

# SemiCCA: Efficient semi-supervised learning of canonical correlations

AKISATO KIMURA<sup>1,a)</sup> MASASHI SUGIYAMA<sup>2</sup> TAKUHO NAKANO<sup>3,1</sup>  
HIROKAZU KAMEOKA<sup>3,1</sup> HITOSHI SAKANO<sup>1</sup> EISAKU MAEDA<sup>1</sup> KATSUHIKO ISHIGURO<sup>1</sup>

**Abstract:** Canonical correlation analysis (CCA) is a powerful tool for analyzing multi-dimensional paired data. However, CCA tends to perform poorly when the number of paired samples is limited, which is often the case in practice. To cope with this problem, we propose a semi-supervised variant of CCA named *SemiCCA* that allows us to incorporate additional unpaired samples for mitigating overfitting. The main contribution of the proposed method against previously proposed methods is its efficiency and intuitive operation: it smoothly bridges the generalized eigenvalue problems of CCA and principal component analysis (PCA), and thus its solution can be computed efficiently just by solving a single eigenvalue problem as the original CCA.

**Keywords:** Canonical correlation analysis, semi-supervised learning, generalized eigenproblem, principal component analysis, multi-label prediction

## 1. Introduction

The goal of dimensionality reduction is to obtain a low-dimensional representation of high-dimensional data samples, while preserving most of the intrinsic information contained in the original data. If dimensionality reduction is carried out appropriately, the compact representation of the data can be used for various tasks, such as visualization, noise reduction and classification.

Analyzing high-dimensional co-occurring data  $(\mathbf{x}, \mathbf{y})$  is an important challenge in machine learning and pattern recognition communities, e.g., in the context of multi-view learning [1], automatic annotation of music, image and video [2], [3], [4], and sensor data mining [5], [6], [7], [8]. Canonical correlation analysis (CCA) [9] is a classical but still powerful tool for analyzing multivariate paired samples. CCA finds projection bases  $\mathbf{w}_x$  and  $\mathbf{w}_y$  so that correlation between projected samples  $\mathbf{w}_x^\top \mathbf{x}$  and  $\mathbf{w}_y^\top \mathbf{y}$  is maximized. However, the performance of CCA tends to be degraded when the number of paired samples  $(\mathbf{x}, \mathbf{y})$  is limited, while a large number of additional *unpaired* samples (i.e.,  $\mathbf{x}$ -only samples and  $\mathbf{y}$ -only samples) are often available a lot in real-world applications. For example, in the case of automatic image annotation, collecting many labeled images (= paired samples  $(\mathbf{x}, \mathbf{y})$ ) is

often hard, while unlabeled images (= unpaired samples  $\mathbf{x}$ ) can be easily obtained a lot. In the case of sensor data mining, data tends to be lost due to faulty devices and unstable transmissions, which produces a lot of unpaired samples.

To utilize such additional unpaired samples, Blaschko et al. [10] proposed a semi-supervised extension of kernelized CCA [11], [12] by the use of Laplacian regularization. This method enables us to find highly correlated directions that are also located on high variance directions along the data manifold. However, it is specialized to kernelized CCA, and deriving semi-supervised variants of the standard (linear) CCA is not necessarily obvious.

This paper proposes quite a simple method to extend linear CCA to semi-supervised one, that we call *SemiCCA*. The proposed method *SemiCCA* utilizes additional unpaired samples by smoothly bridging CCA and principal component analysis (PCA). More specifically, the generalized eigenvalue problems of CCA and PCA are combined using a trade-off parameter. Thus the solution of *SemiCCA* can still be obtained just by solving the combined eigenvalue problem, which is the same computational complexity as the original CCA.

## 2. Reviewing CCA

Consider a set of paired samples of size  $N$ ,  $\mathbf{X}_P = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$  and  $\mathbf{Y}_P = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N)$ , where each sample is a real-valued vector with dimension  $D_x$  and  $D_y$ . Without loss of generality, we assume that  $\mathbf{X}_P$  and  $\mathbf{Y}_P$  are both centered, that can always be achieved by subtracting the sample means from each sample. CCA is a method of finding a pair  $(\mathbf{w}_x, \mathbf{w}_y) \in \mathcal{R}^{D_x} \times \mathcal{R}^{D_y}$  of basis vectors for

<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.

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

a given set  $(\mathbf{X}_P, \mathbf{Y}_P)$  of paired samples so that their normalized correlation is maximized as follows:

$$\begin{aligned} & \rho(\mathbf{X}_P, \mathbf{Y}_P) \\ &= \max_{(\mathbf{w}_x, \mathbf{w}_y) \in \mathcal{R}^{D_x} \times \mathcal{R}^{D_y}} \frac{\langle \mathbf{X}_P^\top \mathbf{w}_x, \mathbf{Y}_P^\top \mathbf{w}_y \rangle}{\|\mathbf{X}_P^\top \mathbf{w}_x\|_F \cdot \|\mathbf{Y}_P^\top \mathbf{w}_y\|_F} \\ &= \max_{(\mathbf{w}_x, \mathbf{w}_y) \in \mathcal{R}^{D_x} \times \mathcal{R}^{D_y}} \frac{\mathbf{w}_x^\top \mathbf{S}_{Pxy} \mathbf{w}_y}{\sqrt{\mathbf{w}_x^\top \mathbf{S}_{Pxx} \mathbf{w}_x} \sqrt{\mathbf{w}_y^\top \mathbf{S}_{Pyy} \mathbf{w}_y}}, \end{aligned}$$

