

# 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>

**Abstract:** 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). 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 images [1], [2], [3], [4] and videos [5], [6] nowadays. However, we are now facing a difficulty in finding an intrinsic trend and nature of such a massive collection of data. Multivariate analysis is traditional, quite simple but might be one of the powerful tools to obtain a hidden structure embedded in the data. Actually, multivariate analysis has been still an important tool, and recent reports showed its effectiveness for several tasks, e.g. human detection [7], image annotation [8], [9], sensor data mining [10], [11], [12].

Principal component analysis (PCA) [13], Fisher discriminant analysis (FDA) [14], multivariate linear regression (MLR), canonical correlation analysis (CCA) [13], and partial least squares (PLS) [15] 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 so-called small sample size problem, i.e., the situation where the number of training samples is small compared with their dimensionality [16], [17], [18], [19], [20], [21]).

Kernel multivariate analysis methods as kernelized extensions of those standard methods have been also developed to deal with non-vector samples and non-linear analysis [22], [23], [24], [25], [26], [27]). 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 [28], [29], [30] and graph Laplacian method [31]. In addition, improvements of robustness against outliers and non-Gaussianity (i.e. multi-dimensional scaling (MDS) [32], locality preserving projection (LPP) [33] and local Fisher discriminant analysis (LFDA) [34]) and their extensions to semi-supervised dimensionality reduction [31], [35], [36] have been considered.

A lot of multivariate analysis methods and several trials to unify these methods have been presented so far. Borge et al [37] and De Bie et al [38] presented that several major linear multivariate analysis method can be formulated by a unified form of generalized eigenvalue problems. Sun et al [39], [40] showed the equivalence between a certain class of generalized eigenvalue problems and least squares ones under a mild assumption. De la Torre [41], [42] extended the work by Sun et al 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 multivariate analysis for a specific purpose or domain still remains an open problem. Until now, researchers have had to choose one of the existing methods that seems best to address the problem of interest, or

<sup>1</sup> NTT Communication Science Laboratories, NTT Corporation, 2-4 Hikaridai, Seika, Soraku, Kyoto, 619-0237 Japan.

<sup>2</sup> Graduate School of Information Science and Engineering, Tokyo Institute of Technology, 2-12-1 Oookayama, Meguro, Tokyo, 152-8552 Japan.

<sup>3</sup> Graduate School of Information Science and Technologies, the University of Tokyo, 7-3-1 Hongo, Bunkyo, Tokyo, 113-8656 Japan.

a) akisato@ieee.org

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 covariance matrices and Gram matrices, which we call *Generalized pairwise expression (GPE)* to make it easy to design a new multivariate analysis method with desired property. The methodology is quite simple: adopting the above mentioned special cases as templates, and generating a new method by combining these templates appropriately. This characteristics has not been discussed yet in any previous researches to our best knowledge. It is also possible to individually select and arrange samples for calculating the scatter matrices of the methods to be combined, which enables us to extend CA methods to semi-supervised ones and multi-modal ones.

Our contributions can be summarized as follows:

- (1) Providing a unified formulation of various multivariate analysis methods via GPE.
- (2) Making it easy to implement a multivariate analysis method with desired property by simply combining the GPEs of existing methods.
- (3) Designing new multivariate analysis methods based on the methodology derived from the GPE.

## 2. Multivariate analysis for vector data

Consider two sets  $\mathbf{X}$  and  $\mathbf{Y}$  of samples<sup>\*1</sup>, where each set contains  $N_x$  and  $N_y$  samples, and each sample can be expressed as a vector with  $d_x$  and  $d_y$  dimensions, respectively, as follows:

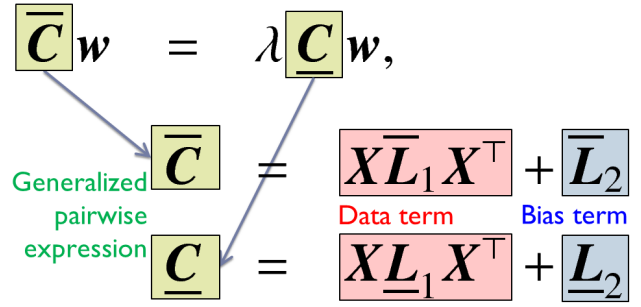
$$\begin{aligned} \mathbf{X} &= \{\mathbf{x}_1, \dots, \mathbf{x}_{N_x}\}, \\ \mathbf{Y} &= \{\mathbf{y}_1, \dots, \mathbf{y}_N, \mathbf{y}_{N_x+1}, \dots, \mathbf{y}_{N_x+N_y-N}\}. \end{aligned}$$

