

Regular Paper

# Designing Various Multivariate Analysis at Will via Generalized Pairwise Expression\*

AKISATO KIMURA<sup>1,a)</sup> MASASHI SUGIYAMA<sup>2</sup> HITOSHI SAKANO<sup>1</sup> HIROKAZU KAMEOKA<sup>3,1</sup>

Received: August 22, 2012, Revised: October 12, 2012,  
Accepted: November 26, 2012

**Abstract:** It is well known that dimensionality reduction based on multivariate analysis methods and their kernelized extensions can be formulated as generalized eigenvalue problems of scatter matrices, Gram matrices or their augmented matrices. This paper provides a generic and theoretical framework of multivariate analysis introducing a new expression for scatter matrices and Gram matrices, called Generalized Pairwise Expression (GPE). This expression is quite compact but highly powerful. The framework includes not only (1) the traditional multivariate analysis methods but also (2) several regularization techniques, (3) localization techniques, (4) clustering methods based on generalized eigenvalue problems, and (5) their semi-supervised extensions. This paper also presents a methodology for designing a desired multivariate analysis method from the proposed framework. The methodology is quite simple: adopting the above mentioned special cases as templates, and generating a new method by combining these templates appropriately. Through this methodology, we can freely design various tailor-made methods for specific purposes or domains.

**Keywords:** multivariate analysis, dimensionality reduction, generalized eigenvalue problem, pairwise expression, kernel method, clustering, semi-supervised learning, regularization

## 1. Introduction

We can easily obtain a massive collection of texts (long articles<sup>\*1</sup>, microblogs<sup>\*2</sup>), images [1], [2], [3], [4], videos [5], [6] and musics [7]<sup>\*3</sup> nowadays. However, we are now facing a difficulty in finding an intrinsic trend and nature of such a massive collection of data. Multivariate analysis [8] is traditional, quite simple but might be one of the powerful tools to obtain a hidden structure embedded in the data, via classification, regression and clustering [9], [10]. Actually, multivariate analysis has been still an important tool, and recent reports showed its effectiveness for several applications, e.g., human detection [11], image annotation [12], [13], and sensor data mining [14], [15], [16], [17].

Principal component analysis (PCA) [18], Fisher discriminant analysis (FDA) [19], multivariate linear regression (MLR), canonical correlation analysis (CCA) [18], and partial least squares (PLS) [20] are well known as standard multivariate analysis methods. These methods can be formulated as a generalized eigenvalue problem of a scatter matrix or an augmented matrix composed of several scatter matrices. Several extended researches tried to tackle the small sample size problem [21], i.e., the situation where the number of training samples is small compared with their dimensionality (e.g., robust PCA [22], [23], [24], [25] and robust FDA [26], [27], [28]). Kernelized extensions of those standard methods have been also developed to

deal with non-vector samples and non-linear analysis (e.g., kernel PCA [29], kernel FDA [30], [31], [32], kernel MLR [33], and kernel CCA [34], [35]). They can be formulated as a generalized eigenvalue problem of an augmented matrix composed of Gram matrices, instead of scatter matrices. Kernel multivariate analysis often needs some regularization techniques such as  $\ell_2$ -norm regularization [36], [37], [38] to inhibit overfitting and the graph Laplacian method [39] to fit underlying data manifolds smoothly. In addition, improvements of robustness against outliers and non-Gaussianity (i.e., multi-dimensional scaling (MDS) [40], locality preserving projection (LPP) [41] and local Fisher discriminant analysis (LFDA) [42]) and their extensions to semi-supervised dimensionality reduction [39], [43], [44] have been considered.

In addition, a lot of multivariate analysis methods and several trials to unify these methods have been presented so far. Borge et al. [45] and De Bie et al. [46] showed that several major linear multivariate analysis method can be formulated by a unified form of generalized eigenvalue problems by introducing the augmented matrix expression. Sun et al. [47], [48] showed the equivalence between a certain class of generalized eigenvalue problems and least squares problems under a mild assumption. De la Torre [49], [50] further extended their work to a various kind of component analysis methods by introducing the formulation of least-squares weighted kernel reduced rank regression (LS-WKRRR). However, freely designing a tailor-made multi-

<sup>1</sup> NTT Communication Science Laboratories, NTT Corporation, Soraku, Kyoto 619-0237, Japan

<sup>2</sup> Graduate School of Information Science and Engineering, Tokyo Institute of Technology, Meguro, Tokyo 152-8552, Japan

<sup>3</sup> Graduate School of Information Science and Technologies, the University of Tokyo, Bunkyo, Tokyo 113-8656, Japan

<sup>a)</sup> akisato@ieee.org

\* A preliminary version of this paper was previously presented in Conference on International Association for Pattern Recognition (ICPR2012).

<sup>\*1</sup> New York Times Article Archive:

<http://www.nytimes.com/ref/membercenter/nytarchive.html>

<sup>\*2</sup> Tweets2011 corpus for TREC2011 microblog track:

<http://trec.nist.gov/data/tweets/>

<sup>\*3</sup> Last.fm: <http://www.lastfm.jp>, Freesound: <http://www.freesound.org>

variate analysis for a specific purpose or domain still remains an open problem. Until now, researchers and engineers have had to choose one of the existing methods that seems best to address the problem of interest, or had to laboriously develop a new analysis method tailored specifically for that purpose.

In view of the above discussions, this paper provides a new expression of second-order statistics including covariance matrices and Gram matrices, which we call the *Generalized Pairwise Expression (GPE)*. GPE is originated in Pairwise Expression (PE) [42], [43], [51] that tries to describe the relation between pairs of samples regarding whether pairs are close together or far apart. Through the PE framework, we can obtain an interpretation of a multivariate analysis method as to how it works. Our GPE framework extends the PE framework to larger classes of multivariate analysis methods, and it can be also diverted to designing a new multivariate analysis method. Our main contributions of this paper can be summarized as follows:

- (1) GPE makes it easy to design a new multivariate analysis method with desired properties without any special knowledge of multivariate analysis.
- (2) The methodology is quite simple: Exploiting the above mentioned existing methods as templates, and constructing a new method by combining these templates appropriately. This property has not been discussed yet in any previous researches to our best knowledge.
- (3) It is also possible to individually select and arrange samples for calculating the second-order statistics of the methods to be combined, which enables us to extend multivariate analysis methods to semi-supervised ones and multi-modal ones, where some parts are calculated from only labeled samples, and the other parts are obtained from both labeled and unlabeled samples.

The rest of this paper is organized as follows: Section 2 defines a class of multivariate analysis methods we are concerned with in this paper. Next, Section 3 describes our proposed framework, GPE, and its fundamental properties. These properties provide a methodology to design multivariate analysis methods with desired characteristics. Then, Section 4 reviews major multivariate analysis methods from the viewpoint of GPE. This review will give us templates of the GPEs for designing desired methods. After the above preparations, Section 5 demonstrates how to design a new multivariate analysis method. By replicating the methodology shown in the preceding sections, we can easily design various multivariate analysis methods at will. Additionally, Section 6 describes a non-linear and/or non-vector extension of GPE with the help of the kernel trick, which is not trivial. With this extension, non-linear dimensionality reduction, and several clustering methods are all in the class of multivariate analysis methods we are concerned with.

## 2. Multivariate Analysis for Vector Data

### 2.1 Preliminaries

Consider two sets  $X$  and  $Y$  of samples<sup>\*4</sup>, where each set contains  $N_x$  and  $N_y$  samples, and each sample can be expressed as a

<sup>\*4</sup> The following discussion can be easily extended to more than 2 sets of samples sets [52].

vector with  $d_x$  and  $d_y$  dimensions, respectively, as follows:

$$X = \{\mathbf{x}_1, \dots, \mathbf{x}_{N_x}\},$$

$$Y = \{\mathbf{y}_1, \dots, \mathbf{y}_N, \mathbf{y}_{N_x+1}, \dots, \mathbf{y}_{N_x+N_y-N}\} \quad (N \leq N_x).$$

For brevity, both of the sample sets  $X$  and  $Y$  are supposed to be centered on the origin by subtracting the mean from each component. Suppose that samples  $\mathbf{x}_n$  and  $\mathbf{y}_n$  with the same suffix are co-occurring. Each set  $X$  and  $Y$  of samples is separated into the following two types: *Complete sample sets*  $X^{(C)}$  and  $Y^{(C)}$  so that every sample  $\mathbf{x}_n$  (resp.  $\mathbf{y}_n$ ) has co-occurring sample  $\mathbf{y}_n$  (resp.  $\mathbf{x}_n$ ), and *incomplete sample sets*  $X^{(I)}$  and  $Y^{(I)}$  so that every sample  $\mathbf{x}_n$  (resp.  $\mathbf{y}_n$ ) cannot find the co-occurring sample.

$$X^{(C)} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\},$$

$$= \{\mathbf{x}_1^{(C)}, \mathbf{x}_2^{(C)}, \dots, \mathbf{x}_N^{(C)}\},$$

$$Y^{(C)} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\},$$

$$= \{\mathbf{y}_1^{(C)}, \mathbf{y}_2^{(C)}, \dots, \mathbf{y}_N^{(C)}\},$$

$$X^{(I)} = \{\mathbf{x}_{N_x+1}, \mathbf{x}_{N_x+2}, \dots, \mathbf{x}_{N_x}\},$$

$$= \{\mathbf{x}_1^{(I)}, \mathbf{x}_2^{(I)}, \dots, \mathbf{x}_{N_x-N}^{(I)}\},$$

$$Y^{(I)} = \{\mathbf{y}_{N_x+1}, \mathbf{y}_{N_x+2}, \dots, \mathbf{y}_{N_x+N_y-N}\},$$

$$= \{\mathbf{y}_1^{(I)}, \mathbf{y}_2^{(I)}, \dots, \mathbf{y}_{N_y-N}^{(I)}\}$$

First, we concentrate on the case that  $N_x = N_y = N$ , namely all the samples are paired, unless otherwise stated.