where  $\langle \mathbf{x}, \mathbf{y} \rangle$  is the inner product of vectors  $\mathbf{x}$  and  $\mathbf{y}$ ,  $\|\mathbf{X}\|_F$  is the Frobenius norm of a matrix  $\mathbf{X}$ ,  $\mathbf{X}^\top$  is a transpose of a matrix  $\mathbf{X}$ ,  $\mathbf{S}_{Pxx}$ ,  $\mathbf{S}_{Pyy}$  and  $\mathbf{S}_{Pxy}$  are sample covariance matrices of paired samples

$$\begin{aligned} \mathbf{S}_{Pxx} &= \mathbf{X}_P \mathbf{X}_P^\top / N, & \mathbf{S}_{Pyy} &= \mathbf{Y}_P \mathbf{Y}_P^\top / N, \\ \mathbf{S}_{Pxy} &= \mathbf{X}_P \mathbf{Y}_P^\top / N. \end{aligned}$$

The maximum of the function  $\rho(\mathbf{X}_P, \mathbf{Y}_P)$  is not affected by re-scaling  $\mathbf{w}_x$  and  $\mathbf{w}_y$  either together or independently. Therefore, the maximization of  $\rho(\mathbf{X}_P, \mathbf{Y}_P)$  is equivalent to the maximizing the numerator  $\mathbf{w}_x^\top \mathbf{S}_{Pxy} \mathbf{w}_y$  of  $\rho(\mathbf{X}_P, \mathbf{Y}_P)$  subject to

$$\mathbf{w}_x^\top \mathbf{S}_{Pxx} \mathbf{w}_x = \mathbf{w}_y^\top \mathbf{S}_{Pyy} \mathbf{w}_y = 1.$$

Taking derivatives of the corresponding Lagrangian with respect to  $\mathbf{w}_x$  and  $\mathbf{w}_y$ , we obtain

$$\begin{aligned} \mathbf{S}_{Pxy} \mathbf{w}_y - \lambda \mathbf{S}_{Pxx} \mathbf{w}_x &= \mathbf{0}, \\ \mathbf{S}_{Pxy}^\top \mathbf{w}_x - \lambda \mathbf{S}_{Pyy} \mathbf{w}_y &= \mathbf{0}. \end{aligned}$$

Therefore, the solution  $(\mathbf{w}_x, \mathbf{w}_y)$  is given as the solution of the following generalized eigenvalue problem:

$$\begin{pmatrix} \mathbf{0} & \mathbf{S}_{Pxy} \\ \mathbf{S}_{Pxy}^\top & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{pmatrix} = \lambda \begin{pmatrix} \mathbf{S}_{Pxx} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{Pyy} \end{pmatrix} \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{pmatrix}. \quad (1)$$

Picking up the top  $D_z$  (should be  $D_z \leq \min(D_x, D_y)$ ) generalized eigenvectors as row vectors, we can obtain  $D_z$ -dimensional mappings  $\mathbf{W}_x$  and  $\mathbf{W}_y$ .

### 3. Proposed method: SemiCCA

#### 3.1 Semi-supervised Setup

As described in the previous section, CCA is only feasible to paired samples. CCA cannot directly deal with unpaired samples. Meanwhile, if only a small number of paired samples are available, CCA tends to overfit them. The main objective of the proposed method *SemiCCA* is to give a new insight that can overcome this weakness.

CCA can be regarded as a method for supervised dimensionality reduction under the task of multi-label prediction, where a companion  $\mathbf{y}$  of a pair is expressed by a binary class vector  $\mathbf{y} \in \{0, 1\}^{D_y}$ . In this set up, SemiCCA is to extend CCA as a supervised method to a semi-supervised one. Namely, preserving the global structure of all the samples in an unsupervised manner can be better than only relying too much on information provided by a small number of labeled samples to avoid overfitting. Not only that, SemiCCA

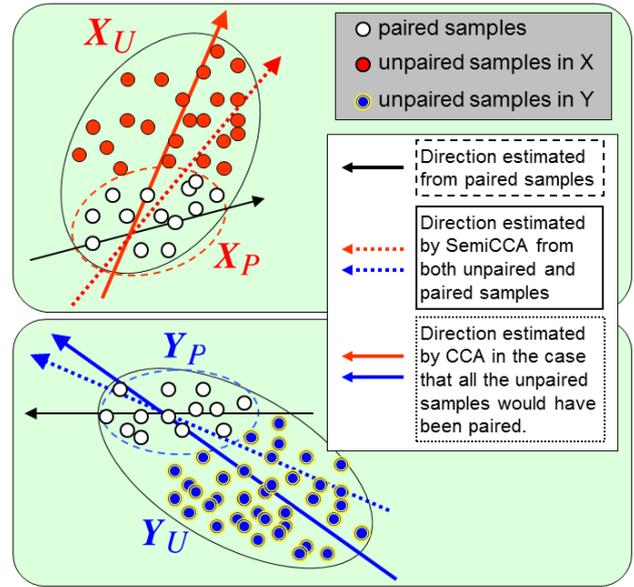


Fig. 1 Effects of unpaired samples in SemiCCA

can be also applied to any types of co-occurring real-valued sample pairs, where several pairs are forced to be unpaired.