For brevity, both of the sample sets  $\mathbf{X}$  and  $\mathbf{Y}$  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  $\mathbf{X}$  and  $\mathbf{Y}$  of samples is separated into the following two types: *Complete sample sets*  $\mathbf{X}^{(C)}$  and  $\mathbf{Y}^{(C)}$  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*  $\mathbf{X}^{(I)}$  and  $\mathbf{Y}^{(I)}$  so that every sample  $\mathbf{x}_n$  (resp.  $\mathbf{y}_n$ ) cannot find the co-occurring sample.

$$\begin{aligned} \mathbf{X}^{(C)} &= \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}, \\ &= \{\mathbf{x}_1^{(C)}, \mathbf{x}_2^{(C)}, \dots, \mathbf{x}_N^{(C)}\}, \\ \mathbf{Y}^{(C)} &= \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}, \\ &= \{\mathbf{y}_1^{(C)}, \mathbf{y}_2^{(C)}, \dots, \mathbf{y}_N^{(C)}\}, \\ \mathbf{X}^{(I)} &= \{\mathbf{x}_{N+1}, \mathbf{x}_{N+2}, \dots, \mathbf{x}_{N_x}\}, \\ &= \{\mathbf{x}_1^{(I)}, \mathbf{x}_2^{(I)}, \dots, \mathbf{x}_{N_x-N}^{(I)}\}, \\ \mathbf{Y}^{(I)} &= \{\mathbf{y}_{N_x+1}, \mathbf{y}_{N_x+2}, \dots, \mathbf{y}_{N_x+N_y-N}\}, \\ &= \{\mathbf{y}_1^{(I)}, \mathbf{y}_2^{(I)}, \dots, \mathbf{y}_{N_y-N}^{(I)}\} \end{aligned}$$

First, we concentrate on the case that  $N_x = N_y = N$ ,

<sup>\*1</sup> The following discussion can be easily extended to more than 2 sets of samples sets [43].



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

namely all the samples are paired, unless otherwise stated.

Many linear multivariate analysis methods developed so far involve an optimization problem of the following form:

$$\begin{aligned} \mathbf{w}^{(\text{opt})} &= \arg \max_{\mathbf{w} \in \mathcal{R}^d} R(\mathbf{w}), \\ R(\mathbf{w}) &= \mathbf{w}^\top \bar{\mathbf{C}}\mathbf{w} (\mathbf{w}^\top \mathbf{C}\mathbf{w})^{-1}, \end{aligned} \quad (1)$$

where  $\bar{\mathbf{C}}$  and  $\mathbf{C}$  are square matrices with certain statistical nature. For example,  $\bar{\mathbf{C}}$  is a scatter matrix of  $\mathbf{X}$  and  $\mathbf{C}$  is an identity matrix in PCA, and  $\bar{\mathbf{C}}$  is a between-class scatter matrix and  $\mathbf{C}$  is a within-class scatter matrix in FDA. Roughly speaking,  $\bar{\mathbf{C}}$  encodes the quantity that we want to increase, and  $\mathbf{C}$  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:

$$\begin{aligned} \mathbf{w}^{(\text{opt})} &= \arg \max_{\mathbf{w} \in \mathcal{R}^d} R_1(\mathbf{w}) \text{ s.t. } R_2(\mathbf{w}) = 1, \\ R_1(\mathbf{w}) &= \mathbf{w}^\top \bar{\mathbf{C}}\mathbf{w}, \quad R_2(\mathbf{w}) = \mathbf{w}^\top \mathbf{C}\mathbf{w}. \end{aligned} \quad (2)$$