### 2.2 Formulation

Many linear multivariate analysis methods developed so far involve an optimization problem of the following form for a  $d$ -dimensional vector  $\mathbf{w}$ :

$$\mathbf{w}^{(\text{opt})} = \arg \max_{\mathbf{w} \in \mathcal{R}^d} R(\mathbf{w}), \tag{1}$$

$$R(\mathbf{w}) = \mathbf{w}^\top \overline{\mathbf{C}} \mathbf{w} (\mathbf{w}^\top \underline{\mathbf{C}} \mathbf{w})^{-1},$$

where  $\overline{\mathbf{C}}$  and  $\underline{\mathbf{C}}$  are square matrices with certain statistical nature. For example,  $\overline{\mathbf{C}}$  is a scatter matrix of  $X$  and  $\underline{\mathbf{C}}$  is an identity matrix in PCA, and  $\overline{\mathbf{C}}$  is a between-class scatter matrix and  $\underline{\mathbf{C}}$  is a within-class scatter matrix in FDA. Roughly speaking,  $\overline{\mathbf{C}}$  encodes the quantity that we want to increase, and  $\underline{\mathbf{C}}$  corresponds to the quantity that we want to decrease. The denominator of the function  $R(\mathbf{w})$  is often normalized to remove scale ambiguity, resulting in the following form:

$$\mathbf{w}^{(\text{opt})} = \arg \max_{\mathbf{w} \in \mathcal{R}^d} R_1(\mathbf{w}) \text{ s.t. } R_2(\mathbf{w}) = 1, \tag{2}$$

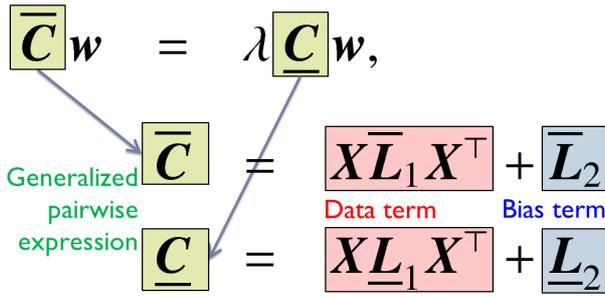
$$R_1(\mathbf{w}) = \mathbf{w}^\top \overline{\mathbf{C}} \mathbf{w}, \quad R_2(\mathbf{w}) = \mathbf{w}^\top \underline{\mathbf{C}} \mathbf{w}.$$

The above optimization problem can be converted to the following generalized eigenvalue problem via the Lagrange multiplier method:

$$\overline{\mathbf{C}} \mathbf{w} = \lambda \underline{\mathbf{C}} \mathbf{w}. \tag{3}$$

The solution  $\overline{\mathbf{w}}_k$  ( $k = 1, 2, \dots, r$ ) of the above generalized eigenvalue problem gives a solution of the original multivariate analysis formulated in Eq. (1).

It can be confirmed that Eq. (1) is invariant against any kinds of linear transformations, i.e., a vector  $\mathbf{U}\mathbf{w}^{(\text{opt})}$  transformed by any



**Fig. 1** Various multivariate analysis methods can be described via generalized pairwise expression (GPE).

$r$ -dimensional unitary matrix  $U$  is also a global solution. This implies that the range of the embedding space can be uniquely determined by Eq. (1), but the metric in the embedding space is arbitrary. A practically useful heuristic is to set

$$U = \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_d}), \tag{4}$$

where  $\text{diag}(a, b, \dots, c)$  denotes the diagonal matrix with the diagonal elements  $a, b, \dots, c$ , and  $\{\lambda_k\}_{k=1}^r$  denotes the generalized eigenvalues. Finally, we obtain the solution as

$$W^{(\text{opt})} = \{\sqrt{\lambda_1}w_1, \sqrt{\lambda_2}w_2, \dots, \sqrt{\lambda_r}w_r\}. \tag{5}$$

Thus, the minor eigenvectors are de-emphasized according to the square root of the eigenvalues.

### 3. Generalized Pairwise Expression

#### 3.1 Definition

When addressing linear multivariate analysis methods, we often deal with the following type of second-order statistics [42], [51] as an extension of scatter matrices, since it is convenient to describe the relation between two features regarding whether they are close together or far apart:

$$S_{Q,xx} = \sum_{n=1}^N \sum_{m=1}^N Q_{n,m} (\mathbf{x}_n - \mathbf{x}_m)(\mathbf{x}_n - \mathbf{x}_m)^T,$$

where  $Q$  is an  $N \times N$  non-negative, positive semi-definite and symmetric matrix<sup>\*5</sup>. A typical example is the scatter matrix:

$$S_{xx} = N^{-1} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T.$$

Let  $D_Q$  be the  $N \times N$  diagonal matrix with

$$D_{Q,n,n} = \sum_{n_2=1}^N Q_{n,n_2},$$

and let  $L_Q$  be  $L_Q = D_Q - Q$ . Then, the matrix  $S_{Q,xx}$  can be expressed in terms of  $L_Q$  as follows:

$$S_{Q,xx} = X L_Q X^T.$$

The above expression is called the *pairwise expression (PE)* of the second-order statistics  $S_{Q,xx}$ . If  $Q$  is a weight matrix for a graph with  $N$  nodes,  $L_Q$  can be regarded as a graph Laplacian matrix in the spectral graph theory. If  $Q$  is symmetric and its elements are all non-negative,  $L_Q$  is known to be positive semi-definite.

<sup>\*5</sup> When dealing with 2 sample sets in this framework, it is sufficient to introduce a concatenated sample set  $Z = (X^T, Y^T)^T$ .

Here, we extend PE to the following expression introducing an additional matrix independent of  $Q$ :

$$\hat{S}_{Q,xx} = X L_{Q,1} X^T + L_2,$$

where  $L_{Q,1}$  is an  $N \times N$  positive semi-definite matrix, and  $L_2$  is a  $d_x \times d_x$  non-negative symmetric matrix. We do not have to explicitly consider the matrix  $Q$  for the following discussions:

$$\hat{S}_{xx} = X L_1 X^T + L_2. \tag{6}$$

After all, we call this expression as the *generalized pairwise expression (GPE)* (See **Fig. 1**). The first term of Eq. (6) is called the *data term* since it depends on the sample data, and the second term is called the *bias term*.

#### 3.2 Properties

We can derive the following fundamental properties of GPE from the definition, if the number of samples,  $N$ , is sufficiently large:

- (1) If  $A$  is GPE and  $\beta > 0$  is a constant, then  $\beta A$  is also GPE.
- (2) If both  $A$  and  $B$  are GPE, then  $A + B$  is also GPE.
- (3) If both  $A$  and  $B$  are GPE, then  $AB$  is also GPE.

*Proof.* The first and second claims can be easily proved, so we concentrate on proving the third one.

First, let us denote  $A$  and  $B$  as follows:

$$A = X L_{A1} X^T + L_{A2},$$

$$B = Y L_{B1} X^T + L_{B2},$$

where  $L_{A1}$  and  $L_{B1}$  are positive semi-definite  $N \times N$  matrices, and  $L_{A2}$  and  $L_{B2}$  are  $d_x \times d_x$  non-negative symmetric matrices. Then, we obtain

$$\begin{aligned} AB &= (X L_{A1} X^T + L_{A2})(X L_{B1} X^T + L_{B2}), \\ &= X(L_{A1} X^T X L_{B1}) X^T \\ &\quad + L_{A2} X L_{B1} X^T + X L_{A1} X^T L_{B2} + L_{A2} L_{B2}. \end{aligned}$$

We can find some matrices  $L_{Ci}$  ( $i = 1, 2, 3$ ) satisfying the following relationships, if  $N \geq d_x$ :

$$L_{C1} = L_{A1} X^T X L_{B1},$$

$$L_{C2} X^T = L_{A1} X^T L_{B2},$$

$$X L_{C3} = L_{A2} X L_{B1}.$$

Here, we will show that those matrices  $L_{Ci}$  ( $i = 1, 2, 3$ ) are all positive semi-definite. For any matrix  $X \in \mathcal{R}^{d_x \times N}$ , we have

$$X L_{C1} X^T = (X L_{A1} X^T)(X L_{B1} X^T),$$

$$X L_{C2} X^T = (X L_{A1} X^T) L_{B2},$$

$$X L_{C3} X^T = L_{A2} (X L_{B1} X^T).$$

Recalling that  $L_{A1}$  and  $L_{B1}$  are positive semi-definite and  $L_{A2}$  and  $L_{B2}$  are non-negative symmetric, we can see that all the above matrix product  $X L_{Ci} X^T$  ( $i = 1, 2, 3$ ) are non-negative, which means that the matrices  $L_{Ci}$  ( $i = 1, 2, 3$ ) are all positive semi-definite. This implies that

$$\begin{aligned} AB &= X(L_{C1} + L_{C2} + L_{C3}) X^T + L_{A2} L_{B2} \end{aligned}$$

$$= XL_{D1}X^T + L_{D2},$$

for a positive semi-definite matrix  $L_{D1}$  and a non-negative symmetric matrix  $L_{D2}$ , which means  $AB$  is also GPE.  $\square$

These fundamental properties of GPE provide us a promising way to design various multivariate analysis methods very easily, namely with addition, weighting and multiplication of GPEs of existing methods with desired characteristics. The rest of the problem is to reveal GPE of existing methods and the function of every type of combinations (addition and/or multiplication), which will be described in the next section.

