Fuzzy Multiple Subspace Fitting for Anomaly Detection

Raissa RELATOR^{†a)}, Nonmember, Tsuyoshi KATO[†], Member, Takuma TOMARU[†], Nonmember, and Naoya OHTA[†], Member

SUMMARY Anomaly detection has several practical applications in different areas, including intrusion detection, image processing, and behavior analysis among others. Several approaches have been developed for this task such as detection by classification, nearest neighbor approach, and clustering. This paper proposes alternative clustering algorithms for the task of anomaly detection. By employing a weighted kernel extension of the least squares fitting of linear manifolds, we develop fuzzy clustering algorithms for kernel manifolds. Experimental results show that the proposed algorithms achieve promising performances compared to hard clustering techniques.

key words: fuzzy algorithm, subspace fitting, kernel vector subspace, kernel affine subspace, anomaly detection

1. Introduction

Anomaly detection is the task of finding patterns within the data that do not conform to the norm. Depending on the specialization domain, the uncovered patterns are usually termed as anomalies, outliers, or novelty points. These often translate to some significant information about the data, which allows us to interpret the data or determine a suitable approach to handle it. Over the years, anomaly detection has gained interest from the scientific community due to its significance in practical applications such as intrusion detection, fraud detection, image processing, sensor networks, traffic networks, and behavior analysis [1]–[4].

Several techniques have been developed addressing anomaly detection problems, depending on their specific application [2], [4]. Some popular straightforward methods include employing classifiers, nearest neighbor-based techniques, and clustering [2], [4]–[6]. While classification is a supervised technique and requires known labels, and nearest neighbor-based classification performs better for semisupervised than unsupervised, clustering may be more advantageous when the data is assumed to be free of anomaly, which is the case for most anomaly detection settings. In this study, clustering is not an objective but a means for fitting multiple simple units to a given dataset to represent a complicated normal class model, as depicted in Fig. 1. The focus is on fuzzy clustering for anomaly detection, in particular, algorithms are developed by applying weighted least squares fitting to update cluster models, and we discuss the

[†]The authors are with the Graduate School of Science and Technology, Gunma University, Kiryu-shi, 376–8515 Japan.

DOI: 10.1587/transinf.2014EDP7027

case where only normal data are given.

Least squares fitting of a linear manifold such as vector subspaces and affine subspaces have been utilized in a variety of applications such as pattern recognition, anomaly detection [2], [4], de-noising [7], data compression [8], and visualization [9]. However, linear manifolds cannot represent curved surfaces, and thus kernel methods have been often employed to address this issue, and other concerns such as handling high-dimensional data. Existing clustering methods related to manifold learning and kernel functions are summarized in Table 1, together with the contributions of this paper. The techniques for fitting vector subspaces and affine subspaces have been used for several decades [10] and were generalized to kernel methods in the late 1990s [11]-[13]. Ho et al. [14] developed a clustering method with vector subspaces, and Lu and Vidal [15] presented its affine subspace version. Ho et al.'s clustering method was kernelized by Li and Fukui [16] and extended to fuzzy clustering by Li et al. [17]. To develop the fuzzy clustering algorithm with vector subspaces, Li et al. devised an elementary technique called the weighted singular value decomposition (WSVD). The WSVD assumes that each entry in the design matrix is to be weighted, unlike in our setting where each input vector is assumed to be weighted. Also, the WSVD is an iterative algorithm and the theoretical guarantee for optimality is not given in their paper, whereas the algorithm presented in this paper is guaranteed to achieve the optimal solution. Although weighted PCA [18] has also been proposed, the deviation for the kernel version of the weighted least squares fitting of the manifold with more than one dimension has not been clearly described.

In this paper, we develop new fuzzy algorithms that learn multiple subspaces, and provide empirical evidence indicating promising performances of the proposed methods in the area of anomaly detection compared to other known algorithms. The rest of the paper is organized as follows. An overview of the problem is presented in Sect. 2. We then proceed to the development of the algorithms for fuzzy clustering of subspaces in Sects. 3 and 4. In Sect. 5, we give the details of the experimental results in applying the aforementioned algorithms to anomaly detection on real-world data including face images, sounds, and amino acid sequences, and provide a performance comparison with hard clustering. Finally, we conclude the paper in Sect. 6. Although all the algorithms presented in this paper are developed in the reproducing kernel Hilbert space, all discussions are de-

Manuscript received January 22, 2014.

Manuscript revised April 17, 2014.

a) E-mail: relator-raissa@kato-lab.cs.gunma-u.ac.jp



(a) Single subspace model.

(b) Multiple subspace model with classical distance

Fig.1 Normal class models. In this study, the normal class is modeled by a set of points $\mathcal{S}(\Theta)$. Classically, the anomalism of an input data point is examined with the distance to its projection onto the point set $S(\Theta)$. Figure (a) describes a model that uses a single affine subspace $S(\Theta) = S_{sa}(U,\mu)$. In (b), the normal class model is given by the union of two affine subspaces $\mathcal{S}(\Theta) = \mathcal{S}_{sa}(U_1, \mu_1) \cup \mathcal{S}_{sa}(U_2, \mu_2)$. The classical distance to the set is the square Euclidean distance to the nearest subspace. In this study, we introduce the κ -distance that is the linear combination of the distances with weights κ_1 and κ_2 , as in (c)

Table 1 Clustering methods with linear manifolds.

