# Designing various component analysis at will

Akisato Kimura, Hitoshi Sakano, Hirokazu Kameoka
*NTT Communication Science Laboratories*
*akisato@ieee.org*

Masashi Sugiyama
*Tokyo Institute of Technology*
*sugi@cs.titech.ac.jp*

## Abstract

*This paper provides a generic framework of component analysis (CA) methods 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 standard CA methods but also (2) several regularization techniques, (3) weighted extensions, (4) some clustering methods, and (5) their semi-supervised extensions. This paper also presents quite a simple methodology for designing a desired CA method from the proposed framework: Adopting the known GPEs as templates, and generating a new method by combining these templates appropriately.*

## 1 Introduction

Component analysis (CA) is traditional, quite simple but might be one of the powerful tools to obtain a hidden structure embedded in the data. Recent reports showed its effectiveness for several tasks in computer vision and pattern recognition. Principal component analysis (PCA), Fisher discriminant analysis (FDA), multiple linear regression (MLR), and canonical correlation analysis (CCA) are well known as standard CA methods [1]. They can be formulated as a generalized eigenvalue problem of a scatter matrix or an augmented matrix composed of several scatter matrices [3, 4]. Kernel CA methods as kernelized extensions of those standard methods have been also developed to deal with non-vector samples and non-linear analysis, which can be formulated as a generalized eigenvalue problem of Gram matrices, instead of scatter matrices [3, 4]. Kernel CA often needs some regularization techniques such as $\ell_2$-norm regularization to inhibit overfitting and Laplacian regularization [2] to fit underlying data manifolds smoothly. In addition, improvements of robustness against outliers and separately distributed samples (e.g. locality preserving projection (LPP) [8] and local FDA (LFDA) [11]) and their extensions to

semi-supervised analysis [12, 2] have been considered.

Although a lot of CA methods and several trials to unify these methods have been presented so far [3, 4, 5, 6], freely designing a tailor-made method of CA 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 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 scatter matrices and Gram matrices, which we call *generalized pairwise expression (GPE)* to make it easy to design a new CA method with desired properties. 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.

## 2 Generalized pairwise expression

Consider two sample sets $\boldsymbol{X} = \{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_{N_x}\}$ and $\boldsymbol{Y} = \{\boldsymbol{y}_1, \ldots, \boldsymbol{y}_N, \boldsymbol{y}_{N_x+1}, \ldots, \boldsymbol{y}_{N_x+N_y-N}\}$ with $N_x$ and $N_y$ samples and $d_x$ and $d_y$ dimensions. For brevity, suppose that both of the sample sets $\boldsymbol{X}$ and $\boldsymbol{Y}$ are centered on the origin, and samples with the same suffix are paired. A pair $(\boldsymbol{X}, \boldsymbol{Y})$ is separated into the following two types: A *complete set* $(\boldsymbol{X}_C, \boldsymbol{Y}_C)$ so that every sample pair $(\boldsymbol{x}_n, \boldsymbol{y}_n)$ co-occurs, and *incomplete sets* $\boldsymbol{X}_I$ and $\boldsymbol{Y}_I$ so that every sample $\boldsymbol{x}_n$ (resp. $\boldsymbol{y}_n$) cannot find the co-occurring sample. Unless otherwise stated, we assume $N_x = N_y = N$.

Many CA methods developed so far involve the following optimization problem:

$$\boldsymbol{w}^{(opt)} \quad := \quad \arg\max_{\boldsymbol{w} \in \mathcal{R}^d} (\boldsymbol{w}^\top \overline{\boldsymbol{C}} \boldsymbol{w})(\boldsymbol{w}^\top \underline{\boldsymbol{C}} \boldsymbol{w})^{-1},$$

where $\overline{C}$ and $\underline{C}$ are symmetric matrices with some statistical natures. Roughly speaking, $\overline{C}$ encodes the quantity that we want to increase, and $\underline{C}$ corresponds to the quantity that we want to decrease. The above optimization can be converted to the following generalized eigenvalue problem via Lagrangian multipliers:

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

The eigenvector $w_k$ $(k = 1, \ldots, r)$ of this generalized eigenvalue problem gives a solution of the original CA.

