mathbf{X})^3 + \varepsilon_2$; (3) $Y_1 = (\eta^T \mathbf{X})^3 + \varepsilon_1$ and $Y_2 = (\eta^T \mathbf{X})^3 + \varepsilon_2$, where $\varepsilon_i \stackrel{i.i.d}{\sim} N(0, 1) \perp \mathbf{X}$ for i = 1, 2.

Because min(*p*, *r*) = 2 in the example, the estimation of \mathbf{M}_x was focused rather than that of \mathbf{M}_y . Then \mathbf{M}_x was estimated by the first case of the seeded CCA with $\hat{\boldsymbol{\Sigma}}_{xy} \in \mathbb{R}^{p \times 2}$, which resulted in $\hat{\mathbf{M}}_x \in \mathbb{R}^{p \times 2}$. To measure the accuracy of the estimation of \mathbf{M}_x , the averages of |r|s were computed. The notation |r| stands for the absolute value of R^2 from a regression of $\eta^T \mathbf{X} | \hat{\mathbf{M}}_x^T \mathbf{X}$. If $\hat{\mathbf{M}}_x$ estimates η well, the averages should be close to one. The results are presented in Table I.

Example 2

In the example, the dimension of **X** was fixed at p = 500 and was generated from $MN(0, \Sigma)$, where $\Sigma = (2/3) \text{diag}(2, \dots, 2, 1, \dots, 1)^{\text{T}}$ with equal multiplicity between 1 and 2.

For the second set **Y**, its dimension *r* varied among 10, 20, 50, 100, and 200. Define that $\mu_1 = 1 + X_1$ and $\mu_2 = 2 + X_2$. Then **Y** was generated as follows: $Y_1 = \mu_1 + \varepsilon_1$; $Y_2 = \mu_2 + \varepsilon_2$; $Y_3 = \mu_1 + \mu_2 + \varepsilon_3$; $Y_4 = \mu_1 - \mu_2 + \varepsilon_4$; $Y_r = \varepsilon_r$ for $r \ge 5$, where $\varepsilon = (\varepsilon_1, \dots, \varepsilon_r)^T$ was independently generated from the standard normal distribution, and $\varepsilon \perp \mathbf{X}$.

In the example, **X** is associated with **Y** only through X_1 and X_2 . Therefore, **M**_X is equal to $\eta = \{(1, 0, ..., 0)^T, (0, 1, 0, ..., 0)^T\}$. On the other hand, **Y** is associated with **X** only through μ_1 and μ_2 . Hence, **M**_Y is equal to $\phi = \{(1, 0, ..., 0)^T, (0, 1, 0, ..., 0)^T\}$.

Table I. Averages of $ r $ s computed from $\eta^T \mathbf{X} \hat{\mathbf{M}}_X^T \mathbf{X}$ for Example 1						
	λ_1		λ2			
models	<i>p</i> = 10	<i>p</i> = 500	<i>p</i> = 10	<i>p</i> = 500		
$Y_1 = \beta^{T} \mathbf{X} + \varepsilon_1, Y_2 = \beta^{T} \mathbf{X} + \varepsilon_2$	0.9789	0.8198	0.9992	0.9421		
$Y_1 = \beta^{T} \mathbf{X} + \varepsilon_1, Y_2 = (\beta^{T} \mathbf{X})^3 + \varepsilon_2$	0.9748	0.7915	0.9709	0.8025		
$Y_1 = (\beta^{T} \mathbf{X})^3 + \varepsilon_1, Y_2 = (\beta^{T} \mathbf{X})^3 + \varepsilon_2$	0.9841	0.855	0.9898	0.928		

Table II. Averages of $ r s$ for Example 2							
r	$\boldsymbol{\eta}_1^T\mathbf{X} \hat{\mathbf{M}}_x^T\mathbf{X}$	$\boldsymbol{\eta}_2^T \mathbf{X} \hat{\mathbf{M}}_x^T \mathbf{X}$	$\boldsymbol{\phi}_1^T\mathbf{Y} \hat{\mathbf{M}}_y^T\mathbf{Y}$	$\boldsymbol{\phi}_2^T \mathbf{Y} \hat{\mathbf{M}}_y^T \mathbf{Y}$			
10	0.8979	0.8901	0.8929	0.8852			
20	0.8923	0.8924	0.8817	0.8831			
50	0.8854	0.8862	0.8632	0.8637			
100	0.8714	0.8676	0.8356	0.8312			
200	0.8332	0.8274	0.7801	0.7828			