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

$$\bar{\mathbf{C}}\mathbf{w} = \lambda \mathbf{C}\mathbf{w}. \quad (3)$$

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

## 3. Generalized pairwise expression

When addressing linear multivariate analysis methods, we often deal with the following type of second-order statistics 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

$$\mathbf{S}_{Q,xy} = \sum_{n=1}^N \sum_{m=1}^N Q_{n,m} (\mathbf{x}_n - \mathbf{x}_m)(\mathbf{y}_n - \mathbf{y}_m)^\top, \quad (4)$$

where  $\mathbf{Q}$  is an  $N \times N$  non-negative, semi-definite and symmetric matrix. A typical example is the scatter matrix<sup>\*2</sup>:

<sup>\*2</sup> Due to the limited space, we describe only the scatter matrix  $\mathbf{S}_{xy}$  and its extensions with the pairwise form. The scatter matrices  $\mathbf{S}_{xx}$  and  $\mathbf{S}_{yy}$ , and their extensions can be easily derived in the same way.

$$\mathbf{S}_{xy} = N^{-1} \sum_{n=1}^N \mathbf{x}_n \mathbf{y}_n^\top.$$

Let  $\mathbf{D}_Q$  be the  $N \times N$  diagonal matrix with

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

and let  $\mathbf{L}_Q$  be  $\mathbf{L}_Q = \mathbf{D}_Q - \mathbf{Q}$ . Then, the matrix  $\mathbf{S}_{Q,xy}$  can be expressed in terms of  $\mathbf{L}_Q$  as follows:

$$\mathbf{S}_{Q,xy} = \mathbf{X} \mathbf{L}_Q \mathbf{Y}^\top.$$

The above expression is called the *pairwise expression (PE)* of the second-order statistics  $\mathbf{S}_{Q,xx}$ [35]. If  $\mathbf{Q}$  is a weight matrix for a graph with  $n$  nodes,  $\mathbf{L}_Q$  can be regarded as a graph Laplacian matrix in the spectral graph theory. If  $\mathbf{Q}$  is symmetric and its elements are all non-negative,  $\mathbf{L}_Q$  is known to be positive semi-definite.

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

$$\hat{\mathbf{S}}_{Q,xy} = \mathbf{X} \mathbf{L}_{Q,1} \mathbf{Y}^\top + \mathbf{L}_2,$$

where  $\mathbf{L}_{Q,1}$  is a  $N \times N$  positive semi-definite matrix, and  $\mathbf{L}_2$  is a  $d_x \times d_y$  non-negative semi-definite matrix. We do not have to explicitly consider the matrix  $\mathbf{Q}$  for the following discussions:

$$\hat{\mathbf{S}}_{xy} = \mathbf{X} \mathbf{L}_1 \mathbf{Y}^\top + \mathbf{L}_2. \quad (5)$$

After all, we call this expression as the *generalized pairwise expression (GPE)*. The first term of Equation (5) is called the *data term* since it depends on the sample data, and the second term is called the *bias term*.

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

- (1) If  $\mathbf{A}$  is GPE and  $\beta > 0$  is a constant, then  $\beta \mathbf{A}$  is also GPE.
- (2) If both  $\mathbf{A}$  and  $\mathbf{B}$  are GPE with  $d_x$  rows and  $d_y$  columns, then  $\mathbf{A} + \mathbf{B}$  is also GPE with  $d_x$  rows and  $d_y$  columns.
- (3) If  $\mathbf{A}$  is GPE with  $d_x$  rows and  $d_y$  columns, and  $\mathbf{B}$  is GPE with  $d_y$  rows and  $d_z$  columns, then  $\mathbf{A}\mathbf{B}$  is also GPE with  $d_x$  rows and  $d_z$  columns.

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

First, let us denote  $\mathbf{A}$  and  $\mathbf{B}$  as follows:

$$\begin{aligned} \mathbf{A} &= \mathbf{X} \mathbf{L}_{A1} \mathbf{Y}^\top + \mathbf{L}_{A2}, \\ \mathbf{B} &= \mathbf{Y} \mathbf{L}_{B1} \mathbf{Z}^\top + \mathbf{L}_{B2}, \end{aligned}$$

where  $\mathbf{L}_{A1}$  (resp.  $\mathbf{L}_{B1}$ ) is a positive semi-definite matrix with  $d_x$  (resp.  $d_y$ ) rows and  $d_y$  (resp.  $d_z$ ) columns, and  $\mathbf{L}_{A2}$  (resp.  $\mathbf{L}_{B2}$ ) is a  $d_x \times d_y$  (resp.  $d_y \times d_z$ ) non-negative matrix. Then, we obtain