When addressing CA 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 pairs of features:

$$S_{Q,xy} := \sum_{n=1}^{N}\sum_{m=1}^{N} Q_{n,m}(x_n - x_m)(y_n - y_m)^{\top},$$

where $Q$ is an $N \times N$ non-negative symmetric matrix. A typical example is the scatter matrix $S_{xy} = N^{-1}\sum_{n=1}^{N} x_n y_n^{\top}$. Let $D_Q$ be the $N \times N$ diagonal matrix with $D_{Q,n,n} = \sum_{m=1}^{N} Q_{n,m}$, and $L_Q = D_Q - Q$ (Laplacian). Then, $S_{Q,xy}$ can be expressed as

$$S_{Q,xy} = XL_Q Y^{\top}.$$

The above expression is called the *pairwise expression (PE)* of the second-order statistics $S_{Q,xy}$[12].

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

$$S_{Q,xy} := XL_{Q,1}Y^{\top} + L_2,$$

where $L_{Q,1}$ is a $N \times N$ positive semi-definite matrix, and $L_2$ is a $d_x \times d_y$ matrix.[1] We call this expression the *generalized pairwise expression (GPE)* of the second-order statistics $S_{Q,xy}$. The first and second terms are called the *data term* and *bias term*, respectively.

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

1. If $A$ is a GPE and $\beta > 0$, then $\beta A$ is also a GPE.

2. If $A$ and $B$ are $d_x \times d_y$ GPEs, then $A + B$ is also a GPE with the same size.

3. If $A$ is a $d_x \times d_y$ GPE with columns, and $B$ is a $d_y \times d_z$ GPE, then $AB$ is a $d_x \times d_z$ GPE.

*Proof.* The 1st and 2nd claims can be easily proved, so we concentrate to prove the 3rd one. First, let us define

$$A := XL_{A1}Y^{\top} + L_{A2}, \quad B := YL_{B1}Z^{\top} + L_{B2},$$

---

[1]In general, we do not have to explicitly consider the matrix $Q$ for the following discussions.

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

$$
\begin{aligned}
AB &= (XL_{A1}Y^{\top} + L_{A2})(YL_{B1}Z^{\top} + L_{B2}), \\
&= X(L_{A1}Y^{\top}YL_{B1})Z^{\top} + (L_{A2}Y)L_{B1}Z^{\top} \\
&\quad + XL_{A1}(Y^{\top}L_{B2}) + L_{A2}L_{B2}.
\end{aligned}
$$

Here, we can find some matrices $L_{Ci}$ $(i = 1, 2, 3)$ satisfying $L_{C1} = L_{A1}Y^{\top}YL_{B1}$, $XL_{C2} = L_{A2}Y$, and $L_{C3}Z^{\top} = Y^{\top}L_{B2}$, if $N > \max(d_x, d_y, d_z)$, which implies that for some matrices $L_{D1}$ and $L_{D2}$

$$
\begin{aligned}
AB &= XL_{C1}Z^{\top} + XL_{C2}L_{B1}Z^{\top} \\
&\quad + XL_{A1}L_{C3}Z^{\top} + L_{A2}L_{B2} \\
&= X(L_{C1} + L_{C2}L_{B1} + L_{A1}L_{C3})Z^{\top} + L_{A2}L_{B2} \\
&= XL_{D1}Z^{\top} + L_{D2}.
\end{aligned}
$$

which means $AB$ is also a GPE. □

Thanks to the above properties, **various CA methods can be easily designed by simply combining (i.e. adding, scaling and multiplying) GPEs of existing methods with desired properties.** The rest of the problem is to reveal GPE of existing methods and to present several significant examples that would bring us quite an important hint when constructing new CA methods.

## 3  Reviewing CA methods via GPE

The GPEs of the standard CA methods are listed in Table 1. Several detailed derivations can be seen in [3, 4]. Instead, this paper provides several significant examples that would be quite an important hint when constructing new CA methods.