Because min(*p*, *r*) is relatively large, the matrices of $\hat{v}_x \in \mathbb{R}^{p \times 2}$ and $\hat{v}_y \in \mathbb{R}^{r \times 2}$ were used in the initialized CCA, where \hat{v}_x and \hat{v}_y were sets of the eigenvectors corresponding to the first two largest eigenvalues of $\hat{\Sigma}_{xy}$ and $\hat{\Sigma}_{yx}$, respectively.

Let $|r_i^x|$ stand for the absolute value of R^2 from a regression of $\eta_i^T X | \hat{M}_x^T X$ for i = 1, 2. To measure the accuracy of the estimation of M_x , the averages of $|r_i^x|$ s were calculated for i = 1, 2. Similarly, for M_y , we computed the averages of $|r_i^y|$ s from a regression of $\phi_j^T X | \hat{M}_y^T X$ for j = 1, 2. The results are summarized in Table II. **Example 3**

Example 2 was slightly modified by setting $\mu_1 = 1.5(5 + |X_1|)(X_2 + X_3)$ and $\mu_2 = 1 + X_1$, while the rest of the setting in Example 2 remained the same. The purpose of the example is to study how the existence of non-linear relationship in **X** affects the estimation of **M**_x and **M**_y. In the example, **X** is associated with **Y** only through X_1 and $X_2 + X_3$, while the association of **Y** with **X** remains the same as Example 2. Hence, η alone is changed to { $(1, 0, \dots, 0)^T$, $(0, 1, 1, 0, \dots, 0)^T$ }. The same statistics summarize the simulation studies, and they are reported in Table III.

According to Table I, the performances of seeded CCA seem quite good for both p = 10 and 500. And, Table II shows that the performances are all notable, although the performances get worse as r increases. Compared with Example 2, the good performances in the estimation of both \mathbf{M}_x and \mathbf{M}_y are observed.

Overall, various numerical studies including these three examples support the proposed method well, and there will be no critical issues in use in practice.

4.2. Real data application

4.2.1. Near-infrared spectroscopy of biscuit doughs data

To illustrate the seeded CCA in practice, near-infrared spectroscopy of biscuit doughs data (Brown *et al.* [11]; NIR) was adopted. Quantitative NIR spectroscopy is used to analyze diverse compositions in food, drink, pharmaceutical products, and petrochemicals. The experiment of the biscuit dough data set was conducted to measure the suitability of NIR spectroscopy for analyzing the composition of biscuit. The percentages of four components of biscuits made by standard recipe-fat, sucrose, dry

Table III. Averages of $ r $ s for Example 3							
r	$\boldsymbol{\eta}_1^T\mathbf{X} \hat{\mathbf{M}}_x^T\mathbf{X}$	$\eta_2^T \mathbf{X} \hat{\mathbf{M}}_x^T \mathbf{X}$	$\boldsymbol{\phi}_1^T\mathbf{Y} \hat{\mathbf{M}}_y^T\mathbf{Y}$	$\boldsymbol{\phi}_2^T \mathbf{Y} \hat{\mathbf{M}}_y^T \mathbf{Y}$			
10	0.8943	0.9912	0.9992	0.8927			
20	0.8918	0.9910	0.9992	0.8905			
50	0.8827	0.9911	0.9992	0.8810			
100	0.8623	0.9910	0.9992	0.8606			
200	0.8307	0.9908	0.9988	0.8341			

flour, and water were calculated from 72 biscuit samples. One set of variables, saying $\mathbf{Y} \in \mathbb{R}^4$, is composed of this four ingredients.

For each sample, there was a wavelength observed by spectroscopy and 700 different points of it had been measured, range from 1100 to 2498 nanometers (NM) at an interval of 2 nm. From this data, the first 140 and the last 49 wavelengths were removed in Brown *et al.* [11], because those figures hardly seemed to contain the useful information, and increased the interval to 4 nm. Consequently, the wavelength had its bounds from 1380 to 2400nm with 256 points, representing another set of variables $\mathbf{X} \in \mathbb{R}^{256}$.