Let us explain the idea of SemiCCA using an illustrative two-dimensional data set depicted in Fig. 1, where paired (resp. unpaired) samples are plotted with white (resp. red and blue). When only the paired samples  $(\mathbf{X}_P, \mathbf{Y}_P)$  are used, poor projection bases may be obtained by CCA due to overfitting, as shown at the black arrows in Fig. 1. In contrast, unpaired samples

$$\begin{aligned} \mathbf{X}_U &= (\mathbf{x}_{N+1}, \mathbf{x}_{N+2}, \dots, \mathbf{x}_{N_x}) \\ &= (\mathbf{x}_{U,1}, \mathbf{x}_{U,2}, \dots, \mathbf{x}_{U, N_x - N}), \\ \mathbf{Y}_U &= (\mathbf{y}_{N+1}, \mathbf{y}_{N+2}, \dots, \mathbf{y}_{N_y}) \\ &= (\mathbf{y}_{U,1}, \mathbf{y}_{U,2}, \dots, \mathbf{y}_{U, N_y - N}) \end{aligned}$$

can be used for revealing the global structure in each domain, as shown at the colored arrows in Fig. 1. Note once a basis in one sample space is rectified, the corresponding bases in the other sample space is also rectified so that correlations between two bases are maximized (cf. the dotted arrows in Fig. 1).

#### 3.2 Algorithm

Motivated by the above illustration, we develop a novel method for effectively incorporating unpaired samples into the original CCA. The proposed method, SemiCCA, combines CCA with only the paired samples and principal component analysis (PCA) with all the samples including paired and unpaired samples. More specifically, we integrate the eigenvalue problems of CCA and PCA since this allows us to compute the combined solution efficiently. The solution of SemiCCA is given by the leading generalized eigenvectors of the following generalized eigenvalue problem:

$$\bar{\mathbf{C}} \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{pmatrix} = \lambda \underline{\mathbf{C}} \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{pmatrix}, \quad (2)$$

where

$$\begin{aligned}\bar{\mathbf{C}} &= \beta \begin{pmatrix} \mathbf{0} & \mathbf{S}^{Pxy} \\ \mathbf{S}_{Pxy}^\top & \mathbf{0} \end{pmatrix} + (1 - \beta) \begin{pmatrix} \mathbf{S}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{yy} \end{pmatrix}, \\ \underline{\mathbf{C}} &= \beta \begin{pmatrix} \mathbf{S}^{Pxx} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}^{Pyy} \end{pmatrix} + (1 - \beta) \begin{pmatrix} \mathbf{I}_{D_x} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{D_y} \end{pmatrix},\end{aligned}$$

$\mathbf{S}_{xx}$  and  $\mathbf{S}_{yy}$  are sample covariance matrices of all the pairs

$$\begin{aligned}\mathbf{S}_{xx} &= (\mathbf{X}_P \mathbf{X}_P^\top + \mathbf{X}_U \mathbf{X}_U^\top) / N_x, \\ \mathbf{S}_{yy} &= (\mathbf{Y}_P \mathbf{Y}_P^\top + \mathbf{Y}_U \mathbf{Y}_U^\top) / N_y,\end{aligned}$$

and  $\beta$  is a constant named a *trade-off parameter* taking a value in  $[0, 1]$ . The parameter  $\beta$  controls the trade-off between CCA and PCA. Namely, when  $\beta = 1$ , Eq. (2) is reduced to the CCA eigenvalue problem Eq. (1), while when  $\beta = 0$  Eq. (2) is reduced to the PCA eigenvalue problem, under the assumption that  $\mathbf{X} = (\mathbf{X}_P, \mathbf{X}_U)$  and  $\mathbf{Y} = (\mathbf{Y}_P, \mathbf{Y}_U)$  are uncorrelated. In general, SemiCCA with a trade-off parameter  $0 < \beta < 1$  inherits the properties of both CCA and PCA so that the global structure in each domain and the co-occurrence information of paired samples are smoothly controlled.

One may use different trade-off parameters in  $\bar{\mathbf{C}}$  and  $\underline{\mathbf{C}}$  to increase the flexibility. However, this in turn makes the trade-off parameter choice laborious. For this reason, we focus on using the single shared trade-off parameter  $\beta$  for both  $\bar{\mathbf{C}}$  and  $\underline{\mathbf{C}}$ , as the first step.