### 3.1  Semi-supervised local FDA (SELF)

Semi-supervised local FDA (SELF) [12] integrates localized FDA (LFDA) [11] as a supervised CA method and PCA as an unsupervised one. SELF is a typical example for designing CA methods via the GPE framework from the following two viewpoints:

1. Combining several CA methods via GPE,

2. Selecting samples to calculate data terms, a key to extend CA methods to semi-supervised ones.

SELF is effective when (1) we have a complete set $(X_C, Y_C)$, where each sample in $Y$ represents a class ID of the paired sample in $X$, and (2) an incomplete set $X_I$ only exists, namely there are at least one unlabeled samples. In such cases, we can regularize the objective function of the LFDA using the additional

## Table 1. GPEs of standard methods

| Method | $\overline{C}$ | $\underline{C}$ |
|---|---|---|
| PCA | $S_{xx}$ | $I_{d_X}$ |
| FDA | $S_{xx}^{(b)}$ | $S_{xx}^{(w)}$ |
| CCA | $\begin{bmatrix} 0 & S_{xy} \\ S_{yx} & 0 \end{bmatrix}$ | $\begin{bmatrix} S_{xx} & 0 \\ 0 & S_{yy} \end{bmatrix}$ |
| MLR | $\begin{bmatrix} 0 & S_{xy} \\ S_{yx} & 0 \end{bmatrix}$ | $\begin{bmatrix} S_{xx} & 0 \\ 0 & I_{d_y} \end{bmatrix}$ |
| PCR[9] | $\begin{bmatrix} 0 & S_{\hat{x}y} \\ S_{y\hat{x}} & 0 \end{bmatrix}$ | $\begin{bmatrix} S_{\hat{x}\hat{x}} & 0 \\ 0 & I_{d_y} \end{bmatrix}$ |
| OPLS | $S_{xy}S_{xy}^\top$ | $S_{xx}$ |
| Ridge regression | $\begin{bmatrix} 0 & S_{xy} \\ S_{yx} & 0 \end{bmatrix}$ | $\begin{bmatrix} S_{xx}+\delta I_{d_x} & 0 \\ 0 & I_{d_y} \end{bmatrix}$ |
| LPP[8] | $XLX^\top$ | $XDX^\top$ |
| LFDA[11] | $S_{Q,xx}^{(lb)}$ | $S_{Q,xx}^{(lw)}$ |

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

$S_{xx}^{(b)}$ and $S_{xx}^{(w)}$: Between-class and within-class scatter matrices of $X$, $\hat{X} = U_K\Sigma_K V_K^\top$: $K$-rank approximation of $X$ by SVD, $I_d$: $d \times d$ identity matrix, $\delta > 0$: constant, $S_Q^{(b)}$ and $S_Q^{(w)}$: between-class and within-class scatter matrices of $X$ weighted by an $N \times N$ non-negative symmetric matrix.

data $X_I$. In detail, SELF integrates the GPE ($S_{Q,C}^{(lb)}$, $S_{Q,C}^{(lw)}$) of LFDA calculated only from the complete set ($X_C, Y_C$) and the GPE $S_{xx}$ of PCA calculated from all the samples $X$, as follows:

$$\overline{C}^{(SELF)} = \beta S_{Q,C}^{(lb)} + (1-\beta)S_{xx},$$
$$\underline{C}^{(SELF)} = \beta S_{Q,C}^{(lw)} + (1-\beta)I_{d_x},$$

where $\beta$ is a scalar satisfying $0 \leq \beta \leq 1$. When $\beta = 1$, SELF is equivalent to LFDA with only the labeled samples ($X_C, Y_C$). Meanwhile, when $\beta = 0$, SELF is equivalent to PCA with all samples in $X$. In general, SELF inherits the properties of both LFDA and PCA, and their influences can be controlled by the scalar $\beta$.