	Vector Subspace	Kernel Vector Subspace	Affine Subspace	Kernel Affine Subspace
Single Cluster	Oja [10]	Tsuda [13]	Oja [10]	Maeda and Murase [12]
Hard Clustering	Ho et al. [14]	Li and Fukui [16]	Lu and Vidal [15]	New
Fuzzy Clustering	(Li et al. [17])	New	New	New

scribed in Euclidean space for the sake of readability. We underline that all the proposed algorithms can be kernelized straightforwardly with the kernel trick [19].

Notation. We denote vectors by bold-faced lower-case letters and matrices by bold-faced upper-case letters. Entries of vectors and matrices are not bold-faced. The transpose of a matrix A is denoted by A^{T} , and the inverse of A by A^{-1} . The $n \times n$ identity matrix is denoted by I_n . The ndimensional vector whose entries are all one is denoted by $\mathbf{1}_n$. We use \mathbb{R} and \mathbb{N} to denote the set of real and natural numbers, \mathbb{R}^n and \mathbb{N}^n to denote the set of *n*-dimensional real and natural vectors, and $\mathbb{R}^{m \times n}$ to denote the set of $m \times n$ real matrices. The set of real nonnegative numbers is represented by \mathbb{R}_+ . For any $n \in \mathbb{N}$, we use \mathbb{N}_n to denote the set of natural numbers less than or equal to n. We use \mathbb{S}^n to denote the set of $n \times n$ symmetric matrices, and $\mathbb{O}^{m \times n}$ the set of $m \times n$ *n* orthonormal matrices, i.e. $\mathbb{O}^{m \times n} \equiv \{A \in \mathbb{R}^{m \times n} | A^{\mathsf{T}}A =$ I_n . This definition implies that $\mathbb{O}^{m \times n} = \emptyset$ if m < n. A *permutation* of the set \mathbb{N}_n is a bijective map from \mathbb{N}_n to \mathbb{N}_n . We use \mathcal{P}_n to denote the family of permutations of \mathbb{N}_n . The *n*-dimensional probabilistic simplex is denoted by $\Delta_n \equiv \{x \in A\}$ $\mathbb{R}^n_+ \,|\, \boldsymbol{x}^\mathsf{T} \boldsymbol{1}_n = 1\}.$

2. Problem Setting

Anomaly detection is a task of learning from a set of normal examples. Given ℓ normal examples $x_i \in \mathbb{R}^d$ $(i \in \mathbb{N}_\ell)$, an anomalous pattern is detected if it is distant from the normal class model learned using the examples. This section focuses on the normal class model which is given by a set of points in a *d*-dimensional space, say $\mathcal{S}(\Theta) \subset \mathbb{R}^d$, where Θ is the set of model parameters. If $x \in \mathbb{R}^d$ is an unknown input vector, its square Euclidean distance to the model $\mathcal{S}(\Theta)$ is given by

$$d_{\text{euc}}(\boldsymbol{x}, \mathcal{S}(\boldsymbol{\Theta})) \equiv \min_{\boldsymbol{y} \in \mathcal{S}(\boldsymbol{\Theta})} \|\boldsymbol{x} - \boldsymbol{y}\|^2.$$

This provides the confidence level for the anomalism of x, where larger distances yield higher confidence levels. Instead of the square Euclidean distance, one may also use the Mahalanobis distance [20] or the distance with a learnable Gram matrix [21]. However, these are not rotationally invariant, thus prohibiting us from kernelizing them.

3. **Distance to Normal Class Models**

We now introduce four normal class models wherein values of the parameter Θ are determined such that the mean square deviation of the training examples given by

$$U(\mathbf{\Theta}) \equiv \frac{1}{\ell} \sum_{i=1}^{\ell} d_{\text{euc}}(\mathbf{x}_i, \mathcal{S}(\mathbf{\Theta}))$$
(1)

is minimized in the conventional method.

Single Vector Subspace Model. When we employ the vector subspace with orthonormal bases u_1, \ldots, u_m , the normal class model is given by

$$S(\Theta) = S_{ss}(U) \equiv \{ y \in \mathbb{R}^d \mid \exists \alpha \in \mathbb{R}^m \text{ s.t. } y = U\alpha \},\$$

and $\Theta = U \equiv [u_1, \ldots, u_m] \in \mathbb{O}^{d \times m}$.

Single Affine Subspace Model. The normal class model using a single affine subspace is described as

$$\mathcal{S}(\Theta) = \mathcal{S}_{\mathrm{sa}}(U,\mu) \equiv \{ y \in \mathbb{R}^a \mid \exists \alpha \in \mathbb{R}^m \text{ s.t. } y - \mu = U\alpha \},\$$

where $\Theta = \{U, \mu\}, U \in \mathbb{O}^{d \times m}$ and $\mu \in \mathbb{R}^d$.

Vector Subspace Set Model. This model is defined by a set of L vector subspaces:

$$S(\mathbf{\Theta}) = S_{\mathrm{ms}}(\mathcal{U}) \equiv \bigcup_{k=1}^{L} S_{\mathrm{ss}}(U_k),$$

where $\mathcal{U} \equiv \{U_k\}_{k \in \mathbb{N}_L}$ and $\Theta = \mathcal{U}$. The clustering method of [14] finds a local optimum of \mathcal{U} . Each subspace is expected to represent a cluster of normal class data.