$$\begin{aligned} \mathbf{A}\mathbf{B} &= (\mathbf{X} \mathbf{L}_{A1} \mathbf{Y}^\top + \mathbf{L}_{A2})(\mathbf{Y} \mathbf{L}_{B1} \mathbf{Z}^\top + \mathbf{L}_{B2}), \\ &= \mathbf{X}(\mathbf{L}_{A1} \mathbf{Y}^\top \mathbf{Y} \mathbf{L}_{B1}) \mathbf{Z}^\top + (\mathbf{L}_{A2} \mathbf{Y}) \mathbf{L}_{B1} \mathbf{Z}^\top \\ &\quad + \mathbf{X} \mathbf{L}_{A1} (\mathbf{Y}^\top \mathbf{L}_{B2}) + \mathbf{L}_{A2} \mathbf{L}_{B2}. \end{aligned}$$

Table 1 GPEs of standard methods

Method	$\bar{\mathbf{C}}$	$\underline{\mathbf{C}}$
PCA	$\mathbf{S}_{xx}$	$\mathbf{I}_{d_x}$
FDA	$\mathbf{S}_{xx}^{(b)}$	$\mathbf{S}_{xx}^{(w)}$
CCA	$\begin{bmatrix} \mathbf{0} & \mathbf{S}_{xy} \\ \mathbf{S}_{yx} & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{S}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{yy} \end{bmatrix}$
MLR	$\begin{bmatrix} \mathbf{0} & \mathbf{S}_{xy} \\ \mathbf{S}_{yx} & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{S}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_y} \end{bmatrix}$
PCR[44]	$\begin{bmatrix} \mathbf{0} & \mathbf{S}_{\hat{x}y} \\ \mathbf{S}_{y\hat{x}} & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{S}_{\hat{x}\hat{x}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_y} \end{bmatrix}$
OPLS	$\mathbf{S}_{xy} \mathbf{S}_{xy}^\top$	$\mathbf{S}_{xx}$
Ridge regression	$\begin{bmatrix} \mathbf{0} & \mathbf{S}_{xy} \\ \mathbf{S}_{yx} & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{S}_{xx} + \delta \mathbf{I}_{d_x} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d_y} \end{bmatrix}$
LPP[33]	$\mathbf{X} \mathbf{L} \mathbf{X}^\top$	$\mathbf{X} \mathbf{D} \mathbf{X}^\top$
LFDA[45]	$\mathbf{S}_{Q,xx}^{(lb)}$	$\mathbf{S}_{Q,xx}^{(lw)}$

PCR: Principal component regression, OPLS: Orthogonal partial least-squares.

$\mathbf{S}_{xx}^{(b)}$  and  $\mathbf{S}_{xx}^{(w)}$ : Between-class and within-class scatter matrices of  $\mathbf{X}$ ,  $\hat{\mathbf{X}} = \mathbf{U}_K \Sigma_K \mathbf{V}_K^\top$ :  $K$ -rank approximation of  $\mathbf{X}$  by SVD,  $\mathbf{I}_d$ :  $d \times d$  identity matrix,  $\delta > 0$ : constant,  $\mathbf{S}_{Q,xx}^{(b)}$  and  $\mathbf{S}_{Q,xx}^{(w)}$ : between-class and within-class scatter matrices of  $\mathbf{X}$  weighted by an  $N \times N$  non-negative symmetric matrix.

Here, we can find some matrices  $\mathbf{L}_{C_i}$  ( $i = 1, 2, 3$ ) satisfying the following relationships, if  $N \geq \max(d_x, d_y, d_z)$ :

$$\begin{aligned} \mathbf{L}_{C1} &= \mathbf{L}_{A1} \mathbf{Y}^\top \mathbf{Y} \mathbf{L}_{B1}, \\ \mathbf{X} \mathbf{L}_{C2} &= \mathbf{L}_{A2} \mathbf{Y}, \\ \mathbf{L}_{C3} \mathbf{Z}^\top &= \mathbf{Y}^\top \mathbf{L}_{B2}. \end{aligned}$$