The data were obtained from ppls-package of statistical language R and is named cookie. Because the 23th and 61st samples in the data set were suspected as outliers, they were deleted from the data set before analysis.

Due to p = 256 > n = 72, it is clear that the standard CCA is not applicable. Because r = 4 is reasonably small compared with n, the first case of the seeded CCA in Section 3.2.1 was implemented. Then $\hat{\mathbf{M}}_{x,0}$ was constructed through taking $\hat{\boldsymbol{\Sigma}}_{xy} \in \mathbb{R}^{256 \times 4}$ as a seed matrix. The proper number of projection turned out to be $u^* = 2$ with simple graphical determination regarding the change in $\hat{\mathbf{R}}_{u-1}$ (not reported).

This initialized CCA resulted in the replacement of the original 256-dimensional **X** by the four-dimensional $\hat{\mathbf{M}}_{x,0}^{\mathsf{T}} \mathbf{X}$. For notational conveniences, let $\hat{\mathbf{M}}_{x,0_i}$ and $\hat{\mathbf{M}}_{x,0_{-i}}$ stand for the *i*th column of $\hat{\mathbf{M}}_{x,0}$ and the $\hat{\mathbf{M}}_{x,0}$ after removing the *i*th column, respectively. Similarly, $\hat{\mathbf{M}}_{y,0_i}$ and $\hat{\mathbf{M}}_{y,0_{-i}}$ are defined.

Next, to complete the CCA dimension reduction, the relationship between $\hat{\mathbf{M}}_{X,0}^{\mathsf{T}} \mathbf{X}$ and \mathbf{Y} was inspected through in a scatterplot matrix in Figure 1a. The plot indicates that two pairs of $(\hat{\mathbf{M}}_{X,0,1}^{\mathsf{T}} \mathbf{X}, Y_1)$ and $\{\hat{\mathbf{M}}_{X,0,-1}^{\mathsf{T}} \mathbf{X}, (Y_2, Y_3, Y_4)\}$ have clear separation by the relations within each pair and between pairs. Along with this information, the finalized CCA was carried out with $\hat{\mathbf{M}}_{X,0}^{\mathsf{T}} \mathbf{X}$ and \mathbf{Y} , and denote $\hat{\mathbf{M}}_X \in \mathbb{R}^{4 \times 4}$ and $\hat{\mathbf{M}}_Y \in \mathbb{R}^{4 \times 4}$ as the finalized canonical coefficient matrices for \mathbf{X} and \mathbf{Y} . To withdraw more information on the relationship between $\hat{\mathbf{M}}_Y^{\mathsf{T}} \mathbf{Y}$ and \mathbf{Y} , the scatterplot matrix in Figure 1(b) was constructed. From the plot, it can be observed that the first canonical variate $\hat{\mathbf{M}}_{y_1}^{\mathsf{T}} \mathbf{Y}$ has strong linear relationships with Y_2 , Y_3 and Y_4 , while the second one $\hat{\mathbf{M}}_{y_2}^{\mathsf{T}} \mathbf{Y}$ can be thought of as Y_1 itself.



Figure 1. Scatterplot matrices for NIR data.



Figure 2. Plots for CCA in nutrimouse data.

By the discussion, the two pairs of canonical variates of $(\hat{\mathbf{M}}_{x_1}^T \mathbf{X}, \hat{\mathbf{M}}_{y_1}^T \mathbf{Y})$ and $(\hat{\mathbf{M}}_{x_2}^T \mathbf{X}, \hat{\mathbf{M}}_{y_2}^T \mathbf{Y})$ should be concluded to be sufficient to summarize the relation between \mathbf{Y} and \mathbf{X} . If one can do further analysis to investigate how the percentages of the four ingredient is changed given the wavelengths, s/he can apply multivariate linear regression of $(\hat{\mathbf{M}}_{y_1}^T \mathbf{Y}, \hat{\mathbf{M}}_{y_2}^T \mathbf{Y}) |(\hat{\mathbf{M}}_{x_1}^T \mathbf{X}, \hat{\mathbf{M}}_{x_2}^T \mathbf{X})$.