Affine Subspace Set Model. The normal class model can also be defined by a set of *L* affine subspaces,

$$S(\mathbf{\Theta}) = S_{\mathrm{ma}}(\mathcal{U}, M) \equiv \bigcup_{k=1}^{L} S_{\mathrm{sa}}(U_k, \mu_k),$$

where $\mathcal{U} = \{U_k\}_{k \in \mathbb{N}_L}$, $M \equiv [\mu_1, \dots, \mu_L]$, and $\Theta = \{\mathcal{U}, M\}$. A local optimum of Θ is found using a trick similar to the *K*-means method, as in the algorithm of Lu and Vidal [15].

4. Learning Normal Class Models

In this section, we present how to determine the model parameter Θ . First, a conventional method, Hard Clustering, is described in 4.1, and then two Fuzzy Clustering Methods are presented in 4.2 and 4.3.

4.1 Hard Clustering

When using a vector or an affine subspace set (i.e. $S(\Theta) = S_{ms}(U)$ or $S(\Theta) = S_{ma}(U, \mu)$), the distance from the model can be expressed as

$$d_{\text{euc}}(\boldsymbol{x}, \mathcal{S}(\boldsymbol{\Theta})) = \min_{k \in \mathbb{N}_L} d_{\text{euc}}(\boldsymbol{x}, \mathcal{S}_k(\boldsymbol{\Theta}_k)),$$
(2)

where $S_k(\Theta_k)$ is the *k*th cluster model with parameter Θ_k . Meanwhile, the cluster model is given by $S_k(\Theta_k) = S_{ss}(U_k)$ with $\Theta_k = \{U_k\}$, when a single vector subspace is used, and $S_k(\Theta_k) = S_{sa}(U_k, \mu_k)$ with $\Theta_k = \{U_k, \mu_k\}$ when using a single affine subspace. The mean square deviation (1) is rewritten as

$$J(\boldsymbol{\Theta}) = \frac{1}{\ell} \sum_{i=1}^{\ell} \min_{k_i \in \mathbb{N}_L} d_{\text{euc}}(\boldsymbol{x}_i, \mathcal{S}_{k_i}(\boldsymbol{\Theta}_{k_i})).$$

This implies that minimizing $J(\Theta)$ with respect to Θ is equivalent to finding the optimal partitions of training examples, index sets $\{I_k\}_{k=1}^L$, such that $\bigcup_{k=1}^L I_k = \mathbb{N}_\ell$, and

$$\forall k, \forall k' \in \mathbb{N}_L \text{ s.t. } k \neq k' : \ I_k \cap I_{k'} = \emptyset.$$

If $\{\mathcal{I}_k^*\}_{k=1}^L$ is the optimal partition, then

$$\min_{\boldsymbol{\Theta}} J(\boldsymbol{\Theta}) = \frac{1}{\ell} \sum_{k=1}^{L} \min_{\boldsymbol{\Theta}_k} \sum_{i \in I_k^*} d_{\text{euc}}(\boldsymbol{x}_i, \mathcal{S}_k(\boldsymbol{\Theta}_k)),$$

which elucidates that, in learning the model parameters Θ , each example contributes to only one out of *L* cluster models. That is, training examples are not shared among cluster models.

A self-organizing map (SOM) [22] is a clustering

model in which clusters share training examples. Usually, SOM is designed so that each cluster is put on the grid in advance, and each example is shared with the clusters near to the winner cluster in the grid. The cluster models of SOM are, thereby, learned as a 'single connected component.' However, clusters in real world data are not always 'connected'. Motivated by this observation, similarly to fuzzy *c*-means method [23], we now propose two fuzzy algorithms that learn multiple subspaces. In the algorithms, cluster models can share examples, but do not necessarily have to be 'connected.'

4.2 Fuzzy Clustering with κ -Distance

Using the set of predefined weights $\kappa \in \mathbb{R}^L_+$ of the clusters, such that

$$\kappa_1 \ge \kappa_2 \ge \dots \ge \kappa_L \ge 0,\tag{3}$$

we replace the distance in (2) with a linear combination of L distances with weights in κ ,

$$d_{\kappa}(\boldsymbol{x}, \mathcal{S}(\boldsymbol{\Theta})) = \min_{\boldsymbol{\pi} \in \mathcal{P}_L} \sum_{k=1}^{L} \kappa_{\pi(k)} d_{\text{euc}}(\boldsymbol{x}, \mathcal{S}_{\pi(k)}(\boldsymbol{\Theta}_{\pi(k)})), \qquad (4)$$

during learning, where \mathcal{P}_L is the set of permutations on the index set \mathbb{N}_L and $\pi(k)$ returns the index value where $\pi \in \mathcal{P}_L$ maps k. Here, we redefine the objective function $J(\Theta)$ (defined previously in (1)) by replacing d_{euc} with d_{κ} . We refer to this distance as the κ -distance. Since the (4) is equal to (2) when

$$\kappa_k = \begin{cases} 1 & \text{for } k = 1, \\ 0 & \text{for } k \ge 2, \end{cases}$$

the κ -distance is a generalization of the square Euclidean distance in (2). From the assumption in (3), if π^* is the optimal permutation, then

$$d_{\mathrm{euc}}(\boldsymbol{x}, \mathcal{S}_{\pi^*(1)}(\boldsymbol{\Theta}_{\pi^*(1)})) \leq \cdots \leq d_{\mathrm{euc}}(\boldsymbol{x}, \mathcal{S}_{\pi^*(L)}(\boldsymbol{\Theta}_{\pi^*(L)})).$$