## 4. Reviewing Multivariate Analysis

### 4.1 Preliminaries

This section reviews major multivariate analysis methods from the viewpoint of GPE. As shown in Sections 2 and 3, the GPEs of PCA and FDA respectively are given by

$$\begin{aligned} \underline{C}^{(PCA)} &= S_{xx}, \quad \underline{C}^{(PCA)} = I_{d_x}, \\ \underline{C}^{(FDA)} &= S_{xx}^{(b)}, \quad \underline{C}^{(FDA)} = S_{xx}^{(w)}, \end{aligned}$$

where  $S_{xx}^{(b)}$  and  $S_{xx}^{(w)}$  are respectively between-class and within-class scatter matrices of  $X$ . From these examples, a scatter matrix  $S_{xx}$  is a typical example of the data term in GPE, and an identity matrix  $I_d$  is a typical example of the bias term. Note that unlike FDA,  $\underline{C}^{(PCA)}$  does not have a PE since  $I_{d_x}$  cannot be expressed in a pairwise form. This indicates the significance of introducing GPE when reviewing various multivariate analysis methods within a unified framework.

### 4.2 Canonical Correlation Analysis (CCA)

Canonical correlation analysis (CCA) [18] is a method of correlating linear relationships between two sample sets. Formally, CCA finds a new coordinate  $(w_x, w_y)$  to maximize the correlation between the two vectors in the new coordinates. In other words, the function  $\rho(w_x, w_y|X, Y)$  to be maximized is

$$\begin{aligned} \rho^{(CCA)}(w_x, w_y|X, Y) &= \frac{\langle X^T w_x, Y^T w_y \rangle}{\|X^T w_x\| \cdot \|Y^T w_y\|} \\ &= \max_{(w_x, w_y)} \frac{\widehat{E}[\langle w_x, x \rangle \langle w_y, y \rangle]}{\sqrt{\widehat{E}[\langle w_x, x \rangle^2] \cdot \widehat{E}[\langle w_y, y \rangle^2]}} \\ &= \max_{(w_x, w_y)} \frac{w_x^T \widehat{E}[xy^T] w_y}{\sqrt{w_x^T \widehat{E}[xx^T] w_x w_y^T \widehat{E}[yy^T] w_y}} \\ &= \frac{w_x^T S_{xy} w_y}{\sqrt{w_x^T S_{xx} w_x w_y^T S_{yy} w_y}}, \end{aligned} \tag{7}$$

where  $\widehat{E}[\cdot]$  denotes an empirical expectation. The maximum of the function  $\rho(X^{(C)}, Y^{(C)})$  is not affected by re-scaling  $w_x$  and  $w_y$  either together or independently. Therefore, the maximization of  $\rho(X^{(C)}, Y^{(C)})$  is equivalent to maximizing the numerator of  $\rho(X^{(C)}, Y^{(C)})$  subject to

$$w_x^T S_{xx} w_x = w_y^T S_{yy} w_y = 1.$$

Taking derivatives of the corresponding Lagrangian with respect

to  $w_x$  and  $w_y$ , we obtain

$$\begin{aligned} S_{xy} w_y - \lambda S_{xx} w_x &= \mathbf{0}, \\ S_{yx} w_x - \lambda S_{yy} w_y &= \mathbf{0}, \end{aligned}$$

where  $\lambda$  is a Lagrange multiplier.

From the above discussion, the GPE of CCA can be obtained as follows:

$$\begin{aligned} \underline{C}^{(CCA)} &= \begin{pmatrix} \mathbf{0} & S_{xy} \\ S_{yx} & \mathbf{0} \end{pmatrix}, \quad \underline{C}^{(CCA)} = \begin{pmatrix} S_{xx} & \mathbf{0} \\ \mathbf{0} & S_{yy} \end{pmatrix}, \\ w &= (w_x^T, w_y^T)^T. \end{aligned}$$

We additionally note that when every sample  $y_n$  in  $Y$  represents a class indicator vectors, namely  $y_n \in \{0, 1\}^M$ ,  $\sum_{m=1}^M y_{n,m} = 1$  and  $M$  is the number of classes, CCA is reduced to FDA [53]<sup>\*6</sup>. Thus, CCA can be regarded as a generalized variant of FDA so that each sample can belong to multiple classes.

### 4.3 Multiple Linear Regression (MLR)

Multiple linear regression (MLR) is a method of finding a projection matrix  $W$  with the minimum squared error between  $y$  and its linear approximation  $Wx$ . For simplicity, we first consider the case that the projection matrix  $W$  is with rank 1, which can be written as a direct product of two bases  $w_x$  and  $w_y$ . This assumption is useful to understand MLR from the viewpoint of GPE. Then, the objective function to be minimized is the following squared error:

$$\begin{aligned} \epsilon^{(MLR)}(w_x, w_y|X, Y) &= \widehat{E}[\|y - \alpha w_y w_x^T x\|^2] \\ &= \widehat{E}[y^T y] - 2\alpha w_y^T \widehat{E}[y^T x] w_x + \alpha^2 w_x^T \widehat{E}[x^T x] w_x \\ &= \widehat{E}[y^T y] - 2\alpha w_y^T S_{xy}^T w_x + \alpha^2 w_x^T S_{xx} w_x. \end{aligned}$$

To get an expression for  $\alpha$ , we calculate the derivative

$$\begin{aligned} \frac{\partial}{\partial \alpha} \epsilon^{(MLR)}(w_x, w_y|X, Y) &= 2(\alpha w_x^T S_{xx} w_x - w_y^T S_{xy}^T w_x) = 0, \end{aligned}$$

which gives

$$\alpha = (w_y^T S_{xy}^T w_x)(w_x^T S_{xx} w_x)^{-1}.$$

Then, we obtain

$$\begin{aligned} \epsilon^{(MLR)}(w_x, w_y|X, Y) &= E[y^T y] - \frac{(w_y^T S_{xy}^T w_x)^2}{w_x^T S_{xx} w_x}. \end{aligned} \tag{8}$$

Since the squared error cannot be negative and the first term of the objective function is independent of the two directions  $w_x$  and  $w_y$ , we can minimize it by maximizing the following generalized Rayleigh quotient:

$$\rho^{(MLR)}(w_x, w_y|X, Y) = \frac{w_x^T S_{xy} w_y}{\sqrt{w_x^T S_{xx} w_x w_y^T w_y}},$$

<sup>\*6</sup> Note that the technical report [53] includes several mistakes in the discussion as to the equivalence between CCA and FDA.

where  $w_x$  and  $w_y$  are supposed to be normalized as  $w_x^T S_{xx} w_x = 1$  and  $w_y^T w_y = 1$ . By comparing the above equation and Eq. (7) and the objective function for CCA, we can see that MLR is a special case of CCA, and

$$\bar{C}^{(MLR)} = \begin{pmatrix} \mathbf{0} & S_{xy} \\ S_{yx} & \mathbf{0} \end{pmatrix}, \quad \underline{C}^{(MLR)} = \begin{pmatrix} S_{xx} & \mathbf{0} \\ \mathbf{0} & I_{d_y} \end{pmatrix},$$

$$w = (w_x^T, w_y^T)^T.$$

The above derivation shows a part of the equivalence between the generalized eigenproblem and the least squares, which have been already revealed by Sun et al. [47], [48]. This equivalence property will be often exploited in the following discussions.

#### 4.4 Principal Component Regression (PCR)

Principal component regression (PCR) [54] is a variant of MLR that uses PCA when estimating regression coefficients  $W$ . It is a procedure used to overcome problems which arise when the exploratory variables are nearly co-linear. In PCR, instead of regressing the dependent variable  $y$  on the independent variables  $x$  directly, the principal components  $Vx$  of the independent variables are used. One typically only uses a subset of the principal components in the regression, making a kind of regularized estimation. Often the principal components with the highest variance are selected. A larger class of multivariate analysis methods that introduces a latent model into the standard linear regression is called latent variable regression (LVR) [55].