4.2.2. Nutrimouse data

For the illustration purpose of the second case in Section 3.2.2, nutrimouse data (Martin *et al.* [12]) was used. The data were collected from a nutrigenomic study in 40 mouse (n = 40) by which investigated the effects of five regimens with contrasted fatty acid compositions on liver lipids and hepatic gene expression.

The following variables were used for two sets of variables: one set $\mathbf{X} \in \mathbb{R}^{120}$ was expressions of 120 genes measured in liver cells,

acquired through microarray technology and selected among about 30 000 as potentially useful in the nutrition study. The other set $\mathbf{Y} \in \mathbb{R}^{21}$ was concentrations of 21 hepatic fatty acids (FA) measured by gas chromatography. Additionally, the 40 mouse are cross-classified based on two factors of genotype an diet: genotype: wild-type (WT) mice and PPAR*alpha* deficient mice (PPAR); diet: corn and colza oils (50/50, REF), hydrogenated coconut oil for a saturated FA diet (COC), sunflower oil for ω 6 FA-rich diet (SUN), linseed oil for ω 3-rich diet (LIN), and corn/colza/enriched fish oils (42.5/42.5/15, FISH). The data are publicly available in the mixOmics-package of statistical language R.

Because min(*p*, *r*) = 21 and was relatively larger compared with n = 40, $\hat{v}_{\chi} \in \mathbb{R}^{120 \times m}$ and $\hat{v}_{y} \in \mathbb{R}^{21 \times m}$ of the eigenvectors corresponding to the *m* largest eigenvalues $\hat{\Sigma}_{yx}$ and $\hat{\Sigma}_{xy}$ were used as the seed matrices in the initialized CCA. By inspecting the scree plot of $\hat{\Sigma}_{xy}$, reported in Figure 2(a), it was determined that m = 4.



Figure 3. Summary plots for nutrimouse data

Additionally, the sums of the first two largest and the first four largest eigenvalues of $\hat{\Sigma}_{xy}$ accounts for 90.1% and 99.4% of the total variablity, respectively. With \hat{v}_x and \hat{v}_y setting m = 4, the initialized CCA yielded $\hat{\mathbf{M}}_{x,0}^{\mathsf{T}} \mathbf{X}$ and $\hat{\mathbf{M}}_{y,0}^{\mathsf{T}} \mathbf{Y}$, where $\hat{\mathbf{M}}_{x,0} \in \mathbb{R}^{120 \times 4}$ and $\hat{\mathbf{M}}_{y,0} \in \mathbb{R}^{21 \times 4}$. Before the finalized CCA, the relation of $\hat{\mathbf{M}}_{x,0}^{\mathsf{T}} \mathbf{X}$ and $\hat{\mathbf{M}}_{y,0}^{\mathsf{T}} \mathbf{Y}$ was investigated through a scatter plot matrix in Figure 2b. The plot indicates that $(\hat{\mathbf{M}}_{x,0_3}, \hat{\mathbf{M}}_{x,0_4})^{\mathsf{T}} \mathbf{X}$ and $(\hat{\mathbf{M}}_{y,0_3}, \hat{\mathbf{M}}_{y,0_4}, \hat{\mathbf{M}}_{y,0_4})^{\mathsf{T}} \mathbf{Y}$ are redundant, and hence $\hat{\mathbf{M}}_{x,0}$ and $\hat{\mathbf{M}}_{y,0}$ were additionally reduced to $(\hat{\mathbf{M}}_{x,0_1}, \hat{\mathbf{M}}_{x,0_2}) \in \mathbb{R}^{120 \times 2}$ and $(\hat{\mathbf{M}}_{y,0_1}, \hat{\mathbf{M}}_{y,0_2}) \in \mathbb{R}^{21 \times 2}$. So, it was concluded that two pairs of the canonical covariates could represent the relation of \mathbf{X} and \mathbf{Y} . At last, to have $\hat{\mathbf{M}}_x^{\mathsf{T}} \mathbf{X}$ and $(\hat{\mathbf{M}}_{x,0_1}, \hat{\mathbf{M}}_{x,0_2})^{\mathsf{T}} \mathbf{X}$ and $(\hat{\mathbf{M}}_{y,0_1}, \hat{\mathbf{M}}_{y,0_2})^{\mathsf{T}} \mathbf{X}$.