Using the optimal permutation $\pi_i^* \forall i \in \mathbb{N}_\ell$, the minimum of the mean square deviation is expressed as

$$\min_{\boldsymbol{\Theta}} J(\boldsymbol{\Theta}) = \frac{1}{\ell} \sum_{k=1}^{L} \min_{\boldsymbol{\Theta}_{k}} \sum_{i=1}^{\ell} v_{k,i} d_{\text{euc}}(\boldsymbol{x}_{i}, \mathcal{S}_{k}(\boldsymbol{\Theta}_{k})), \quad (5)$$

where $v_{k,i}$ is defined as

$$\forall i \in \mathbb{N}_{\ell}, \, \forall k \in \mathbb{N}_L : \quad v_{\pi^*(k),i} \equiv \kappa_k.$$
(6)

The *k*th cluster model is learned from every training example with weight $v_{k,i} > 0$, hence, training examples can be shared with multiple cluster models, resulting to 'naturally connected' cluster models.

The block coordinate ascent method [24] for minimizing $J(\Theta)$ with respect to Θ is given in the following algorithm. **Algorithm 1.** Alternate the following steps until convergence.

Step 1: For each example x_i, i ∈ N_ℓ, update the permutation such that ∀k ∈ N_{L-1}:

$$d_{\text{euc}}(\mathbf{x}_{i}, S_{\pi^{(t+1)}(k)}(\mathbf{\Theta}_{\pi^{(t+1)}(k)})) \\ \leq d_{\text{euc}}(\mathbf{x}_{i}, S_{\pi^{(t+1)}(k+1)}(\mathbf{\Theta}_{\pi^{(t+1)}(k+1)}))$$

by sorting the L distances and compute v^(t+1)_{k,i} using (6).
Step 2: For each cluster model, update the parameters

as: $\forall k \in \mathbb{N}_L$,

$$\mathbf{\Theta}_{k}^{(t+1)} := \arg\min_{\mathbf{\Theta}_{k}} \sum_{i=1}^{\ell} v_{k,i}^{(t+1)} d_{\text{euc}}(\mathbf{x}_{i}, \mathcal{S}_{k}(\mathbf{\Theta}_{k}))$$

The definition of the algorithm guarantees that the sequence $\{J(\Theta^{(t)})\}_{t\in\mathbb{N}}$ is monotonically decreasing. Furthermore, the values of model parameters become unchanged within a finite number of iterations since the cardinality of the permutation set \mathcal{P}_L is finite.

In updating the model parameters in Step 2, Theorem 1 and Corollary 2 (in Appendix) are applied to each cluster model to update Θ_k for the vector subspace set model. For the affine subspace set model, Theorem 4 and Corollary 5 (in Appendix) are employed.

4.3 Fuzzy Clustering with Bezdek Distance

Another approach to sharing the training examples with multiple clusters is to extend the square Euclidean distance (2) to the *Bezdek distance* [23] defined by

$$d_{\text{bez}}(\boldsymbol{x}, \mathcal{S}(\boldsymbol{\Theta})) = \min_{\boldsymbol{w} \in \Delta_L} \sum_{k=1}^L w_k^{b_{\text{bez}}} d_{\text{euc}}(\boldsymbol{x}, \mathcal{S}_k(\boldsymbol{\Theta}_k)),$$
(7)

where $b_{bez} (\ge 1)$ is constant. A similar distance is used in fuzzy *c*-means method [23] and is equivalent to (2) when $b_{bez} = 1$. If $b_{bez} > 1$, the optimal weight w^* can be derived using the method of Lagrange multipliers, and the *k*th entry is given by

$$w_k^* \propto \frac{1}{(d_{\text{euc}}(\boldsymbol{x}, \mathcal{S}_k(\boldsymbol{\Theta}_k)))^{1/(b_{\text{bez}}-1)}}.$$
(8)

The clustering algorithm using the Bezdek distance is given as follows.

Algorithm 2. Alternate the following steps until convergence.

• Step 1: For each example \mathbf{x}_i , $i \in \mathbb{N}_\ell$, update weight $w_i^{(t+1)}$ using (8), and set

$$\forall k \in \mathbb{N}_L, \, \forall i \in \mathbb{N}_\ell : \, v_{k,i}^{(t+1)} = \left([\boldsymbol{w}_i^{(t+1)}]_k \right)^{b_{\text{bez}}}.$$

• Step 2: For each cluster model, update the parameters as: $\forall k \in \mathbb{N}_{I}$,

$$\mathbf{\Theta}_{k}^{(t+1)} := \arg\min_{\mathbf{\Theta}_{k}} \sum_{i=1}^{\ell} v_{k,i}^{(t+1)} d_{\text{euc}}(\mathbf{x}_{i}, \mathcal{S}_{k}(\mathbf{\Theta}_{k}))$$



Fig.2 Examples of face images. (a) Normal data. (b) Anomalous data. (c) Anomalous data that were detected by MA- κ C, but were not detected by other existing methods when the specificity is 0.95. (d) Normal data that was not detected as an anomaly by MA- κ C, but falsely detected by existing methods.