This implies that

$$\begin{aligned} \mathbf{A}\mathbf{B} &= \mathbf{X} \mathbf{L}_{C1} \mathbf{Z}^\top + \mathbf{X} \mathbf{L}_{C2} \mathbf{L}_{B1} \mathbf{Z}^\top \\ &\quad + \mathbf{X} \mathbf{L}_{A1} \mathbf{L}_{C3} \mathbf{Z}^\top + \mathbf{L}_{A2} \mathbf{L}_{B2} \\ &= \mathbf{X}(\mathbf{L}_{C1} + \mathbf{L}_{C2} \mathbf{L}_{B1} + \mathbf{L}_{A1} \mathbf{L}_{C3}) \mathbf{Z}^\top + \mathbf{L}_{A2} \mathbf{L}_{B2} \\ &= \mathbf{X} \mathbf{L}_{D1} \mathbf{Z}^\top + \mathbf{L}_{D2}, \end{aligned}$$

for some matrices  $\mathbf{L}_{D1}$  and  $\mathbf{L}_{D2}$ , which means  $\mathbf{A}\mathbf{B}$  is also GPE.  $\square$

Recall that the class of multivariate analysis we are dealing with can be expressed as  $\bar{\mathbf{C}}\mathbf{w} = \lambda \underline{\mathbf{C}}\mathbf{w}$ , and both  $\bar{\mathbf{C}}$  and  $\underline{\mathbf{C}}$  can be expressed by GPEs or their augmented matrices. The notable point is that various multivariate analysis methods can be easily designed with the help of these GPE properties, namely by combining GPEs of existing methods with desired properties. 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

The GPEs of the standard CA methods are listed in Table 1. Several detailed derivations can be seen in [37], [38]. Instead, this paper provides several significant examples that would be quite an important hint when constructing new CA methods.

#### 4.2 Locality preserving projection (LPP)

Locality preserving projections (LPP) [33] 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 [46], 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 [46]. 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

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

#### 4.3 Local Fisher discriminant analysis (LFDA)

Local Fisher discriminant analysis (LFDA) [34] 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  $Q$  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  $Q^{(\text{lb})}$  and  $Q^{(\text{lw})}$  are the  $N \times N$  matrices with

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

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

where  $N_c$  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:

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

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

Semi-supervised local fisher discriminant analysis, called SELF [35], integrates LFDA as a supervised dimensionality reduction and PCA as a 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,
- (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  $\mathbf{X}^{(I)}$  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 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:

$$\overline{\mathbf{C}}_Q^{(\text{SELF})} = \beta \mathbf{S}_Q^{(\text{C,lb})} + (1 - \beta) \mathbf{S}_{xx},$$

$$\underline{\mathbf{C}}_Q^{(\text{SELF})} = \beta \mathbf{S}_Q^{(\text{C,lw})} + (1 - \beta) \mathbf{I}_{d_x},$$

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  $\beta$ .

#### 4.5 Semi-supervised CCA

In a similar way to that of SELF, a semi-supervised extension of CCA can be derived, which is called SemiCCA [36].

Assume that there are two samples sets  $\mathbf{X}$  and  $\mathbf{Y}$ , and each includes incomplete sample set  $\mathbf{X}^{(I)}$  and  $\mathbf{Y}^{(I)}$  only exists, 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:

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

When  $\beta = 1$ , SemiCCA is equivalent to CCA with only the complete samples  $(\mathbf{X}^{(C)}, \mathbf{Y}^{(C)})$ .

## 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 correlations 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 equation:

$$\overline{\mathbf{C}}_Q^{(\text{CFDA})} = \beta \begin{pmatrix} \mathbf{0} & \mathbf{S}_{xy} \\ \mathbf{S}_{yx} & \mathbf{0} \end{pmatrix} + (1 - \beta) \mathbf{S}_Q^{(\text{lb})}, \quad (6)$$

$$\underline{\mathbf{C}}_Q^{(\text{CFDA})} = \beta \begin{pmatrix} \mathbf{S}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{yy} \end{pmatrix} + (1 - \beta) \mathbf{S}_Q^{(\text{lw})}. \quad (7)$$

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

## 6. Kernelized extensions

### 6.1 Kernelization of standard methods

A lot of methods in the GPE framework can be kernelized in a similar manner to the existing ones. The GPEs of major kernelized CA methods are listed in Table 1. By introducing kernelized expression, several methods for clustering and local embedding can be included in this framework, e.g.