### 3.3 Some extensions

We have focused on the case where two sets of samples are given so far. However, the proposed method SemiCCA can be easily extended to multiple data sets by considering correlations over all pairs of samples [13]. For example, we can formulate SemiCCA for a triad  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$  of sample sets, as follows:

$$\bar{\mathbf{C}} \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \\ \mathbf{w}_z \end{pmatrix} = \lambda \underline{\mathbf{C}} \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \\ \mathbf{w}_z \end{pmatrix},$$

where

$$\begin{aligned}\bar{\mathbf{C}} &= \beta \begin{pmatrix} \mathbf{0} & \mathbf{S}^{(P)xy} & \mathbf{S}^{(P)xz} \\ \mathbf{S}_{(P)xy}^\top & \mathbf{0} & \mathbf{S}^{(P)yz} \\ \mathbf{S}_{(P)xz}^\top & \mathbf{S}_{(P)yz}^\top & \mathbf{0} \end{pmatrix} \\ &\quad + (1 - \beta) \begin{pmatrix} \mathbf{S}_{xx} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{yy} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{S}_{zz} \end{pmatrix}, \\ \underline{\mathbf{C}} &= \beta \begin{pmatrix} \mathbf{S}^{(P)xx} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}^{(P)yy} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{S}^{(P)zz} \end{pmatrix} \\ &\quad + (1 - \beta) \begin{pmatrix} \mathbf{I}_{D_x} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{D_y} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_{D_z} \end{pmatrix}.\end{aligned}$$

Of course, the above discussion can be applied to more than 3 sample set in the same way.

We can also obtain a kernelized variant of SemiCCA by

using the standard kernel trick and the technique of pairwise expression [14]. A covariance matrix  $\mathbf{S}_{xy}$  can be converted to the following pairwise expression (see [14] for details):

$$\mathbf{S}_{xy} = \mathbf{X}(\mathbf{D} - \mathbf{W})\mathbf{X}^\top = \mathbf{X}\mathbf{L}\mathbf{X}^\top,$$

where  $\mathbf{W}$  is a matrix so that all the elements are 1,  $\mathbf{D}$  is a diagonal matrix so that the  $n$ -th diagonal element is  $D_{n,n} = \sum_{m=1}^N W_{n,m}$ , and  $\mathbf{L}$  is called a graph Laplacian matrix  $\mathbf{L} = \mathbf{D} - \mathbf{W}$ . In the same way, the identity matrix can be expressed with the following pairwise form:

$$\mathbf{I}_{D_x} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top,$$

where  $\mathbf{X}^\dagger$  denotes the Moore-Penrose generalized inverse of a matrix  $\mathbf{X}$ . Therefore, we can express the eigenvalue problem (Eq. (2)) solved in SemiCCA as

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \bar{\mathbf{L}} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}^\top \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{pmatrix} = \lambda \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \underline{\mathbf{L}} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}^\top \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{pmatrix}, \quad (3)$$

where

$$\begin{aligned}\bar{\mathbf{L}} &= \beta \begin{pmatrix} \mathbf{0} & \mathbf{L}^{(P)xy} \\ \mathbf{L}_{(P)xy}^\top & \mathbf{0} \end{pmatrix} + (1 - \beta) \begin{pmatrix} \mathbf{L}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{L}_{yy} \end{pmatrix}, \\ \underline{\mathbf{L}} &= \begin{pmatrix} \mathbf{L}^{(P)xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{L}^{(P)yy} \end{pmatrix} + (1 - \beta) \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\mathbf{Y}^\top \mathbf{Y})^\dagger \end{pmatrix}.\end{aligned}$$

Here, we introduce the following expressions with appropriate vectors  $\boldsymbol{\alpha}_X, \boldsymbol{\alpha}_Y \in \mathcal{R}^N$  as follows:

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}^\top \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{pmatrix} = \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}^\top \begin{pmatrix} \mathbf{X} \boldsymbol{\alpha}_x \\ \mathbf{Y} \boldsymbol{\alpha}_y \end{pmatrix} = \begin{pmatrix} \mathbf{K}_x \boldsymbol{\alpha}_x \\ \mathbf{K}_y \boldsymbol{\alpha}_y \end{pmatrix}$$

where  $\mathbf{K}_x = \{K_{x(i,j)}\}_{i,j=1}^N$  and  $\mathbf{K}_y = \{K_{y(i,j)}\}_{i,j=1}^N$  are  $N \times N$  matrices with

$$K_{x(i,j)} = \mathbf{x}_i^\top \mathbf{x}_j, \quad K_{y(i,j)} = \mathbf{y}_i^\top \mathbf{y}_j.$$

Then, multiplying Eq. (3) by  $(\mathbf{X}^\top, \mathbf{Y}^\top)$  from the left-hand side yields

$$\begin{pmatrix} \mathbf{K}_x \\ \mathbf{K}_y \end{pmatrix} \bar{\mathbf{L}} \begin{pmatrix} \mathbf{K}_x \\ \mathbf{K}_y \end{pmatrix}^\top \begin{pmatrix} \boldsymbol{\alpha}_x \\ \boldsymbol{\alpha}_y \end{pmatrix} = \lambda \begin{pmatrix} \mathbf{K}_x \\ \mathbf{K}_y \end{pmatrix} \underline{\mathbf{L}} \begin{pmatrix} \mathbf{K}_x \\ \mathbf{K}_y \end{pmatrix}^\top \begin{pmatrix} \boldsymbol{\alpha}_x \\ \boldsymbol{\alpha}_y \end{pmatrix}.$$