Algorithm 2 is formed such that the sequence $\{J(\mathbf{\Theta}^{(t)})\}_{t\in\mathbb{N}}$ is decreasing monotonically, where $J(\mathbf{\Theta})$ (firstly defined in (1)) is redefined by changing d_{euc} to d_{bez} . Moreover, from the definition of $J(\cdot)$, the sequence is non-negative. Hence, the algorithm must converge.

5. Experiments and Results

To investigate the anomaly detection performance of each model presented in the previous section, experiments on face images, sounds, and string patterns were conducted.

The face images data are from the Extended Yale Face Database B which contains 2,350 192 × 168 gray-scaled images of 39 people. The images of the first three people are used as normal data, and the remaining images are regarded as anomalous data (See Fig. 2 (a),(b)), yielding 192 normal images and 2,158 anomalous images. The kernel function defined by $K(I, I') = \text{tr}(I^T I')$, for any $I, I' \in \mathbb{R}^{192 \times 168}$, is used.

For the second data, sounds are recorded in a bathroom and the power spectra are extracted to obtain 256dimensional input vectors. The linear kernel is used to obtain the kernel values. This dataset contains 1,406 normal data and 94 anomalous data.

Lastly, for string patterns, we used 3,427 amino-acid sequences in 12 folds classified in the protein structure database SCOP [25]. The sequences in the first four folds are assumed to be in the normal class, and the remaining data in the anomalous class, yielding 909 normal data and 2,518 anomalous data. Kernel values are obtained via a string kernel by Lodhi et al. [26] that are capable of efficient inner product computation without explicit extraction of very high-dimensional input vectors.

Eighty percent of the normal data are randomly selected and used as training data, while the remaining normal data together with the anomalous data are used for performance evaluation. For the single vector and single affine subspace models, we determine the number of dimensions m of the manifold so that m is the maximum dimension for which the ratio of the cumulative variances is below 0.95. To obtain the value of m for the subspace set models, we performed single-linkage clustering to divide the training data

(a) Face Images										
	SS	MS-KC	MS-BC	MS-HC	SA	МА-кС	MA-BC	MA-HC		
Highest	0.757 (0.034)	<u>0.954</u> (0.010)	<u>0.952</u> (0.004)	0.921 (0.018)	0.835 (0.028)	0.955 (0.007)	0.950 (0.007)	0.911 (0.040)		
Average	0.757 (0.034)	<u>0.924</u> (0.018)	<u>0.918</u> (0.020)	0.842 (0.026)	0.835 (0.028)	0.936 (0.012)	<u>0.930</u> (0.008)	0.860 (0.025)		
(b) Sounds										
	SS	MS-KC	MS-BC	MS-HC	SA	МА-кС	MA-BC	MA-HC		
Highest	0.846 (0.019)	0.955 (0.007)	0.944 (0.004)	0.964 (0.004)	0.846 (0.019)	0.993 (0.004)	0.993 (0.001)	0.990 (0.006)		
Average	0.846 (0.019)	0.911 (0.014)	0.943 (0.004)	0.908 (0.015)	0.846 (0.019)	<u>0.988</u> (0.005)	0.990 (0.001)	0.976 (0.004)		
(c) Strings										
	SS	MS-KC	MS-BC	MS-HC	SA	МА-кС	MA-BC	MA-HC		
Highest	0.811 (0.014)	0.838 (0.017)	0.841 (0.017)	0.813 (0.015)	0.813 (0.014)	0.836 (0.017)	0.845 (0.017)	0.814 (0.015)		
Average	0.811 (0.014)	0.836 (0.016)	0.834 (0.015)	0.811 (0.015)	0.813 (0.014)	0.831 (0.017)	0.838 (0.016)	0.813 (0.015)		

 Table 2
 Performances on anomaly detections.

into the predetermined number of clusters and searched for the maximum number of dimensions for which the ratio of the cumulative variances is below 0.95 in any cluster. We varied the number of clusters with L = 10, 20, 30. The weights κ in the κ -distance are set as $\kappa_1 = 0.9, \kappa_2 = 0.1$, and $\kappa_3 = 0$, while the parameter of the Bezdek distance is set to $b_{\text{bez}} = 2$.

We tested eight methods: SS, MS- κ C, MS-BC, MS-HC, SA, MA- κ C, MA-BC, and MA-HC. SS and SA, respectively, are the single vector and the single affine subspace models, while MS and MA correspond to the vector and the affine subspace set models. We have κ C, BC, and HC indicating the types of learning methods, where κ C and BC are clustering methods with the κ -distance and the Bezdek distance, and HC is classical hard clustering such as the *K*-means method. The square Euclidean distance is used in the prediction stage, while the κ -distance or the Bezdek distance is employed during learning. Accordingly, MS- κ C, MA- κ C, MS-BC, and MA-BC are the new methods proposed in this paper. Furthermore, to the best of our knowledge, no work uses a kernel version of MA-HC (See Table 1).