In the same way as MLR, we assume that the projection matrix  $W$  is with rank 1, namely  $W = w_y w_x^T$ . A rank- $K$  approximation  $\hat{X}$  of the data matrix  $X$  can be obtained by singular value decomposition as

$$\hat{X} = U_K \Sigma_K V_K^T, \tag{9}$$

where  $\Sigma_K$  is a  $K \times K$  diagonal matrix whose diagonal components are top- $K$  eigenvalues obtained by PCA of  $X$ , and  $V_K$  is a  $K \times d_x$  matrix whose columns are the top- $k$  eigenvectors. Then, the objective function of PCR to be minimized can be obtained by substituting  $\hat{X}$  into  $X$  in the objective function of MLR, as follows:

$$\begin{aligned} \epsilon^{(PCR)}(w_x, w_y | X, Y) &= \epsilon^{(MLR)}(w_x, w_y | \hat{X}, Y) \\ &= \hat{E} \left[ \|y - \alpha w_y w_x^T \hat{x}\|^2 \right] \\ &= \hat{E} \left[ y^T y \right] - 2\alpha w_y^T \hat{E} \left[ y^T \hat{x} \right] w_x + \alpha^2 w_x^T \hat{E} \left[ y^T \hat{x} \right] w_x \\ &= \hat{E} \left[ y^T y \right] - 2\alpha w_y^T S_{xy}^T w_x + \alpha^2 w_x^T S_{\hat{x}\hat{x}} w_x, \end{aligned}$$

where  $S_{\hat{x}y}$  and  $S_{\hat{x}\hat{x}}$  can be obtained as follows:

$$S_{\hat{x}y} = \frac{1}{N} \sum_{n=1}^N u_{K,n} \Sigma_K V_K^T y_n^T,$$

$$S_{\hat{x}\hat{x}} = \frac{1}{N} \sum_{n=1}^N u_{K,n} \Sigma_K V_K^T V_K \Sigma_K u_{K,n},$$

$$= \frac{1}{N} \sum_{n=1}^N u_{K,n} (\Sigma_K)^2 u_n^T.$$

From the description of the previous subsection, we can obtain

$$\bar{C}^{(PCR)} = \begin{pmatrix} \mathbf{0} & S_{\hat{x}y} \\ S_{xy}^T & \mathbf{0} \end{pmatrix}, \quad \underline{C}^{(PCR)} = \begin{pmatrix} S_{\hat{x}\hat{x}} & \mathbf{0} \\ \mathbf{0} & I_{d_y} \end{pmatrix},$$

$$w = (w_x^T, w_y^T)^T.$$

#### 4.5 Partial Least Squares (PLS)

Partial Least Squares (PLS) [20] (or sometimes called PLS regression) belongs to a family of latent variable regression (LVR), and tries to find a direction for the observable sample set  $X$  that explains the maximum variance direction for the predicted sample set  $Y$ . The contribution of PLS against the standard MLR and PCR is to simultaneously estimate the latent model and regression from the latent space to the predicted space, which leads to robust regression against noisy observations.

Although PLS cannot be formulated as a generalized eigenproblem in general, orthogonal PLS (OPLS) [56], [57] as a variant of the original PLS has a form of generalized eigenproblem. This improves the interpretability (but not the predictivity) of the original PLS. OPLS can be formulated as follows:

$$XY^T YX^T w = \lambda X X^T w,$$

$$w = (w_x^T, w_y^T)^T$$

meaning,

$$\bar{C}^{(OPLS)} = XY^T YX^T \quad \propto S_{xy} S_{xy}^T,$$

$$\underline{C}^{(OPLS)} = X X^T \quad \propto S_{xx}.$$

When every sample  $y_n$  in  $Y$  represents a class indicator vectors (cf. Section 4.2), OPLS is called OPLS-discriminant analysis (OPLS-DA) [57], which has been often used for the task of bio-marker identification [58].

#### 4.6 $\ell_2$ -norm Regularization

$\ell_2$ -norm regularization is a popular regularization technique for various optimization problems including multivariate analysis. In the area of statistics or machine learning, this is sometimes called Tikhonov regularization [9], [37]. The most popular method that utilizes  $\ell_2$ -norm regularization is ridge regression [36], which combines MLR and  $\ell_2$ -norm regularization. The objective function to be minimized is the following squared error:

$$\begin{aligned} \epsilon^{(Ridge)}(w_x, w_y | X, Y) &= \hat{E} \left[ \|y - \alpha w_y w_x^T x\|^2 \right] + \delta \|w_x\|^2 \\ &= \hat{E} \left[ y^T y \right] - 2\alpha w_y^T S_{xy}^T w_x + \alpha^2 w_x^T S_{xx} w_x + \delta \|w_x\|^2 \\ &= \hat{E} \left[ y^T y \right] - 2\alpha w_y^T S_{xy}^T w_x + \alpha^2 w_x^T (S_{xx} + \hat{\delta} I_{d_x}) w_x, \end{aligned}$$

where  $\hat{\delta} = \delta/\alpha^2$ . From the above equation and the objective function of MLR, the GPE of ridge regression can be derived as

$$\bar{C}^{(Ridge)} = \begin{pmatrix} \mathbf{0} & S_{xy} \\ S_{yx} & \mathbf{0} \end{pmatrix},$$

$$\underline{C}^{(Ridge)} = \begin{pmatrix} S_{xx} + \hat{\delta} I_{d_x} & \mathbf{0} \\ \mathbf{0} & I_{d_y} \end{pmatrix},$$

$$w = (w_x^T, w_y^T)^T.$$

In a similar way to ridge regression, we can derive the GPE of

CCA with  $\ell_2$ -norm regularization [38], [59] as

$$\bar{\mathbf{C}}^{(\text{CCA}-\ell_2)} = \begin{pmatrix} \mathbf{0} & \mathbf{S}_{xy} \\ \mathbf{S}_{yx} & \mathbf{0} \end{pmatrix},$$

$$\underline{\mathbf{C}}^{(\text{CCA}-\ell_2)} = \begin{pmatrix} \mathbf{S}_{xx} + \hat{\delta}\mathbf{I}_{d_x} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{yy} + \hat{\delta}\mathbf{I}_{d_y} \end{pmatrix}.$$

In addition, we can incorporate  $\ell_1$ -norm regularization into the GPE framework only if the objective generalized eigenproblem has the following form:

$$\mathbf{X}\mathbf{L}_Q\mathbf{X}^\top \mathbf{w} = \lambda \mathbf{X}\mathbf{X}^\top \mathbf{w},$$

meaning

$$\mathbf{S}_{Q,xx}\mathbf{w} = \lambda \mathbf{S}_{xx}\mathbf{w}.$$

PCA, FDA, MLR, CCA, OPLS and several variants can be included in this form. The details can be found in the previous work [47].

As shown in the above discussion, one of the major motivations that introduce the bias term of GPE is to integrate some regularization techniques within the framework of GPE.

#### 4.7 Locality Preserving Projection (LPP)

Locality preserving projections (LPP) [41] seeks for an embedding transformation such that nearby data pairs in the original space close in the embedding space. Thus, LPP can reduce the dimensionality without losing the local structure.

Let  $\mathbf{A}$  be an affinity matrix, that is, the  $N$ -dimensional matrix with the  $(n, m)$ -th element  $A_{n,m}$  being the affinity between  $\mathbf{x}_n$  and  $\mathbf{x}_m$ . We assume that  $A_{n,m} \in [0, 1]$ ;  $A_{n,m}$  is large if  $\mathbf{x}_n$  and  $\mathbf{x}_m$  are close and  $A_{n,m}$  is small if  $\mathbf{x}_n$  and  $\mathbf{x}_m$  are far apart. There are several different manners of defining  $\mathbf{A}$ , such as using the local scaling heuristics [60], i.e.,

$$A_{n,m} = \exp \left\{ -\frac{\|\mathbf{x}_n - \mathbf{x}_m\|^2}{\sigma_n \sigma_m} \right\},$$

$$\sigma_n = \|\mathbf{x}_n - \mathbf{x}_n^{(k)}\|,$$

where  $\mathbf{x}_n^{(k)}$  is the  $k$ -th nearest neighbor of  $\mathbf{x}_n$ . A heuristic choice of  $k = 7$  was shown to be useful through experiments [60]. The objective function to be minimized is the following weighted squared error:

$$\epsilon^{(\text{LPP})}(\mathbf{w}|\mathbf{X}) = \sum_{n=1}^N \sum_{m=1}^N A_{n,m} \|\mathbf{w}^\top \mathbf{x}_n - \mathbf{w}^\top \mathbf{x}_m\|^2$$

s.t.  $\mathbf{w}^\top \mathbf{X}\mathbf{D}_A\mathbf{X}^\top \mathbf{w} = 1.$

In the same way as the derivation of GPE (see Section 3), the above minimization can be converted to the following generalized eigenvalue problem:

$$\mathbf{X}\mathbf{L}_A\mathbf{X}^\top \mathbf{w} = \lambda \mathbf{X}\mathbf{D}_A\mathbf{X}^\top \mathbf{w}.$$

Thus, the GPE of LPP can be obtained as

$$\bar{\mathbf{C}}^{(\text{LPP})} = \mathbf{X}\mathbf{L}_A\mathbf{X}^\top, \quad \underline{\mathbf{C}}^{(\text{LPP})} = \mathbf{X}\mathbf{D}_A\mathbf{X}^\top.$$