The above equation implies that the samples appear only via their inner products, which means that  $\mathbf{K}_x = \{K_{x(i,j)}\}_{i,j=1}^N$  and  $\mathbf{K}_y = \{K_{y(i,j)}\}_{i,j=1}^N$  can be replaced by Gram matrices, each of whose component can be decomposed with a pair  $(\phi_x, \phi_y)$  of functions as

$$\begin{aligned}K_{x(i,j)} &= K_x(\mathbf{x}_i, \mathbf{x}_j) = \phi_x(\mathbf{x}_i)^\top \phi_x(\mathbf{x}_j), \\ K_{y(i,j)} &= K_y(\mathbf{y}_i, \mathbf{y}_j) = \phi_y(\mathbf{y}_i)^\top \phi_y(\mathbf{y}_j).\end{aligned}$$

The kernelized version of SemiCCA can be integrated into the work by Blaschko et al [10] with the introduction of Laplacian regularization to inhibit overfitting.

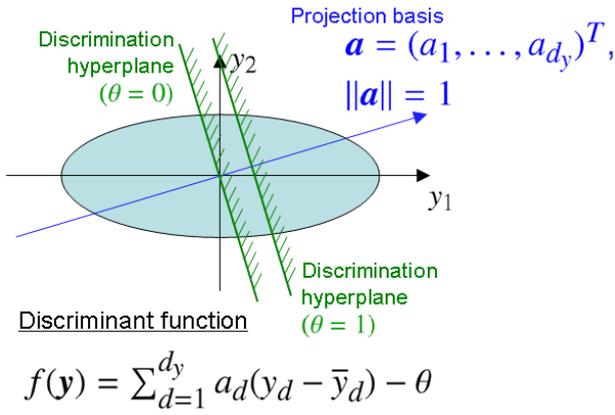


Fig. 2 How to generate artificial data

#### 4. Fundamental evaluations

We first evaluated the performance of the proposed method using the artificial data set created as follows: Consider a simple Gaussian latent model, where the latent random variable (corresponding to a canonical variable in the framework of CCA) is denoted by  $Z$  and the observable random variables are  $X$  and  $Y$ . We drew samples  $\{\mathbf{z}_i\}_{i=1}^{N_z}$  from a standard normal distribution independently,  $\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{D_z})$ , where  $D_z = 10$  is the dimension of the random variable  $Z$ . The number  $N_z$  of samples was set to  $N_z = 10000$ . The means and covariance matrices of the conditional (Gaussian) densities  $p(X|Z)$  and  $p(Y|Z)$  were determined randomly. More specifically, we randomly generated each component of transformation matrices  $\mathbf{T}_x$  and  $\mathbf{T}_y$  and means  $\bar{\mathbf{x}}$  and  $\bar{\mathbf{y}}$  from  $\mathcal{N}(0, 1)$ . Then complete paired samples  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{N_z}$  were created as

$$\begin{aligned} \mathbf{x}_i &= \mathbf{T}_x \mathbf{z}_i + \bar{\mathbf{x}} + \delta_{x,i}, & \delta_{x,i} &\sim \mathcal{N}(\mathbf{0}, \Sigma_{X|Z}), \\ \mathbf{y}_i &= \mathbf{T}_y \mathbf{z}_i + \bar{\mathbf{y}} + \delta_{y,i}, & \delta_{y,i} &\sim \mathcal{N}(\mathbf{0}, \Sigma_{Y|Z}), \end{aligned}$$

where each component of  $\Sigma_{X|Z}$  and  $\Sigma_{Y|Z}$  was generated from the folded standard normal distribution. The dimensions of the samples were set to  $D_x = 15$  and  $D_y = 20$ .

Then, we removed several samples from  $\{\mathbf{y}_i\}_{i=1}^{N_z}$  to artificially generate unpaired samples, as depicted in Fig. 2. Here, we used the following linear discriminant function  $f(\cdot)$  to remove samples:

$$f(\mathbf{y}) = \sum_{d=1}^{d_y} a_d(y_d - \bar{y}_d) - \theta, \quad (4)$$

where  $\mathbf{a} = (a_1, \dots, a_{d_y})^\top$  is a coefficient vector satisfying  $\|\mathbf{a}\| = 1$ , and  $\theta$  is a threshold value such that the larger  $\theta$  is, the more samples are removed. A sample  $(\mathbf{x}_i, \mathbf{y}_i)$  was kept paired if  $f(\mathbf{y}_i) > 0$ , and  $\mathbf{y}_i$  was removed otherwise.

We compare the proposed SemiCCA with the standard CCA. We evaluated the performance of (Semi)CCA by the weighted sum of cosine distances defined as follows:

$$\sum_{i=1}^r \lambda_i^* \frac{\mathbf{w}_{x,i}^\top \mathbf{w}_{x,i}^*}{\|\mathbf{w}_{x,i}\| \cdot \|\mathbf{w}_{x,i}^*\|},$$