### 3.2 Semi-supervised CCA (SemiCCA)

In a similar manner to SELF, we can derive a new extension of CCA that can deal with not only a complete set ($X_C, Y_C$) and also incomplete sets ($X_I, Y_I$). This extensive method is generally called SemiCCA [10]. The GPE of SemiCCA can be described as

$$\overline{C}^{(sCCA)} = \beta \begin{bmatrix} 0 & S_{Cxy} \\ S_{Cyx} & 0 \end{bmatrix} + (1-\beta)\begin{bmatrix} S_{xx} & 0 \\ 0 & S_{yy} \end{bmatrix},$$
$$\underline{C}^{(sCCA)} = \beta \begin{bmatrix} S_{Cxx} & 0 \\ 0 & S_{Cyy} \end{bmatrix} + (1-\beta)\begin{bmatrix} I_{d_x} & 0 \\ 0 & I_{d_y} \end{bmatrix}.$$

## Table 2. GPEs of kernelized CA methods

| Method | $\overline{C}$ | $\underline{C}$ |
|---|---|---|
| kPCA | $K_x$ | $I_N$ |
| kFDA | $K_x^{(b)}$ | $K_x^{(w)}$ |
| kCCA | $\begin{bmatrix} 0 & K_xK_y \\ K_yK_x & 0 \end{bmatrix}$ | $\begin{bmatrix} K_x^2 & 0 \\ 0 & K_y^2 \end{bmatrix}$ |
| kMLR | $\begin{bmatrix} 0 & K_xK_y \\ K_yK_x & 0 \end{bmatrix}$ | $\begin{bmatrix} K_x^2 & 0 \\ 0 & I_N \end{bmatrix}$ |
| kCCA+$\ell_2$[7] | $\begin{bmatrix} 0 & K_xK_y \\ K_yK_x & 0 \end{bmatrix}$ | $\begin{bmatrix} K_x^{(\ell_2)} & 0 \\ 0 & K_y^{(\ell_2)} \end{bmatrix}$ |
| L-kCCA[2] | $\begin{bmatrix} 0 & K_xK_y \\ K_yK_x & 0 \end{bmatrix}$ | $\begin{bmatrix} K_x^{(L)} & 0 \\ 0 & K_x^{(L)} \end{bmatrix}$ |
| LE, SC | $L_x$ | $D_x$ |
| LLE | $K_x^{(LL)}K_x^{(LL)\top}$ | $I_N$ |
| NC, nSC | $D_x^{-1/2}L_xD_x^{-1/2}$ | $I_N$ |

L-kCCA: Laplacian-regularized kernel CCA, $K_x$, $K_x$: Gram matrices, $K_x^{(\ell_2)} = K_x^2 + \delta_x K_x$, $K_x^{(L)} = K_x^2 + \gamma_x R_x$, $R_x = K_x L_x K_x$m, LE: Laplacian eigenmap, SC: Spectral clustering, NC: Normalized cuts, nSC: normalized SC, $K_x^{(LL)} = I_N - K_x$

### 3.3 Kernelized extensions

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

## 4 Designing new methods

Summarizing the discussions so far, we describe (1) GPEs of standard CA methods, (2) the way for integrating several GPEs and (3) some semi-supervised extensions by changing samples for calculating GPEs. This section shows that we can easily design new CA methods at will by replicating those steps.

Consider a problem of video categorization, where its training data includes image features $X$, audio features $Y$ and class indexes. Finding appropriate correlations of such three different modals would be still challenging. Here, we consider an integration (CFDA) of CCA and FDA which enables us to extract class-wise differences of multiple feature correlations as well as to achieve discriminative embedding simultaneously.
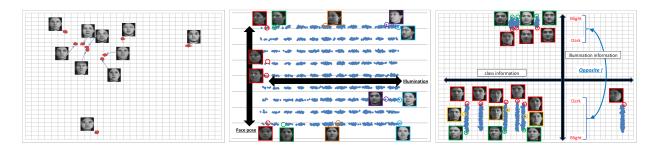
**Figure 1. 2D embedding of MIT CBCL dataset. (Left) FDA with person IDs as classes, (middle) CCA with illuminations and face poses as side info, (right) CCA-FDA with person IDs as classes and illuminations and face poses as side information**

CCA-FDA can be formulated as

$$
\overline{\boldsymbol{C}}_{\boldsymbol{Q}}^{(CFDA)} = \beta \begin{bmatrix} \boldsymbol{0} & \boldsymbol{S}_{xy} \\ \boldsymbol{S}_{yx} & \boldsymbol{0} \end{bmatrix} + (1 - \beta) \boldsymbol{S}_{\boldsymbol{Q}}^{(lb)},
$$

$$
\underline{\boldsymbol{C}}_{\boldsymbol{Q}}^{(CFDA)} = \beta \begin{bmatrix} \boldsymbol{S}_{xx} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{S}_{yy} \end{bmatrix} + (1 - \beta) \boldsymbol{S}_{\boldsymbol{Q}}^{(lw)}.
$$

Figure 1 shows an example how CCA-FDA works with MIT CBCL face dataset, which implies that CCA-FDA obtains a specific property combining CCA and FDA.