#### 4.8 Local Fisher Discriminant Analysis (LFDA)

Local Fisher discriminant analysis (LFDA) [42] is a method for supervised dimensionality reduction, and an extension of Fisher discriminant analysis (FDA). LFDA can overcome the weakness of the original FDA against outliers. The point is the introduction of between-sample similarity matrix  $\mathbf{Q}$  obtained from the affinity matrix, for calculating the between-class scatter matrix  $\mathbf{S}_Q^{(\text{lb})}$  and the within-class scatter matrix  $\mathbf{S}_Q^{(\text{lw})}$ .

$$\mathbf{S}_Q^{(\text{lb})} = \sum_{n=1}^N \sum_{m=1}^N Q_{n,m}^{(\text{lb})} (\mathbf{x}_n - \mathbf{x}_m)(\mathbf{x}_n - \mathbf{x}_m)^\top,$$

$$\mathbf{S}_Q^{(\text{lw})} = \sum_{n=1}^N \sum_{m=1}^N Q_{n,m}^{(\text{lw})} (\mathbf{x}_n - \mathbf{x}_m)(\mathbf{x}_n - \mathbf{x}_m)^\top,$$

where  $\mathbf{Q}^{(\text{lb})}$  and  $\mathbf{Q}^{(\text{lw})}$  are the  $N \times N$  matrices with

$$Q_{n,m}^{(\text{lb})} = \begin{cases} A_{n,m}(1/N - 1/N_c) & \text{if } y_n = y_m = c, \\ 1/N & \text{if } y_n \neq y_m, \end{cases}$$

$$Q_{n,m}^{(\text{lw})} = \begin{cases} A_{n,m}/N_c & \text{if } y_n = y_m = c, \\ 1/N & \text{if } y_n \neq y_m, \end{cases}$$

and  $N_c$  is the number of samples in class  $c$ . Note that the local scaling is computed in a class-wise manner in LFDA, since we want to preserve the within-class local structure. This also contributes to reducing the computational cost for nearest neighbor search when computing the local scaling.

From the above discussion, the GPE of LFDA can be obtained as follows:

$$\bar{\mathbf{C}}^{(\text{LFDA})} = \mathbf{S}_Q^{(\text{lb})}, \quad \underline{\mathbf{C}}^{(\text{LFDA})} = \mathbf{S}_Q^{(\text{lw})}.$$

#### 4.9 Semi-supervised LFDA (SELF)

Semi-supervised local Fisher discriminant analysis, called SELF [43], integrates LFDA as a supervised dimensionality reduction and PCA as an unsupervised dimensionality reduction. SELF brings us one example for designing multivariate analysis methods via the GPE framework from the following two viewpoints:

- (1) combining several multivariate analysis methods via GPE by following the properties shown in Section 3.2,
- (2) changing sample sets to calculate the data term in GPE, which provides us to extend the method to a semi-supervised one.

Assume that there are two samples sets  $\mathbf{X}$  and  $\mathbf{Y}$ , each sample in  $\mathbf{Y}$  represents a class indicator vector, and an incomplete sample set  $\mathbf{X}^{(l)}$  only exists, namely there are at least one unlabeled samples in the sample set  $\mathbf{X}$ . In such cases, we can search for solutions that lie in the span of the larger sample set  $\mathbf{X}$ , and regularize the solution using the additional data. SELF looks for solutions that lie along an empirical estimate of the subspace spanned by all the samples. This gives increased robustness to the algorithm, and increases class separability in the absence of label information. In detail, SELF integrates the GPE ( $\mathbf{S}_Q^{(\text{C,lb})}$  and  $\mathbf{S}_Q^{(\text{C,lw})}$ ) of LFDA calculated only from the labeled samples (in other words, complete sample sets) and the GPE  $\mathbf{S}_{xx}$  of PCA calculated from all the samples, as follows:

$$\begin{aligned} \overline{\mathbf{C}}_Q^{(\text{SELF})} &= \beta \mathbf{S}_Q^{(\text{C,lb})} + (1 - \beta) \mathbf{S}_{xx}, \\ \underline{\mathbf{C}}_Q^{(\text{SELF})} &= \beta \mathbf{S}_Q^{(\text{C,lw})} + (1 - \beta) \mathbf{I}_{d_x}, \end{aligned}$$

where  $\beta$  is a hyper parameter satisfying  $0 \leq \beta \leq 1$ . When  $\beta = 1$ , SELF is equivalent to LFDA with only the labeled samples  $(\mathbf{X}^{(\text{C})}, \mathbf{Y}^{(\text{C})})$ . Meanwhile, when  $\beta = 0$ , SELF is equivalent to PCA with all samples in  $\mathbf{X}$ . Generally speaking, SELF inherits the properties of both LFDA and PCA, and their influences can be controlled by the parameter  $\beta$ .

#### 4.10 Semi-supervised CCA

In a similar way to that of SELF, a semi-supervised extension of CCA can be derived, which is called SemiCCA [44].

Assume that there are two samples sets  $\mathbf{X}$  and  $\mathbf{Y}$ , and each includes incomplete sample set  $\mathbf{X}^{(\text{l})}$  and  $\mathbf{Y}^{(\text{l})}$ , namely there are at least one unpaired samples in both  $\mathbf{X}$  and  $\mathbf{Y}$ . SemiCCA integrates the GPE of CCA calculated only from the complete sample sets and the GPE of PCA calculated from the complete and incomplete sample sets, as follows:

$$\begin{aligned} \overline{\mathbf{C}}^{(\text{SemiCCA})} &= \beta \begin{pmatrix} \mathbf{0} & \mathbf{S}_{xy}^{(\text{C})} \\ \mathbf{S}_{yx}^{(\text{C})} & \mathbf{0} \end{pmatrix} + (1 - \beta) \begin{pmatrix} \mathbf{S}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{yy} \end{pmatrix}, \\ \underline{\mathbf{C}}^{(\text{SemiCCA})} &= \beta \begin{pmatrix} \mathbf{S}_{xx}^{(\text{l})} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{yy}^{(\text{l})} \end{pmatrix} + (1 - \beta) \begin{pmatrix} \mathbf{I}_{d_x} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_y} \end{pmatrix}, \end{aligned}$$

When  $\beta = 1$ , SemiCCA is equivalent to CCA with only the complete samples  $(\mathbf{X}^{(\text{C})}, \mathbf{Y}^{(\text{C})})$ . Meanwhile, when  $\beta = 0$ , SemiCCA is equivalent to PCA with all samples in  $\mathbf{X}$  and  $\mathbf{Y}$  under the assumption that  $\mathbf{X}$  and  $\mathbf{Y}$  are uncorrelated with each other.

Another type of semi-supervised extension of CCA has been developed by Blaschko et al. [39]. Please see the detail in Section 6.

### 5. How to Design New Methods

To summarize the discussions so far, we describe (1) GPEs of major existing methods, (2) the way for integrating several GPEs and (3) some semi-supervised extensions by changing the sample sets for calculating GPEs. This section shows that we can easily design new multivariate analysis methods at will by replicating those steps. Note that another way to generate new methods would be possible, and the following one is only one example.

One of the simple extensions is to integrate FDA as supervised dimensionality reduction and CCA as unsupervised dimensionality reduction with a latent model. Consider a problem of video categorization, where its training data includes image features  $\mathbf{X}$ , audio features  $\mathbf{Y}$  and class indexes. Finding appropriate correlations of such three different modals would be still challenging. Several approaches might be possible: (1) FDA for concatenated features  $(\mathbf{X}^\top, \mathbf{Y}^\top)^\top$ , which cannot obtain appropriate correlations between two different types of feature vectors, (2) CCA for two features  $(\mathbf{X}, \mathbf{Y})$  followed by FDA on the compressed domain, which cannot find class-wise differences of correlations.

Here, we newly introduce an integration of CCA and FDA, which enables us to extract class-wise differences of feature cor-

relations as well as to achieve discriminative embedding simultaneously. In the following, we call this method CFDA for the simplicity. CFDA can be formulated by the following equations:

$$\overline{\mathbf{C}}_Q^{(\text{CFDA})} = \beta \begin{pmatrix} \mathbf{0} & \mathbf{S}_{xy} \\ \mathbf{S}_{yx} & \mathbf{0} \end{pmatrix} + (1 - \beta) \mathbf{S}_Q^{(\text{lb})}, \quad (10)$$

$$\underline{\mathbf{C}}_Q^{(\text{CFDA})} = \beta \begin{pmatrix} \mathbf{S}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{yy} \end{pmatrix} + (1 - \beta) \mathbf{S}_Q^{(\text{lw})}. \quad (11)$$

When  $\beta = 1$  CFDA is equivalent to CCA, while when  $\beta = 0$  CFDA is equivalent to FDA for concatenated features  $(\mathbf{X}^\top, \mathbf{Y}^\top)^\top$ .

We note that we do not have to explicitly consider GPEs of FDA and CCA when constructing CFDA. All what we need to design a new multivariate analysis method are that existing methods to be combined can be described by GPE and operations for the combination are shown in Section 3.2.

## 6. Kernelized Extensions

### 6.1 Kernelization of Standard Methods

Almost all the methods in the GPE framework can be kernelized in a similar manner to the existing ones. First, we describe kernel CCA [34], [35] and related regularization techniques.