The detection performances are summarized in Table 2. The values are the area under the ROC curve (AUC) where 'Highest' refers to the maximum among the three AUCs obtained when L = 10, 20, 30, and 'Average' is their mean. Four random partitions are made dividing the data into training and test sets. The average and the standard deviation of the AUCs are used to determine the performance quality given Table 2. We employed the one-sample *t*-test to detect statistical significance of the differences among the detection performances, and set the significance level to 0.01. In Table 2, bold-faced figures represent the best AUC, and underlined figures indicate performances with no significant difference from the best AUC.

In the first experiment using face images, MA- κ C achieves the best performance when the Highest values are considered in all methods. The Highest AUC of MA-HC is 0.044 lower than that of MA- κ C, however, they are not significantly different since standard deviation for MA-HC is large (0.040). The Highest AUC of MA- κ C is also the best among the eight methods. In general, the values of the hyper-parameters, such as the number of clusters, can be determined by using cross-validation in supervised learning, although the use of the cross validation method is not easy

in the scenario of anomaly detection. For MA-HC, Average is 0.051 lower than Highest, which is larger compared to the difference between Average and Highest in MA- κ C (0.019). For MS- κ C, MS-BC, and MA-BC, respectively, Average and Highest values differ by 0.030, 0.034, and 0.020, while the difference is 0.059 for MS-HC. This suggests that performances using fuzzy clustering do not change drastically with different number of clusters compared to the hard clustering. Similar observations are derived using the dataset of sounds, while MA-BC exceeded the performance of other methods using the string patterns.

For MS- κ C, MS-BC, MA- κ C, and MA-BC, the results presented in Table 2 make use of the square Euclidean distance in prediction, and the fuzzy distances — the κ -distance and the Bezdek distance — in learning, as described previously. Experiments using the generalized distances, d_{κ} and d_{bez} , both in learning and in prediction are also done. The prediction performances when the square Euclidean distance d_{euc} is used with the manifold learned using the generalized distances d_{κ} and d_{bez} are slightly better than the performances of the methods using the κ -distance d_{κ} and Bezdek distance d_{bez} .

Kernel methods are often employed to obtain nonlinear learning machines. Any kernel in our experiments enjoys no nonlinear effect, although the pre-computed kernel values can be transformed easily to nonlinear kernels: $(K(x, x') + c)^p$ for polynomial kernel, and $\exp(-\gamma(K(x, x) + K(x', x') - 2K(x, x')))$ for RBF kernel, where *c*, *p*, and γ are constants. We used these tricks to test the nonlinear kernels, but no improvements were exhibited.

We also applied the one-class SVM to anomaly detection for performance comparison, employing both linear and RBF kernels. Similar settings as above were employed while the value of the parameter ν of the one-class SVM was varied as 0.1, 0.2, ..., 0.7. We report here the best AUC values among all trials for each dataset. For all datasets, Highest AUC is obtained using RBF kernel: 0.868 for the face images, 0.890 for the sound data, and 0.836 for the strings data, and their respective Averages are given by 0.751, 0.867, and 0.789. While the Highest AUC value in the last experiment using string patterns has no significant difference from the best AUC obtained using the proposed algorithm according to statistical tests performed, comparing the said values to the those from the proposed algorithms as presented in Table 2 reveals that the proposed methods outperform one-class SVM.

6. Conclusions

We have developed fuzzy clustering algorithms with vector and affine subspaces in a reproducing kernel Hilbert space. The performances of the algorithms in detecting anomalous patterns proved to be stable against using different number of clusters. Utilization of fuzzy multiple subspace fitting is not limited to anomaly detection, so our technique can also be applied to several other tasks. We consider de-noising as one promising application. Future work includes performance evaluation on other applications.

Acknowledgments

The work of TK is supported by MEXT KAKENHI Grant number 23500373. RR is supported by the Ministry of Education, Culture, Sports, Science and Technology of Japan.

References

- P. Chan, M. Mahoney, and M. Arshad, "Learning rules and clusters for anomaly detection in network traffic," Managing Cyber Threats: Issues, Approaches and Challenges, pp.81–99, 2005.
- [2] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," ACM Comput. Surv., vol.41, pp.15:1–15:58, July 2009.
- [3] J.F. Nieves, "Data clustering for anomaly detection in network intrusion detection," Research Alliance in Math and Science, pp.1–12, Aug. 2009.
- [4] A. Patcha and J.M. Park, "An overview of anomaly detection techniques: Existing solutions and latest technological trends," Comput. Netw., vol.51, pp.3448–3470, Aug. 2007.
- [5] M. Markou and S. Singh, "Novelty detection: A review part 1," Signal Processing, vol.83, pp.2481–2497, 2003.
- [6] M. Markou and S. Singh, "Novelty detection: A review part 2," Signal Processing, vol.83, pp.2499–2521, 2003.
- [7] S. Mika, B. Schölkopf, A. Smola, K.R. Müller, M. Scholz, and G. Rätsch, "Kernel PCA and de-noising in feature spaces," Proc. 1998 Conference on Advances in Neural Information Processing Systems II, pp.536–542, Cambridge, MA, USA, 1999.
- [8] A. Abadpour and S. Kasaei, "Color PCA eigenimages and their application to compression and watermarking," Image Vision Comput., vol.26, pp.878–890, July 2008.
- [9] C.M. Bishop and M.E. Tipping, "A hierarchical latent variable model for data visualization," IEEE Trans. Pattern Anal. Mach. Intell., vol.20, pp.281–293, March 1998.
- [10] E. Oja, "New aspects on the subspace methods of pattern recognition," in Electron. Electr. Eng. Res. Stud. Pattern Recognition and Image Processing Ser. 5, pp.55–64, Letchworth, UK, 1984.
- [11] 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.
- [12] E. Maeda and H. Murase, "Multi-category classification by kernel based nonlinear subspace method," Proc. 1999 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), pp.1025–1028, 1999.
- [13] K. Tsuda, "Subspace classifier in the Hilbert space," Pattern Recognit. Lett., vol.20, pp.513–519, May 1999.
- [14] J. Ho, M.H. Yang, J. Lim, K.C. Lee, and D. Kriegman, "Clustering appearances of objects under varying illumination conditions," CVPR, pp.11–18, 2003.