Integrating two methods within the kernelized GPE framework is not obvious, since a simple addition of Gram matrices is not a GPE. One example can be seen in a kernelized extension of SELF, called kernel SELF [12]. Remember that the original SELF integrates LFDA with labeled samples and PCA with all the samples (see Section 3.1), and it can be formulated by a localized between-class scatter matrix $\boldsymbol{S}_{\boldsymbol{Q},C}^{(lb)}$, localized within-class matrix $\boldsymbol{S}_{\boldsymbol{Q},C}^{(lw)}$ and the ordinary scatter matrix $\boldsymbol{S}_{xx}$. Kernel SELF can be formulated via their Laplacian matrices $\boldsymbol{L}_{\boldsymbol{Q},C}^{(lb)}$, $\boldsymbol{L}_{\boldsymbol{Q},C}^{(lw)}$, $\boldsymbol{L}_{xx}$, as follows:

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

$$
\underline{\boldsymbol{C}}^{(kSELF)} = \beta \boldsymbol{K}_x \boldsymbol{L}_{\boldsymbol{Q},C}^{(lb)} \boldsymbol{K}_x + (1 - \beta) \boldsymbol{K}_x.
$$

From this formulation, when dealing with kernelized CA, we have to explicitly derive GPEs of existing methods, and replace the data matrix into its Gram matrix.

## 5 Concluding remarks

This paper presented a new expression of scatter matrices and Gram matrices called generalized pairwise expression (GPE). The GPE not only provided a unified insight into various CA methods and their extensions, but also made it eacy to design new CA methods with desired properties. The methodology is quite simple: adopting GPEs of existing methods as templates,

and combining (adding, scaling and multiplying those templates according to the properties you want.

The GPE framework covers a wide variety of CA methods, and thus the way we have presented in this paper for designing new methods is still one example. Developing more general guidelines would be one of the important future work. We will disclose the details just before the conference in arxiv.org.

## References

[1] T. Anderson. *An Introduction to Multivariate Statistical Analysis.* Wiley-Interscience, 2003.
[2] M. Blaschko+. Semi-supervised laplacian regularization of kernel canonical correlation analysis. In *Proc. ECML-PKDD2008*
[3] M. Borga+. A unified approach to PCA, PLS, MLR and CCA. Technical Report LiTH-ISY-R-1992, Linkoping University, 1997.
[4] T. De Bie+. Eigenproblems in pattern recognition. In *Handbook of Geometric Computing: Applications in Pattern Recognition, Computer Vision, Neuralcomputing, and Robotics*, Springer, 2005.
[5] F. De la Torre. A unification of component analysis methods. In *Handbook of Pattern Recognition and Computer Vision*, World Scientific Pub, 2010.
[6] F. De la Torre. A least-squares framework for component analysis. *IEEE Trans. PAMI*, 2012.
[7] D. Hardoon+. Canonical Correlation Analysis: An Overview with Application to Learning Methods. *Neural Computation*, 2004.
[8] X. He and P. Niyogi. Locality Preserving Projections. In *Proc. NIPS2003*.
[9] I. Jolliffe. A Note on the Use of Principal Components in Regression. *JRSS*, Series C (Applied Statistics), 1982.
[10] A. Kimura+. SemiCCA: Efficient semi-supervisedl learning of canonical correlations. In *Proc. ICPR2010*.
[11] M. Sugiyama. Dimensionality Reduction of Multimodal Labeled Data by Local Fisher Discriminant Analysis. *JMLR*, 2007.
[12] M. Sugiyama+. Semi-supervised local Fisher discriminant analysis for dimensionality reduction. *Machine learning*, 2010.