The original CCA can be extended to, e.g., non-vectorial domains by defining kernels over  $\mathbf{x}$  and  $\mathbf{y}$ ,

$$\begin{aligned} k_x(\mathbf{x}_n, \mathbf{x}_m) &= \langle \phi_x(\mathbf{x}_n), \phi_x(\mathbf{x}_m) \rangle, \\ k_y(\mathbf{y}_n, \mathbf{y}_m) &= \langle \phi_y(\mathbf{y}_n), \phi_y(\mathbf{y}_m) \rangle, \end{aligned}$$

and searching for solutions that lie in the span of  $\phi_x(\mathbf{x})$  and  $\phi_y(\mathbf{y})$

$$\mathbf{w}_x = \sum_{n=1}^N \alpha_n \phi_x(\mathbf{x}_n), \quad \mathbf{w}_y = \sum_{n=1}^N \beta_n \phi_y(\mathbf{y}_n).$$

In this setting, we use the following empirical scatter matrix,

$$\hat{\mathbf{S}}_{xy} = \sum_{n=1}^N \phi_x(\mathbf{x}_n) \phi_y(\mathbf{y}_n)^\top.$$

Denoting the Gram matrices defined by the samples as  $\mathbf{K}_x$  and  $\mathbf{K}_y$ , we can obtain the solution from the following optimization problem with respect to coefficient vectors,  $\alpha$  and  $\beta$ ,

$$\rho^{(\text{kCCA})}(\mathbf{w}_x, \mathbf{w}_y | \mathbf{X}, \mathbf{Y}) = \frac{\alpha^\top \mathbf{K}_x \mathbf{K}_y \beta}{\sqrt{\alpha^\top \mathbf{K}_x^2 \alpha \beta^\top \mathbf{K}_y^2 \beta}}.$$

In the same way as CCA, the optimization can be achieved by solving the following generalized eigenvalue problem:

$$\begin{aligned} \overline{\mathbf{C}}^{(\text{kCCA})} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} &= \lambda \underline{\mathbf{C}}^{(\text{kCCA})} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}, \\ \overline{\mathbf{C}}^{(\text{kCCA})} &= \begin{pmatrix} \mathbf{0} & \mathbf{K}_x \mathbf{K}_y \\ \mathbf{K}_y \mathbf{K}_x & \mathbf{0} \end{pmatrix}, \\ \underline{\mathbf{C}}^{(\text{kCCA})} &= \begin{pmatrix} \mathbf{K}_x^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_y^2 \end{pmatrix}. \end{aligned}$$

Although the bases  $(\mathbf{w}_x, \mathbf{w}_y)$  cannot be explicitly obtained, the projection to those bases can be calculated with the help of the kernel trick:

$$\mathbf{w}_x^\top \phi_x(\mathbf{x}) = \sum_{n=1}^N \alpha_n k_x(\mathbf{x}, \mathbf{x}_n),$$

$$\mathbf{w}_y^\top \phi_y(\mathbf{y}) = \sum_{n=1}^N \beta_n k_y(\mathbf{y}, \mathbf{y}_n).$$

As discussed in Ref. [38], this optimization leads to degenerate solutions in the case that either  $\mathbf{K}_x$  or  $\mathbf{K}_y$  is not invertible. Therefore, the following  $\ell_2$ -regularized formulation should be necessary in general:

$$\overline{\mathbf{C}}^{(\text{kCCA}-\ell_2)} = \begin{pmatrix} \mathbf{0} & \mathbf{K}_x \mathbf{K}_y \\ \mathbf{K}_y \mathbf{K}_x & \mathbf{0} \end{pmatrix},$$

$$\underline{\mathbf{C}}^{(\text{kCCA}-\ell_2)} = \begin{pmatrix} \mathbf{K}_x^2 + \delta_x \mathbf{K}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_y^2 + \delta_y \mathbf{K}_y \end{pmatrix}.$$

Another popular regularization technique is the graph Laplacian method [39], [61]. By using Laplacian regularization, we are able to learn directions that tend to lie along the data manifold estimated from a collection of data. Denoting the empirical graph Laplacian  $\hat{\mathbf{L}}_x$  and  $\hat{\mathbf{L}}_y$  obtained from  $\mathbf{K}_x$  and  $\mathbf{K}_y$ , the formulation is replaced by the following equations:

$$\overline{\mathbf{C}}^{(\text{kCCA-Lap})} = \begin{pmatrix} \mathbf{0} & \mathbf{K}_x \mathbf{K}_y \\ \mathbf{K}_y \mathbf{K}_x & \mathbf{0} \end{pmatrix},$$

$$\underline{\mathbf{C}}^{(\text{kCCA-Lap})} = \begin{pmatrix} \mathbf{K}_x^2 + \gamma_x \mathbf{R}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_y^2 + \gamma_y \mathbf{R}_y \end{pmatrix},$$

$$\mathbf{R}_x = \mathbf{K}_x \hat{\mathbf{L}}_x \mathbf{K}_x, \quad \mathbf{R}_y = \mathbf{K}_y \hat{\mathbf{L}}_y \mathbf{K}_y.$$

**6.2 Non-linear Embedding Methods**

With the kernelized extension, non-linear dimensionality reduction such as locally linear embedding [62] and Laplacian eigenmaps [51] are also in the GPE framework.

**6.2.1 Laplacian Eigenmaps**

Laplacian eigenmaps [51] is one of the popular methods for non-linear embedding. The goal of Laplacian eigenmaps is to find an embedding that preserves the local structure of nearby high-dimensional samples. Laplacian eigenmaps exploits graph Laplacian of a neighborhood graph on the samples  $\mathbf{X}$ , where each edge measures the affinity between two samples. Since a set of edge weights can be expressed by a Gram matrix  $\mathbf{K}_x$ , the objective function of Laplacian eigenmaps to be minimized is

$$\rho^{(\text{LE})}(\mathbf{w}_x|\mathbf{X}) = (\alpha^\top \hat{\mathbf{L}}_x \alpha)(\alpha^\top \hat{\mathbf{D}}_x \alpha)^{-1},$$

where  $\hat{\mathbf{D}}_x$  is a diagonal matrix satisfying  $\hat{\mathbf{D}}_x = \mathbf{K}_x + \hat{\mathbf{L}}_x$ . Therefore, the GPE of Laplacian eigenmaps can be obtained as

$$\overline{\mathbf{C}}^{(\text{LE})} = \hat{\mathbf{L}}_x, \quad \underline{\mathbf{C}}^{(\text{LE})} = \hat{\mathbf{D}}_x.$$

**6.2.2 Locally Linear Embedding (LLE)**

Locally linear embedding (LLE) [62] finds an embedding of the samples  $\mathbf{X}$  that preserves the local structure of nearby samples in the high-dimensional space. LLE builds the embedding by preserving the geometry of pairwise relations between samples in the high-dimensional manifold. LLE first computes a Gram matrix  $\mathbf{K}_x$  containing the structural information of the embedding by minimizing the following function:

$$\rho^{(\text{LLE1})}(\mathbf{K}_x|\mathbf{X}) = \|\mathbf{X}(\mathbf{I}_N - \mathbf{K}_x)\|_F^2$$

s.t.  $\mathbf{K}_x \mathbf{1}_N = \mathbf{1}_N$ ,

where each column of the Gram matrix  $\mathbf{K}_x$  has  $k$  non-zero values. This minimization can be solved via a linear system of equations. Once  $\mathbf{K}_x$  is calculated, LLE next finds a base that minimizes

$$\rho^{(\text{LLE2})}(\mathbf{w}_x|\mathbf{X}) = \alpha^\top (\mathbf{I}_N - \mathbf{K}_x)^2 \alpha \quad \text{s.t. } \alpha^\top \alpha = 1.$$

Therefore, the GPE of LLE can be obtained as

$$\overline{\mathbf{C}}^{(\text{LLE})} = (\mathbf{I}_N - \mathbf{K}_x)^2, \quad \underline{\mathbf{C}}^{(\text{LLE})} = \mathbf{I}_N.$$

**6.3 Clustering Methods**

With the benefit of the kernelized extension of the GPE framework, several clustering methods such as spectral clustering (SC) [13], [63], [64] and normalized cuts (NC) [65] can also be included in the class of multivariate analysis we are concerned with.

$$\overline{\mathbf{C}}^{(\text{SC})} = \mathbf{L}_x, \quad \underline{\mathbf{C}}^{(\text{SC})} = \mathbf{D}_x,$$

$$\overline{\mathbf{C}}^{(\text{NC})} = \mathbf{D}_x^{-1/2} \mathbf{L}_x \mathbf{D}_x^{-1/2}, \quad \underline{\mathbf{C}}^{(\text{NC})} = \mathbf{I}_N.$$

Kernel k-means [66] is also known to belong to this family if we admit introducing a certain iterative procedure [67]. The details can be seen in the preceding work by De la Torre [49], [50].

**6.4 How to Design New Kernelized Methods**

Integrating two methods within the kernelized GPE framework is not straightforward, since a simple addition of Gram matrices is not GPE.