**Table 2** GPEs of kernelized CA methods

Method	$\overline{\mathbf{C}}$	$\underline{\mathbf{C}}$
kPCA	$\mathbf{K}_x$	$\mathbf{I}_N$
kFDA	$\mathbf{K}_x^{(b)}$	$\mathbf{K}_x^{(w)}$
kCCA	$\begin{bmatrix} \mathbf{0} & \mathbf{K}_x \mathbf{K}_y \\ \mathbf{K}_y \mathbf{K}_x & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{K}_x^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_y^2 \end{bmatrix}$
kMLR	$\begin{bmatrix} \mathbf{0} & \mathbf{K}_x \mathbf{K}_y \\ \mathbf{K}_y \mathbf{K}_x & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{K}_x^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_N \end{bmatrix}$
kCCA+ $\ell_2$ [30]	$\begin{bmatrix} \mathbf{0} & \mathbf{K}_x \mathbf{K}_y \\ \mathbf{K}_y \mathbf{K}_x & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{K}_x^{(\ell_2)} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_y^{(\ell_2)} \end{bmatrix}$
L-kCCA[31]	$\begin{bmatrix} \mathbf{0} & \mathbf{K}_x \mathbf{K}_y \\ \mathbf{K}_y \mathbf{K}_x & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \mathbf{K}_x^{(L)} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_x^{(L)} \end{bmatrix}$
LE, SC	$\mathbf{L}_x$	$\mathbf{D}_x$
LLE	$\mathbf{K}_x^{(LL)} \mathbf{K}_x^{(LL)\top}$	$\mathbf{I}_N$
NC, nSC	$\mathbf{D}_x^{-1/2} \mathbf{L}_x \mathbf{D}_x^{-1/2}$	$\mathbf{I}_N$

L-kCCA: Laplacian-regularized kernel CCA,  $\mathbf{K}_x, \mathbf{K}_y$ : Gram matrices,  $\mathbf{K}_x^{(\ell_2)} = \mathbf{K}_x^2 + \delta_x \mathbf{K}_x$ ,  $\mathbf{K}_x^{(L)} = \mathbf{K}_x^2 + \gamma_x \mathbf{R}_x$ ,  $\mathbf{R}_x = \mathbf{K}_x \mathbf{L}_x \mathbf{K}_x$ , LE: Laplacian eigenmap, SC: Spectral clustering, NC: Normalized cuts, nSC: normalized SC,  $\mathbf{K}_x^{(LL)} = \mathbf{I}_N - \mathbf{K}_x$

Laplacian eigenmap (LE), locally linear embedding (LLE), spectral clustering (SC) and normalized cuts (NC).

### 6.2 How to design new kernelized methods

Integrating two methods within the kernelized GPE framework is not obvious, since a simple addition of Gram matrices is not GPE. One example can be seen in a kernelized extension of SELF, called kernel SELF [35]. Remember that the original SELF integrates LFDA with labeled samples and PCA with all the samples (see Section 4.4), and it can be formulated by a localized between-class scatter matrix  $\mathbf{S}_Q^{(\text{C,lb})}$ , localized within-class matrix  $\mathbf{S}_Q^{(\text{C,lw})}$  and the ordinary scatter matrix  $\mathbf{S}_{xx}$ . Kernel SELF can be formulated via their Laplacian matrices  $\mathbf{L}_Q^{(\text{C,lb})}$ ,  $\mathbf{L}_Q^{(\text{C,lw})}$ ,  $\mathbf{L}_{xx}$ , as follows:

$$\overline{\mathbf{C}}^{(\text{kSELF})} = \mathbf{K}_x \{ \beta \mathbf{L}_Q^{(\text{C,lb})} + (1 - \beta) \mathbf{L}_{xx} \} \mathbf{K}_x,$$

$$\underline{\mathbf{C}}^{(\text{kSELF})} = \beta \mathbf{K}_x \mathbf{L}_Q^{(\text{C,lb})} \mathbf{K}_x + (1 - \beta) \mathbf{K}_x.$$

From this formulation, we can see that a weighted sum of GPEs in original multivariate analysis corresponds to a weighted sum of Laplacian matrices in kernelized multivariate analysis. Namely, when dealing with kernelized multivariate analysis, we have to explicitly derive GPEs of existing methods, and replace the data matrix into its Gram matrix.