- [15] L. Lu and R. Vidal, "Combined central and subspace clustering for computer vision applications," Proc. 23rd International Conference on Machine Learning, ICML '06, pp.593–600, New York, NY, USA, 2006.
- [16] X. Li and K. Fukui, "Nonlinear k-subspaces based appearances clustering of objects under varying illumination conditions," ACCV Workshop Subspace 2007, pp.46–52, 2007.
- [17] X. Li, Z. Ning, and L. Xiang, "Robust multi-body motion segmentation based on fuzzy k-subspace clustering," IEICE Trans. Inf. & Syst., vol.E88-D, no.11, pp.2609–2614, Nov. 2005.
- [18] C. Alzate and J.A.K. Suykens, "Multiway spectral clustering with out-of-sample extensions through weighted kernel PCA," IEEE Trans. Pattern Anal. Mach. Intell., vol.32, pp.335–347, Feb. 2010.
- [19] B. Schölkopf and A.J. Smola, Learning with kernels, MIT Press, Cambridge, MA, 2002.
- [20] G.R. Lanckriet, L.E. Ghaoui, and M.I. Jordan, "Robust novelty detection with single-class MPM," In advances in neural information processing systems, pp.905–912, MIT Press, 2003.
- [21] K.Q. Weinberger and L.K. Saul, "Fast solvers and efficient implementations for distance metric learning," Proc. 25th International Conference on Machine Learning, ICML '08, New York, NY, USA, pp.1160–1167, 2008.
- [22] T. Kohonen, M.R. Schroeder, and T.S. Huang, eds., Self-Organizing Maps, 3rd ed., Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2001.
- [23] R.O. Duda, P.E. Hart, and D.G. Stork, Pattern Classification, Second Edition, John Wiley & Sons, 2001.
- [24] J.M. Ortega and W.C. Rheinboldt, Iterative Solution of Nonlinear Equations in Several Variables, Academic Press, New York, NY, 1970.
- [25] A. Andreeva, D. Howorth, S.E. Brenner, T.J.P. Hubbard, C. Chothia, and A.G. Murzin, "SCOP database in 2004: refinements integrate structure and sequence family data," Nuclear Acid Research, vol.32, pp.D226–D229, 2004.
- [26] H. Lodhi, C. Saunders, J. Shawe-Taylor, N. Cristianini, and C. Watkins, "Text classification using string kernels," J. Mach. Learn. Res., vol.2, pp.419–444, March 2002.

Appendix A: Fitting of Subspaces

Weighted least square fitting of cluster component is a necessary step in fuzzy clustering. In this Appendix we give algorithms for weighted least squares fitting of vector and affine subspaces in a kernelizable manner (See Fig. A \cdot 1).

Fitting of Vector Subspaces. The square Euclidean distance from any input vector $\mathbf{x} \in \mathbb{R}^d$ to a linear subspace $S_{ss}(U)$ is given by

$$d(\mathbf{x}, \mathcal{S}_{ss}(\mathbf{U})) = \min_{\mathbf{y} \in \mathcal{S}_{ss}(\mathbf{U})} \|\mathbf{x} - \mathbf{y}\|^{2}$$
$$= \left\| (\mathbf{I}_{d} - \mathbf{U}\mathbf{U}^{\mathsf{T}})\mathbf{x} \right\|^{2}.$$
(A·1)

For a training set $\{x_i\}_{i=1}^{\ell}$, it is known that a value of the parameter $U \in \mathbb{O}^{d \times m}$ minimizing the unweighted mean square deviation is the matrix whose columns are the *m* major eigenvectors of

$$\sum_{i=1}^{\ell} \boldsymbol{x}_i \boldsymbol{x}_i^{\mathsf{T}}.$$

Our goal is to give a kernelizable solution to minimize the weighted mean square deviation



Fig. A•1 Fitting of one-dimensional linear manifolds in a twodimensional toy problem. Plots (a) and (b) are the fitting results using the unweighted least squares and the weighted least squares, respectively. The weights here are given by $\exp(-(x - 0.3)^2/8)$. The darkness of each point indicates the value of the weight. It is observed in (b) that darker points tend to be nearer to the manifold.

$$J_{\rm ss}(\boldsymbol{U};\boldsymbol{v}) \equiv \frac{\sum_{i=1}^{\ell} v_i d(\boldsymbol{x}_i, \boldsymbol{\mathcal{S}}_{\rm ss}(\boldsymbol{U}))}{\sum_{i=1}^{\ell} v_i}$$

with respect to U, where $v_i \in \mathbb{R}_+$ is the weight of the *i*th example x_i .