For example, let us consider the following generalized eigenproblem characterized by two matrices of GPE-style second-order statistics:

$$\overline{\mathbf{C}} \mathbf{w} = \lambda \underline{\mathbf{C}} \mathbf{w}, \tag{12}$$

$$\overline{\mathbf{C}} = \mathbf{X} \overline{\mathbf{L}}_1 \mathbf{X}^\top + \overline{\mathbf{L}}_2, \quad \underline{\mathbf{C}} = \mathbf{X} \underline{\mathbf{L}}_1 \mathbf{X}^\top + \underline{\mathbf{L}}_2,$$

where  $\overline{\mathbf{L}}_1$  and  $\underline{\mathbf{L}}_1$  are both positive semi-definite matrices and  $\overline{\mathbf{L}}_2$  and  $\underline{\mathbf{L}}_2$  are both non-negative symmetric matrices, from the definition of GPE. The above generalized eigenproblem can be replaced with the following one without essentially changing the solution [43]

$$\overline{\mathbf{C}}_2 \mathbf{w} = \lambda \underline{\mathbf{C}}_2 \mathbf{w},$$

$$\overline{\mathbf{C}}_2 = \mathbf{X} \{ \overline{\mathbf{L}}_1 + (\mathbf{X}^\top \overline{\mathbf{L}}_2 \mathbf{X})^\dagger \} \mathbf{X}^\top = \mathbf{X} \overline{\mathbf{L}}_3 \mathbf{X}^\top,$$

$$\underline{\mathbf{C}}_2 = \mathbf{X} \{ \underline{\mathbf{L}}_1 + (\mathbf{X}^\top \underline{\mathbf{L}}_2 \mathbf{X})^\dagger \} \mathbf{X}^\top = \mathbf{X} \underline{\mathbf{L}}_3 \mathbf{X}^\top,$$

where  $\dagger$  denotes the Moore-Penrose pseudo-inverse [68]. In this derivation, we used the fact that the matrices  $\overline{\mathbf{C}}$  and  $\overline{\mathbf{C}}_2$  (resp.  $\underline{\mathbf{C}}$  and  $\underline{\mathbf{C}}_2$ ) that characterize the generalized eigenproblems share the same range. Following the procedure shown in Sugiyama et al. [43], we can obtain a kernelized variant of the multivariate analysis method formulated by Eq. (12).

From the above discussion, we can see that when dealing with kernelized multivariate analysis, we have to explicitly derive GPEs of existing methods, and replace the data matrix with its Gram matrix.

**7. Concluding Remarks**

This paper provided a new expression of covariance matrices

and Gram matrices, which we call generalized pairwise expression (GPE). This provided a unified insight into various multivariate analysis methods and their extensions. GPE made it easy to design desired multivariate analysis methods by simple combinations of GPEs of existing methods as templates. According to this methodology, we designed several new multivariate analysis methods.

The GPE framework covers a wide variety of multivariate analysis methods, and thus the way we have presented in this paper for designing new methods is still one of the examples. Developing more general guidelines would be promising future work.

## References

- [1] Russell, B.C., Torralba, A., Murphy, K.P. and Freeman, W.T.: LabelMe: A database and web-based tool for image annotation, *International Journal on Computer Vision*, Vol.77, No.5, pp.157–173 (2008).
- [2] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K. and Fei-Fei, L.: ImageNet: A large-scale hierarchical image database, *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.248–255 (2009).
- [3] Torralba, A., Fergus, R. and Freeman, W.T.: 80 million tiny images: A large data set for nonparametric object and scene recognition, *IEEE Trans. Pattern Analysis and Machine Intelligence*, Vol.30, No.11, pp.1958–1970 (2008).
- [4] Wang, X.J., Zhang, L., Liu, M., Li, Y. and Ma, W.Y.: ARISTA - Image search to annotation on billions of web photos, *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.2987–2994 (2010).
- [5] Smeaton, A.F., Over, P. and Kraaij, W.: Evaluation campaigns and trecvid, *Proc. ACM International Workshop on Multimedia Information Retrieval (MIR)*, pp.321–330 (2006).
- [6] Yuen, J., Russell, B.C., Liu, C. and Torralba, A.: LabelMe video: Building a video database with human annotations, *Proc. International Conference on Computer Vision (ICCV)*, pp.1451–1458 (2009).
- [7] Bertin-Mahieux, T., Ellis, D.P., Whitman, B. and Lamere, P.: The million song dataset, *Proc. International Conference on Music Information Retrieval (ISMIR)*, pp.591–596 (2011).
- [8] Anderson, T.W.: *An Introduction to Multivariate Statistical Analysis (Wiley Series in Probability and Statistics)*, 3rd ed., Wiley-Interscience (2003).
- [9] Bishop, C.M.: *Pattern Recognition and Machine Learning*, Springer (2006).
- [10] Hastie, T., Tibshirani, R. and Friedman, J.: *The Elements of Statistical Learning: Data Mining, Inference, and Prediction (Springer Series in Statistics)*, 2nd ed., Springer-Verlag (2009).
- [11] Schwartz, W.R., Kembhavi, A., Harwood, D. and Davis, L.: Human detection using partial least squares analysis, *Proc. International Conference on Computer Vision (ICCV)*, pp.24–31 (2009).
- [12] Nakayama, H., Harada, T. and Kuniyoshi, Y.: Evaluation of dimensionality reduction methods for image auto-annotation, *Proc. British Machine Vision Conference (BMVC)*, pp.94.1–94.12 (2010).
- [13] Blaschko, M.B. and Lampert, C.H.: Correlational spectral clustering, *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.1–8 (2008).
- [14] Wang, M., Perera, A. and Gutierrez-Osuna, R.: Principal discriminants analysis for small-sample-size problems: Application to chemical sensing, *Proc. IEEE Sensors*, pp.591–594 (2004).
- [15] Pezeshki, A., Azimi-Sadjadi, M.R. and Scharf, L.L.: Undersea target classification using canonical correlation analysis, *IEEE Journal of Oceanic Engineering*, Vol.32, No.4, pp.948–955 (2007).
- [16] Nielsen, A.A.: Multiset canonical correlations analysis and multispectral, truly multi-temporal remote sensing data, *IEEE Trans. Image Processing*, Vol.11, No.3, pp.293–305 (2002).
- [17] Schizas, I., Giannakis, G. and Luo, Z.Q.: Distributed estimation using reduced-dimensionality sensor observations, *IEEE Trans. Signal Processing*, Vol.55, No.8, pp.4284–4299 (2007).
- [18] Hotelling, H.: Analysis of a complex of statistical variables into principal components, *Journal of Educational Psychology*, Vol.24 (1933).
- [19] Fisher, R.A.: The use of multiple measurements in taxonomic problems, *Annals Eugen.*, Vol.7, pp.179–188 (1936).
- [20] Wold, H.: Estimation of principal components and related models by iterative least squares, *Multivariate Analysis*, pp.391–420, Academic Press (1966).
- [21] Kanal, L. and Chandrasekaran, B.: On dimensionality and sample size in statistical pattern classification, *Pattern Recogn.*, Vol.3, No.3, pp.225–234 (1971).
- [22] Campbell, N.A.: Robust procedures in multivariate analysis I: Robust covariance estimation, *Applied Statistics*, Vol.29, No.3, pp.231–237 (1980).
- [23] De la Torre, F. and Black, M.: Robust principal component analysis for computer vision, *Proc. IEEE International Conference on Computer Vision (ICCV)*, pp.362–369 (2001).
- [24] De La Torre, F. and Black, M.J.: A framework for robust subspace learning, *International Journal of Computer Vision*, Vol.54, p.2003 (2003).
- [25] Inoue, K., Hara, K. and Urahama, K.: Robust multilinear principal component analysis, *Proc. IEEE International Conference on Computer Vision (ICCV)*, pp.591–597 (2009).
- [26] Lu, J., Plataniotis, K. and Venetianopoulos, A.: Regularization studies of linear discriminant analysis in small sample size scenarios with application to face recognition, *Pattern Recogn. Lett.*, Vol.26, No.2, pp.181–191 (2005).
- [27] Zhu, M. and Martinez, A.: Subclass discriminant analysis, *IEEE Trans. Pattern Analysis and Machine Intelligence*, Vol.28, No.8, pp.1274–1286 (2006).
- [28] Gkalelis, N., Mezaris, V. and Kompatsiaris, I.: Mixture subclass discriminant analysis, *IEEE Signal Processing Letters*, Vol.18, No.5, pp.319–322 (2011).
- [29] Schölkopf, B., Smola, A. and Müller, K.R.: Nonlinear component analysis as a kernel eigenvalue problem, *Neural Computation*, Vol.10, No.5, pp.1299–1319 (1998).
- [30] Mika, S., Ratsch, G., Weston, J., Schölkopf, B. and Mullers, K.R.: Fisher discriminant analysis with kernels, *Proc. IEEE Workshop on Neural Networks for Signal Processing (NNSP)*, No.8, pp.41–48 (1999).
- [31] Baudat, G. and Anouar, F.: Generalized discriminant analysis using a kernel approach, *Neural Computation*, Vol.12, No.10, pp.2385–2404 (2000).
- [32] Dai, G., Yan Yeung, D. and Chang, H.: Extending kernel Fisher discriminant analysis with the weighted pairwise Chernoff criterion, *Proc. European Conference on Computer Vision (ECCV)*, pp.308–320 (2006).
- [33] Bishop, C.: Linear models for regression, *Pattern Recognition and Machine Learning*, ch. 3, Springer (2006).
- [34] Akaho, S.: A kernel method for canonical correlation analysis, *Proc. International Meeting of the Psychometric Society (IMPS)*, pp.1–5, Springer-Verlag (2001).
- [35] Lai, P.L. and Fyfe, C.: Kernel and nonlinear canonical correlation analysis, *Proc. International Joint Conference on Neural Networks (IJCNN)*, p.4614 (2000).
- [36] Hoerl, A.E.: Application of ridge analysis to regression problem, *Chemical Engineering Progress*, Vol.58, pp.54–59 (1962).
- [37] Tikhonov, A.: On the stability of inverse problems, *Dokl. Akad. Nauk SSSR*, Vol.39, No.5, pp.195–198 (1943).
- [38] Haroon, D.R., Szedmak, S. and Shawe-Taylor, J.: Canonical correlation analysis: An overview with application to learning methods, *Neural Computation*, Vol.16, No.12, pp.2639–2664 (2004).
- [39] Blaschko, M., Lampert, C. and Gretton, A.: Semi-supervised Laplacian regularization of kernel canonical correlation analysis, *Proc. European Conference on Machine Learning and Knowledge Discovery in Databases (ECML-PKDD)*, pp.133–145 (2008).
- [40] Cox, T. and Cox, M.: *Multidimensional Scaling, Monographs on Statistics and Applied Probability*, No.1, Chapman & Hall (1994).
- [41] He, X. and Niyogi, P.: Locality preserving projections, *Advances in Neural Information Processing Systems (NIPS)*, pp.1–8 (2003).
- [42] Sugiyama, M.: Dimensionality reduction of multimodal labeled data by local Fisher discriminant analysis, *Journal of Machine Learning Research*, Vol.8, No.5, pp.1027–1061 (2007).
- [43] Sugiyama, M., Idé, T., Nakajima, S. and Sese, J.: Semi-supervised local Fisher discriminant analysis for dimensionality reduction, *Machine Learning*, Vol.78, No.1–2, pp.35–61 (2010).
- [44] Kimura, A., Kameoka, H., Sugiyama, M., Nakano, T., Maeda, E., Sakano, H. and Ishiguro, K.: SemiCCA: Efficient semi-supervised learning of canonical correlations, *Proc. IAPR International Conference on Pattern Recognition (ICPR)*, pp.2933–2936 (2010).
- [45] Borga, M., Landelius, T. and Knutsson, H.: A unified approach to PCA, PLS, MLR and CCA, Report LiTH-ISY-R-1992, ISY, SE-581 83 Linköping, Sweden (Nov. 1997).
- [46] De Bie, T., Cristianini, N. and Rosipal, R.: Eigenproblems in pattern recognition, *Handbook of Geometric Computing: Applications in Pattern Recognition, Computer Vision, Neural Computing, and Robotics*, pp.129–170, Springer (2005).
- [47] Sun, L., Ji, S. and Ye, J.: A least squares formulation for a class of generalized eigenvalue problems in machine learning, *Proc. International Conference on Machine Learning (ICML)*, pp.977–984 (2009).
- [48] Sun, L., Ji, S. and Ye, J.: Canonical correlation analysis for multil-