## 7. Concluding remarks

This paper provided a new theoretical 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 significant future work.

## References

- [1] B.C. Russell, A. Torralba, K.P. Murphy, and W.T. Freeman, "Labelme: A database and web-based tool for image annotation," *IJCV*, vol.77, pp.157–173, 2008.
- [2] J. Deng, W. Dong, R. Socher, L.J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," *Proc. CVPR*, 2009.
- [3] A. Torralba, R. Fergus, and W.T. Freeman, "80 million tiny images: A large data set for nonparametric object and scene recognition," *IEEE Trans. PAMI*, vol.30, no.11, pp.1958–1970, 2008.
- [4] X.J. Wang, L. Zhang, M. Liu, Y. Li, and W.Y. Ma, "ARISTA - image search to annotation on billions of web photos," *Proc. CVPR*, pp.2987–2994, 2010.
- [5] A.F. Smeaton, P. Over, and W. Kraaij, "Evaluation campaigns and trecvid," *Proc. MIR*, pp.321–330, 2006.
- [6] J. Yuen, B.C. Russell, C. Liu, and A. Torralba, "Labelme video: Building a video database with human annotations," *Proc. ICCV*, pp.1451–1458, 2009.
- [7] W.R. Schwartz, A. Kembhavi, D. Harwood, and L.S. Davis, "Human detection using partial least squares analysis," *Proc. ICCV*, 2009.
- [8] H. Nakayama, T. Harada, and Y. Kuniyoshi, "Evaluation of dimensionality reduction methods for image auto-annotation," *Proc. BMVC*, pp.1–12, 2010.
- [9] M.B. Blaschko and C.H. Lampert, "Correlational spectral clustering," *Proc. CVPR*, pp.1–8, 2008.
- [10] A. Pezeshki, M.R. Azimi-Sadjadi, and L.L. Scharf, "Undersea target classification using canonical correlation analysis," *IEEE Journal of Oceanic Engineering*, vol.32, no.4, pp.948–955, 2007.
- [11] A.A. Nielsen, "Multiset canonical correlations analysis and multispectral, truly multi-temporal remote sensing data," *IEEE Trans IP*, vol.11, no.3, pp.293–305, 2002.
- [12] I. Schizas, G. Giannakis, and Z.Q. Luo, "Distributed estimation using reduced-dimensionality sensor observations," *IEEE Trans. SP*, vol.55, no.8, pp.4284–4299, aug. 2007.
- [13] H. Hotelling, "Analysis of a complex of statistical variables into principal components," *Journal of Educational Psychology*, vol.24, 1933.
- [14] R.A. Fisher, "The use of multiple measurements in taxonomic problems," *Annals Eugen.*, vol.7, pp.179–188, 1936.
- [15] H. Wold, *Estimation of Principal Components and Related Models by Iterative Least squares*, pp.391–420, Academic Press, 1966.
- [16] F. De la Torre and M. Black, "Robust principal component analysis for computer vision," *Proc. ICCV*, pp.362–369, 2001.
- [17] D. Fernando and M.J. Black, "A framework for robust subspace learning," *IJCV*, vol.54, p.2003, 2003.
- [18] K. Inoue, K. Hara, and K. Urahama, "Robust multilinear principal component analysis," *Proc. ICCV*, pp.591–597, 2009.
- [19] J. Lu, K. Plataniotis, and A. Venetsanopoulos, "Regularization studies of linear discriminant analysis in small sample size scenarios with application to face recognition," *Pattern Recognition Letters*, vol.26, no.2, pp.181–191, 2005.
- [20] M. Zhu and A. Martinez, "Subclass discriminant analysis," *IEEE Trans PAMI*, vol.28, no.8, pp.1274–1286, 2006.
- [21] N. Gkalelis, V. Mezaris, and I. Kompatsiaris, "Mixture subclass discriminant analysis," *IEEE Signal Processing Letters*, vol.18, no.5, pp.319–322, 2011.
- [22] B. Schölkopf, A. Smola, and K.R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural Computation*, vol.10, no.5, pp.1299–1319, 1998.
- [23] S. Mika, G. Ratsch, J. Weston, B. Schölkopf, and K.R. Müllers, "Fisher discriminant analysis with kernels," *Proc. NNSP*, pp.41–48, 1999.