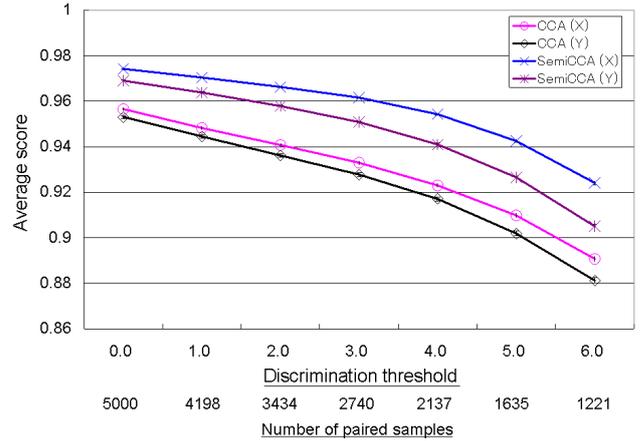


Fig. 3 Average evaluation score for artificial data

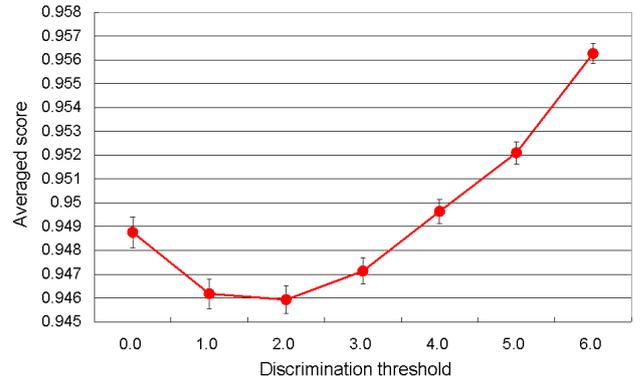


Fig. 4 Average trade-off parameter taking the highest score

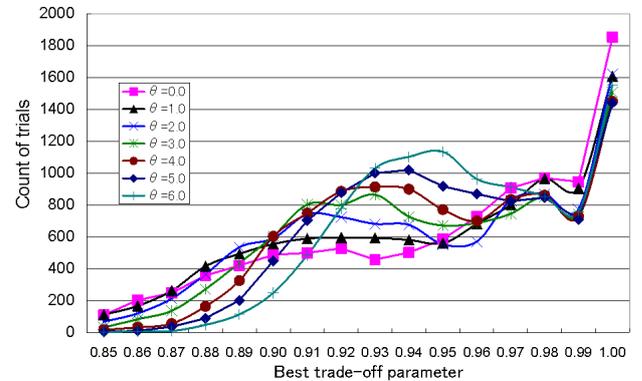


Fig. 5 Histogram of trade-off parameters taking the highest score.

where  $\mathbf{w}_{x,i}^*$  and  $\lambda_i^*$  are the “true” eigenvectors and eigenvalues. We took an oracle setting for selecting the trade-off parameter  $\beta$ . Namely, we adopted the trade-off parameter  $\beta$  marking the highest score for each trial.

Fig. 3 shows the evaluation scores averaged over 10000 independent trials for several discrimination thresholds  $\theta$ , and also shows the average number of paired samples for each discrimination threshold. The results indicate that SemiCCA outperforms the standard CCA; it is note worthy that even when the number of unpaired samples is not so large, SemiCCA performs better than the original CCA.

Fig. 4 shows the trade-off parameter taking the highest score averaged over all the trials, and Fig. 5 depicts the

histogram of the best trade-off parameters. The results imply that the best trade-off parameters have a concave profile with respect to the number of paired samples. Since standard errors of the best trade-off parameters were relatively small, we expect to obtain similar results not only for oracle settings but also for cross validation scenarios. The results also indicate that the best trade-off parameters were usually close to 1, i.e., the effect of PCA is only mildly incorporated. Nevertheless, the performance is much improved, as shown in Fig. 3.

## 5. Applications to multi-label prediction

### 5.1 Method

We applied the proposed method SemiCCA to multi-label prediction, and evaluated its performance with automatic audio annotation. The baseline was proposed by Nakayama et al [15] and Harada et al [16], which is based on a simple latent model with the same structure as *probabilistic Latent Semantic Analysis* (pLSA) [17], [18].

Feature vectors were extracted from audios  $\mathbf{G} = \{\mathbf{g}_n\}_{n=1}^{N_x}$  and associated text labels  $\mathbf{W} = \{\mathbf{w}_n\}_{n=1}^N$ , where  $N$  is the number of labeled samples and  $N_x$  is the total number of samples including labeled and unlabeled samples (should be  $N \leq N_x$  and in most cases  $N \ll N_x$ ). Each text label  $\mathbf{w}_n$  was composed of text words selected from a word set given in advance. We utilized *Bag of Features* (BoF) with dimension  $D_x = 1024$  as audio features  $\mathbf{X} = \{\mathbf{x}_n\}_{n=1}^{N_x}$ . We adopted word existence vectors  $\mathbf{Y} = \{\mathbf{y}_n\}_{n=1}^N$  as text features, where each element represents an existence or absence of a specific word and thus the dimension  $D_y$  of text features was equal to the number of classes.