abel classification: A least-squares formulation, extensions, and analysis, *IEEE Trans. Pattern Analysis and Machine Intelligence*, Vol.33, pp.194–200 (2011).

[49] De la Torre, F.: A unification of component analysis methods, *Handbook of Pattern Recognition and Computer Vision*, 4th ed., Chen, C. (Ed.), ch. 1, pp.3–22, World Scientific Pub Co Inc (2010).

[50] De la Torre, F.: A least-squares framework for component analysis, *IEEE Trans. Pattern Analysis and Machine Intelligence*, Vol.34, No.6, pp.1041–1055 (2012).

[51] Belkin, M. and Niyogi, P.: Laplacian eigenmaps for dimensionality reduction and data representation, *Neural Computation*, Vol.15, pp.1373–1396 (2002).

[52] Yanai, H. and Puntanen, S.: Partial canonical correlation associated with the inverse and some generalized inverse of a partitioned dispersion matrix, *Proc. Pacific Area Statistical Conference on Statistical Sciences and Data Analysis*, pp.253–264 (1993).

[53] Bach, F.R. and Jordan, M.I.: A probabilistic interpretation of canonical correlation analysis, Tech. Rep. 688, Department of Statistics, University of California, Berkeley (2005).

[54] Jolliffe, I.T.: A note on the use of principal components in regression, *Journal of the Royal Statistical Society*, Vol.31, No.3 (1982).

[55] Burnham, A.J., MacGregor, J.F. and Viveros, R.: Latent variable multivariate regression modeling, *Chemometrics and Intelligent Laboratory Systems*, pp.167–180 (Aug. 1999).

[56] Worsley, K.J., Poline, J.b., Friston, K.J. and Evans, A.C.: Characterizing the response of PET and fMRI data using multivariate linear models, *NeuroImage*, Vol.6, pp.305–319 (1997).

[57] Trygg, J. and Wold, S.: Orthogonal projections to latent structures (o-pls), *Journal of Chemometrics*, Vol.16, No.3, pp.119–128 (2002).

[58] Wang, H., Gottfries, J., Barrenäs, F. and Benson, M.: Identification of novel biomarkers in seasonal allergic rhinitis by combining proteomic, multivariate and pathway analysis, *PLoS ONE*, Vol.6, No.8, p.e23563 (2011).

[59] Bach, F.: Kernel independent component analysis, *Journal of Machine Learning Research*, Vol.3, pp.1–48 (2002).

[60] Zelnik-manor L. and Perona, P.: Self-tuning spectral clustering, *Advances in Neural Information Processing Systems (NIPS)*, pp.1601–1608 (2004).

[61] Belkin, M. Niyogi, P. and Sindhvani, V.: Manifold regularization: A geometric framework for learning from labeled and unlabeled examples, *Journal of Machine Learning Research*, Vol.7, pp.2399–2434 (Dec. 2006).

[62] Roweis, S.T. and Saul, L.K.: Nonlinear dimensionality reduction by locally linear embedding, *Science*, Vol.290, pp.2323–2326 (Dec. 2000).

[63] Weiss, Y.: Segmentation using eigenvectors: A unifying view, *Proc. IEEE International Conference on Computer Vision (ICCV)*, pp.975–982 (1999).

[64] Yu, S. and Shi, J.: Multiclass spectral clustering, *Proc. IEEE International Conference on Computer Vision (ICCV)*, No.10, pp.313–319 (2003).

[65] Shi, J. and Malik, J.: Normalized cuts and image segmentation, *IEEE Trans. Pattern Analysis and Machine Intelligence*, Vol.22, No.8, pp.888–905 (2000).

[66] Dhillon, I.S., Guan, Y. and Kulis, B.: Kernel k-means: spectral clustering and normalized cuts, *Proc. ACM International Conference on Knowledge Discovery and Data Mining (KDD)*, pp.551–556, ACM (2004).

[67] Zass, R. and Shashua, A.: A unifying approach to hard and probabilistic clustering, *Proc. IEEE International Conference on Computer Vision (ICCV)*, pp.294–301 (2005).

[68] Albert, A.: *Regression and the Moore-Penrose pseudoinverse*, *Mathematics in Science and Engineering*, Elsevier Science (1972).



**Akisato Kimura** received his B.E., M.E. and D.E. degrees in communications and integrated systems from Tokyo Institute of Technology, Japan in 1998, 2000 and 2007, respectively. Since 2000, he has been with NTT Communication Science Laboratories, NTT Corporation, where he is currently a senior research scientist in

Innovative Communication Laboratory. He has been engaged in work on multimedia content identification, automatic multimedia annotation, human visual attention modeling and social media mining. His research interests include pattern recognition, computer vision, image processing, human visual perception, statistical signal processing, data mining and social media. He is a member of IEICE, IEEE and ACM SIGMM/SIGKDD.



**Masashi Sugiyama** received his B.E., M.E., and Ph.D. degrees from Department of Computer Science, Tokyo Institute of Technology, Tokyo, Japan, in 1997, 1999, and 2001, respectively. In 2001, he was appointed as a Research Associate in the same institute, and from 2003, he is an Associate Professor. His research interests

include theory and application of machine learning.



**Hitoshi Sakano** received his B.S. degree in physics from Chuo University Tokyo and the M.S. degrees in physics from Saitama University, and the Ph.D. in applied physics from Waseda University, Tokyo in 1988, 1990 and 2008, respectively. He joined NTT Communication Science Laboratories in 2008 and studied pattern recognition technology. He is a member of IEEE, IEICE and Physical Society of Japan (PSJ).



**Hirokazu Kameoka** received his B.E., M.E. and Ph.D. degrees all from the University of Tokyo, Japan, in 2002, 2004 and 2007, respectively. He is currently a research scientist at NTT Communication Science Laboratories and an Adjunct Associate Professor at the University of Tokyo. His research interests include computational auditory scene analysis, statistical signal processing, speech and music processing, and machine learning.

He is a member of IEEE, IEICE, IPSJ and ASJ. He received 13 awards over the past 9 years, including the IEEE Signal Processing Society 2008 SPS Young Author Best Paper Award.