- [24] G. Dai, D. yan Yeung, and H. Chang, "Extending kernel Fisher discriminant analysis with the weighted pairwise Chernoff criterion," *Proc. ECCV*, pp.308–320, 2006.
- [25] C. Bishop, "Linear models for regression," in *Pattern Recognition and Machine Learning*, ch. 3, Springer, 2006.
- [26] S. Akaho, "A kernel method for canonical correlation analysis," *Proc. IMPS2001*, 2001.
- [27] S. Harmeling, A. Ziehe, M. Kawanabe, and K.R. Müller, "Kernel-based nonlinear blind source separation," *Neural Computation*, vol.15, pp.1089–1124, 2003.
- [28] A.E. Hoerl, "Application of ridge analysis to regression problem," *Chemical Engineering Progress*, vol.58, pp.54–59, 1962.
- [29] A. Tikhonov, "On the stability of inverse problems," *Dokl. Akad. Nauk SSSR*, vol.39, no.5, pp.195–198, 1943.
- [30] D.R. Hardoon, S. Szedmak, and J. Shawe-Taylor, "Canonical correlation analysis: An overview with application to learning methods," *Neural Computation*, vol.16, no.12, pp.2639–2664, 2004.
- [31] M. Blaschko, C. Lampert, and A. Gretton, "Semi-supervised Laplacian regularization of kernel canonical correlation analysis," *Proc. ECML-PKDD*, pp.133–145, Springer-Verlag, 2008.
- [32] T. Cox and M. Cox, *Multidimensional scaling*, Monographs on statistics and applied probability, no.1, Chapman & Hall, 1994.
- [33] X. He and P. Niyogi, "Locality preserving projections," *Proc. NIPS*, 2003.
- [34] M. Sugiyama, "Dimensionality reduction of multimodal labeled data by local Fisher discriminant analysis," *JMLR*, vol.8, pp.1027–1061, 2007.
- [35] M. Sugiyama, T. Idé, S. Nakajima, and J. Sese, "Semi-supervised local Fisher discriminant analysis for dimensionality reduction," *Machine Learning*, vol.78, no.1–2, pp.35–61, 2010.
- [36] A. Kimura, H. Kameoka, M. Sugiyama, T. Nakano, E. Maeda, H. Sakano, and K. Ishiguro, "SemiCCA: Efficient semi-supervised learning of canonical correlations," *Proc. ICPR*, pp.2933–2936, 2010.
- [37] M. Borga, T. Landelius, and H. Knutsson, "A unified approach to PCA, PLS, MLR and CCA," Report LiTH-ISY-R-1992, ISY, SE-581 83 Linköping, Sweden, November 1997.
- [38] T. De Bie, N. Cristianini, and R. Rosipal, "Eigenproblems in pattern recognition," in *Handbook of Geometric Computing: Applications in Pattern Recognition, Computer Vision, Neural computing, and Robotics*, pp.129–170, Springer, 2005.
- [39] L. Sun, S. Ji, and J. Ye, "A least squares formulation for a class of generalized eigenvalue problems in machine learning," *Proc. ICML*, pp.977–984, 2009.
- [40] L. Sun, S. Ji, and J. Ye, "Canonical correlation analysis for multilabel classification: A least-squares formulation, extensions, and analysis," *IEEE Trans PAMI*, vol.33, pp.194–200, 2011.
- [41] F. De la Torre, "A unification of component analysis methods," in *Handbook of Pattern Recognition and Computer Vision*, ed. C. Chen, ch. 1, pp.3–22, World Scientific Pub Co Inc, 4 ed., 2010.
- [42] F. De la Torre, "A least-squares framework for component analysis," *IEEE Trans PAMI*, vol.34, no.6, pp.1041–1055, 2012.
- [43] H. Yanai and S. Puntanen, "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.
- [44] I.T. Jolliffe, "A note on the use of principal components in regression," *Journal of the Royal Statistical Society*, vol.31, no.3, 1982.
- [45] M. Sugiyama, "Dimensionality reduction of multimodal labeled data by local Fisher discriminant analysis," *JMLR*, vol.8, pp.1027–1061, 2007.
- [46] L. Zelnik-manor and P. Perona, "Self-tuning spectral clustering," *Proc. NIPS*, pp.1601–1608, 2004.