Next, a latent model was estimated from feature vectors  $(\mathbf{X}, \mathbf{Y})$  with the help of (Semi)CCA. The first step was to generate latent variables  $\mathbf{Z} = \{\mathbf{z}_n\}_{n=1}^{N_x}$  with (Semi)CCA. More specifically, a function  $f_x : \mathcal{R}^{D_x} \rightarrow \mathcal{R}^{D_z}$  was derived from  $(\mathbf{X}, \mathbf{Y})$  as training samples with SemiCCA, and latent variables  $\mathbf{Z}$  are generated from  $(\mathbf{X}, \mathbf{Y})$  with  $f_x$ . The dimension  $D_z$  of latent variables was experimentally determined. We set the function  $f_x$  as  $f_x(\mathbf{x}) = \Lambda^{1/2} \mathbf{W}_x \mathbf{x}$ . The second step was to set up a latent model. The latent model was described by the following equations:

$$\begin{aligned} p(\mathbf{x}, \mathbf{y}) &= \frac{1}{N_x} \sum_{n=1}^{N_x} p(\mathbf{x}|\mathbf{z}_n)p(\mathbf{y}|\mathbf{z}_n), \\ p(\mathbf{x}|\mathbf{z}_n) &\propto \exp\left(-\frac{\|f_x(\mathbf{x}) - \mathbf{z}_n\|^2}{2\gamma^2}\right), \\ p(\mathbf{y}|\mathbf{z}_n) &= \prod_{d=1}^{d_y} p(y_d|\mathbf{z}_n), \\ p(y_d = 1|\mathbf{z}_n) &= \mu\delta(1 - y_{n,d}) + (1 - \mu)N_d/N, \\ p(y_d = 0|\mathbf{z}_n) &= 1 - p(y_d = 1|\mathbf{z}_n), \end{aligned}$$

where  $y_{n,d}$  is the  $d$ -th element of  $\mathbf{y}_n$ ,  $N_d$  is the number of samples containing the  $d$ -th word in labeled samples,  $\mu$  is a parameter representing how reliable a given label is,  $\delta$  is the Kronecker delta, an operator  $\propto$  stands for proportion, and  $\gamma$  is a positive constant. According to the preceding study

[15], we set  $\mu = 0.99$  and  $\gamma = 1.0$ .

Once the model estimation has been finished, we can execute annotation within the same framework through *maximum a posteriori* (MAP) estimation. More specifically, the text feature  $\hat{\mathbf{y}}$  of the most probable text label  $\hat{\mathbf{w}}$  can be derived by using a feature  $\mathbf{x}^{(g)}$  extracted from a given audio  $\mathbf{g}^{(g)}$ , as follows:

$$\begin{aligned} \hat{\mathbf{y}} &= \operatorname{argmax}_{\mathbf{y} \in [0,1]^{D_y}} p(\mathbf{y}|\mathbf{x}^{(g)}) \\ &= \operatorname{argmax}_{\mathbf{y} \in [0,1]^{D_y}} \frac{\sum_{n=1}^{N_x} p(\mathbf{x}^{(g)}|\mathbf{z}_n)p(\mathbf{y}|\mathbf{z}_n)}{\sum_{n=1}^{N_x} p(\mathbf{x}^{(g)}|\mathbf{z}_n)}. \end{aligned}$$

Since a conditional density  $p(\mathbf{y}|\mathbf{z}_n)$  for a text feature  $\mathbf{y}$  is modeled as an element-wise independent distribution

$$p(\mathbf{y}|\mathbf{z}_n) = \prod_{d=1}^{D_y} p(y_d|\mathbf{z}_n),$$

the annotation problem can be rewritten to the following simple form:

$$\hat{y}_d = \frac{\sum_{n=1}^{N_x} p(\mathbf{x}^{(g)}|\mathbf{z}_n)p(y_d = 1|\mathbf{z}_n)}{\sum_{n=1}^{N_x} p(\mathbf{x}^{(g)}|\mathbf{z}_n)}.$$

When  $\hat{y}_d$  exceeds a pre-defined threshold  $\theta_d$ , the text word of index  $d$  is provided to the given image  $\mathbf{g}^{(g)}$ .

### 5.2 Automatic audio annotation

For experiments of automatic audio annotation, we use the data collected from a audio sharing service called *Freesound* <sup>\*1</sup>, which consists of various audio files annotated with word tags such as "people", "noisy", and "restaurant". The goal is to predict the existence of each tag for a new audio clip. We downloaded 2012 audio clips from among all files containing any of pre-defined 230 text labels, 3-60 seconds in length and with a sampling rate 44.1kHz. We then randomly selected 1000 clips as labeled training samples, 912 clips as unlabeled training samples and the rest (= 100 samples) as samples for evaluation. As a fundamental feature comprising a BoF, *Mel-frequency Cepstral Coefficients* (MFCCs), the first and second instantaneous derivatives ( $\Delta$ - and  $\Delta\Delta$ -MFCC) were used for audio frame features. We adopted the precision rate and recall rate as the evaluation measures.

Fig. 6 shows the experimental results, where the horizontal axis stands for the recall rate, and the vertical axis represents the precision rate. We compared the proposed method SemiCCA utilizing both labeled and unlabeled samples with the standard CCA utilizing only labeled samples and the standard CCA in the case that all the unlabeled samples would have been labeled <sup>\*2</sup> Fig. 6 indicates that latent space extraction based on SemiCCA was effective for automatic audio annotation.

<sup>\*1</sup> <http://www.freesound.org>

<sup>\*2</sup> All the audio clips in the dataset used in this experiment have text labels, and a part of them were considered as unlabeled to build a semi-supervised setup.