In Figure 3(a)–(d), the scatter plots of the first and second pairs of the finalized canonical variates marked by genotype and diet of mouse used are reported as summary. According to Figure 3(a) and (c), the first canonical covariate is distinguished better by genotype than the second one, while the latter is better grouped for diet than the former in Figures 3(b) and (d).

5. Conclusion

Canonical correlation analysis is one of popular statistical methodologies in multivariate analysis, especially, in studying relation of two sets of variables. However, when sample sizes are smaller than the maximum of the dimensions of two sets of variables, it is plausible to implement canonical correlation analysis, because the sample covariance matrices are not invertible in practice.

In this article, by adopting seeded dimension reduction, a twostep canonical correlation procedure is proposed and is called seeded canonical correlation analysis. In the first step, the dimensions of the original two sets are initially reduced without losing information on the canonical correlation analysis perspective. Then, in the second step, the standard canonical correlation application with the initially reduced sets is carried out to complete the dimension reduction. Numerical studies confirm the approach, and two real data analyses are presented for illustration of the purpose.

It is believed that our works can make canonical correlation analysis more fruitful as statistical methodologies in high-throughput data analysis. The codes for the seeded canonical correlation analysis are publicly available on the following website of the authors:

http://home.ewha.ac.kr/~yjkstat/publication.html.

Acknowledgements

The authors are also grateful to the associate editor and the referee for many insightful and helpful comments. The corresponding author, Jae Keun Yoo appreciates the Department of Statistics, University of Washington to wrap this research up during the visit in summer, 2013. For Jae Keun Yoo (the corresponding author), this work was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Korean Ministry of Education (NRF-2014R1A2A1A11049389). For Yunju Im and Heyin Gang, this work was supported by the BK21 Plus Project through the National Research Foundation of Korea (NRF) funded of Korea (NRF) funded by the BK21 Plus Project through the National Research Foundation of Korea (NRF) funded by the Korean Ministry of Education (22A20130011003). Currently, the author, Yunju Im is a PhD candidate of the graduate program at the Department of Statistics and Actuarial Science, University of Iowa.

REFERENCES

- 1. Johnson RA, Wichern DW. Applied Multivariate Statistical Analysis (6th edn). Pearson Prentice Hall: New Jersey, USA, 2007, 539–574.
- 2. Parkhomenko E, Tritchler D, Beyene J. Genome-wide sparse canonical correlation of gene expression with genotypes. *BMC Proc.* 2007; 1: S119.

- Waaijenborg S, Verselewel de Witt Hamer P, Zwinderman A. Quantifying the association between gene expressions and DNA-markers by penalized canonical correlation analysis. *Stat. Appl. Genet. Mol.* 2008; 7: 3.
- Zou H, Hastie T. Regularization and variable selection via the elastic net. J. R. Statist. Soc. B 2005; 67: 301–320.
- 5. Le Cao K, Pascal M, Robert-Granie C, Philippe B. Sparse canonical methods for biological data integration: application to a cross-platform study. *BMC Bioinformatics* 2009; **10**: 34.
- 6. Witten D, Tibshirani R, Hastie T. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics* 2009; **10**: 515–534.
- Cook RD, Li B, Chiaromonte F. Dimension reduction in regression without matrix inversion. *Biometrika* 2007; 94: 569–584.
- 8. Searle A. Matrix algebra useful for statistics. Wiley: New York, USA, 1982, 212–213.
- 9. Yoo JK. Advances in seeded dimension reduction: bootstrap criteria and extensions. *Comput. Stat. Data. An.* 2013; **60**: 70–79.
- 10. Lee K, Yoo JK. Canonical correlation analysis through linear modeling. *Aust. Nz. J. Stat.* 2014; **56**: 59–72.
- 11. Brown PJ, Fearn T, Vannucci M. Bayesian wavelet regression on curves with applications to a spectroscopic calibration problem. *J. Am. Stat. Assoc.* 2001; **96**: 398–408.
- Martin PGP, Guillou H, Lasserre F, Djean S, Lan A, Pascussi J-M, San Cristobal M, Legrand P, Besse P, Pineau T. Novel aspects of PPAR-mediated regulation of lipid and xenobiotic metabolism revealed through a multrigenomic study. *Hepatology* 2007; 54: 767–777.

Supporting information

Additional supporting information may be found in the online version of this article at the publisher's web-site.