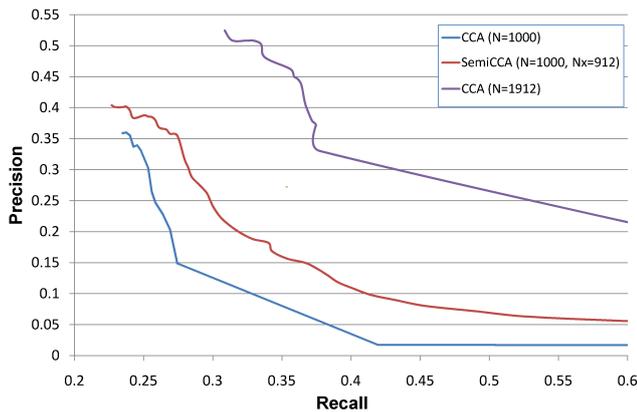


Fig. 6 Precision-recall curve for automatic audio annotation with Freesound dataset.

## 6. Concluding Remarks

In this paper, we proposed a semi-supervised extension of CCA that we call *SemiCCA*. Our formulation is quite simple and also intuitively understandable. Namely, *SemiCCA* smoothly bridges CCA with paired samples and PCA with paired and unpaired samples by a trade-off parameter. We evaluated its experimental performance, and revealed the effectiveness of *SemiCCA* against the original CCA.

In our future work, we will clarify some relationships between the proposed method *SemiCCA* and Bayesian modeling [19], [20], [21], and apply *SemiCCA* to other challenging real-world problems such as multi-modal event correlation analysis for audio-video synchronization, audio-visual speech recognition and sensor data mining.

## References

- [1] 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.
- [2] J.S. Downie, "The music information retrieval evaluation exchange (2005–2007): A window into music information retrieval research," *Acoustical Science and Technology*, vol.29, no.4, pp.247–255, 2008.
- [3] M. Everingham, L. Van Gool, C.K.I. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results." <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>, 2012.
- [4] A.F. Smeaton, P. Over, and W. Kraaij, "Evaluation campaigns and trecvid," *Proc. ACM International Workshop on Multimedia Information Retrieval (MIR)*, pp.321–330, 2006.
- [5] N. Correa, T. Adali, Y.O. Li, and V. Calhoun, "Canonical correlation analysis for data fusion and group inferences," *IEEE Signal Processing Magazine*, vol.27, no.4, pp.39–50, 2010.
- [6] 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.
- [7] A.A. Nielsen, "Multiset canonical correlations analysis and multispectral, truly multi-temporal remote sensing data," *IEEE Transactions on Image Processing*, vol.11, no.3, pp.293–305, 2002.
- [8] I. Schizas, G. Giannakis, and Z.Q. Luo, "Distributed estimation using reduced-dimensionality sensor observations," *IEEE Transactions on Signal Processing*, vol.55, no.8, pp.4284–4299, aug. 2007.
- [9] H. Hotelling, "Analysis of a complex of statistical variables into principal components," *Journal of Educational Psychology*, vol.24, 1933.
- [10] M.B. Blaschko, C.H. Lampert, and A. Gretton, "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.
- [11] P.L. Lai and C. Fyfe, "Kernel and nonlinear canonical correlation analysis," *Proc. International Joint Conference on Neural Networks (IJCNN)*, Los Alamitos, CA, USA, p.4614, 2000.
- [12] S. Akaho, "A kernel method for canonical correlation analysis," *Proc. International Meeting of the Psychometric Society (IMPS)*, 2001.
- [13] 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.
- [14] M. Sugiyama, T. Idé, S. Nakajima, and J. Sese, "Semi-supervised local Fisher discriminant analysis for dimensionality reduction," *Machine Learning*, vol.78, no.1, pp.35–61, 2010.
- [15] H. Nayayama, T. Harada, Y. Kuniyoshi, and N. Otsu, "High-performance image annotation and retrieval for weakly labeled images," *Proc. Pasific-Rim Conference on Multimedia (PCM)*, pp.601–610, 2008.
- [16] T. Harada, H. Nakayama, and Y. Kuniyoshi, "Image annotation retrieval based on efficient learning of contextual latent space," *Proc. IEEE International Conference on Multimedia and Expo (ICME)*, pp.858–861, 2009.
- [17] T. Hofmann, "Probabilistic latent semantic indexing," *Proc. International ACM SIGIR Conference on Research and development in Information Retrieval (SIGIR)*, pp.50–57, 1999.
- [18] T. Hofmann, "Unsupervised learning by probabilistic latent semantic analysis," *Machine Learning*, vol.42, no.1, pp.177–196, 2001.
- [19] F.R. Bach and M.I. Jordan, "A probabilistic interpretation of canonical correlation analysis," *Tech. Rep. 688*, Department of Statistics, University of California, Berkeley, 2005.
- [20] C. Wang, "Variational bayesian approach to canonical correlation analysis," *IEEE Transactions on Neural Networks*, vol.18, no.3, pp.905–910, 2007.
- [21] S. Virtanen, A. Klami, and S. Kaski, "Bayesian cca via group sparsity," *Proc. IEEE International Conference on Machine Learning (ICML)*, pp.457–464, 2011.