Let $X \in \mathbb{R}^{d \times \ell}$ be the design matrix

$$X \equiv [x_1,\ldots,x_\ell].$$

For arbitrary vectors $x, y \in \mathbb{R}^d$, the function K(x, y) is given by

$$K(\boldsymbol{x},\boldsymbol{y}) \equiv \langle \boldsymbol{x},\boldsymbol{y} \rangle$$

In kernel methods, various analyses are performed by rewriting the definition of K(x, y) in terms of a positive definite kernel. The inner-product matrix for the input data $\{x_i\}_{i=1}^{\ell}$ is defined as the $\ell \times \ell$ matrix $K \in \mathbb{S}^{\ell}$ whose entries are $K_{ij} \equiv K(x_i, x_j)$. Moreover, such matrices are always positive semidefinite. We also define the vector-valued function $k : \mathbb{R}^{\ell} \mapsto \mathbb{R}^{\ell}$ as

$$\boldsymbol{k}(\boldsymbol{x}) \equiv [\langle \boldsymbol{x}_1, \boldsymbol{x} \rangle, \dots, \langle \boldsymbol{x}_{\ell}, \boldsymbol{x} \rangle]^{\mathsf{T}}.$$

The following theorem gives the algorithm for the weighted least squares estimation of a vector subspace.

Theorem 1. Define a symmetric matrix $\tilde{K} \equiv D_v^{1/2} K D_v^{1/2}$. Assume that $r \ge m$, where $r \equiv \operatorname{rank}(\tilde{K})$. Denote the m major eigenvectors of \tilde{K} by $\lambda_1, \ldots, \lambda_m \in \mathbb{R}_{++}$ and the corresponding eigenvalues by $b_1, \ldots, b_m \in \mathbb{R}^\ell$. Let $\lambda \equiv [\lambda_1, \ldots, \lambda_m]^\mathsf{T}$ and $B \equiv [b_1, \ldots, b_m]$, and define

$$\hat{\boldsymbol{U}}_{ss} \equiv \boldsymbol{X} \boldsymbol{D}_{\boldsymbol{\nu}}^{1/2} \boldsymbol{B} \text{diag}(\boldsymbol{\lambda})^{-1/2}. \tag{A.2}$$

Then

$$\hat{U}_{\mathrm{ss}} \in \operatorname*{arg\,min}_{U \in \mathbb{O}^{d imes m}} J_{\mathrm{ss}}(U; \mathbf{v}).$$

Proof. The weighted mean deviation can be written as

$$J_{\rm ss}(\boldsymbol{U};\boldsymbol{\nu}) = \frac{1}{2} {\rm tr}(\boldsymbol{X}\boldsymbol{D}_{\boldsymbol{\nu}}\boldsymbol{X}^{\mathsf{T}}) - \frac{1}{2} {\rm tr}(\boldsymbol{U}^{\mathsf{T}}\boldsymbol{X}\boldsymbol{D}_{\boldsymbol{\nu}}\boldsymbol{X}^{\mathsf{T}}\boldsymbol{U}). \quad (A \cdot 3)$$

Hence, the columns of the optimal U are the *m* major eigenvectors of a symmetric matrix $XD_{\nu}X^{\mathsf{T}}$. Moreover, each

column lies in the span of the input vectors with positive weights. Therefore, there exists a matrix $A \in \mathbb{R}^{\ell \times m}$ such that $U = XD_{\nu}^{1/2}A$. Substituting this to $J_{ss}(U;\nu)$, the second term in (A·3) becomes

$$\frac{1}{2} \operatorname{tr}(\boldsymbol{U}^{\mathsf{T}} \boldsymbol{X} \boldsymbol{D}_{\boldsymbol{\nu}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{U})$$

= $\frac{1}{2} \operatorname{tr}(\boldsymbol{A}^{\mathsf{T}} \boldsymbol{D}_{\boldsymbol{\nu}}^{1/2} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X} \boldsymbol{D}_{\boldsymbol{\nu}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X} \boldsymbol{D}_{\boldsymbol{\nu}}^{1/2} \boldsymbol{A})$
= $\frac{1}{2} \operatorname{tr}(\boldsymbol{A}^{\mathsf{T}} \tilde{\boldsymbol{K}}^{2} \boldsymbol{A})$

From the orthonormality of U, A must satisfy $A^{\mathsf{T}}\tilde{K}A = I_m$. Hence, the optimal matrix A is the solution of the following constraint optimization problem:

$$\max \operatorname{tr}(A^{\mathsf{T}} \tilde{K}^{2} A)$$
wrt $A \in \mathbb{R}^{\ell \times m}$
subj to $A^{\mathsf{T}} \tilde{K} A = I_{m}$.

Furthermore, the columns of A should be proportional to the m major eigenvectors of \tilde{K} . Thus, we have

$$\forall h \in \mathbb{N}_m, \exists t_h \in \mathbb{R}_+ : \boldsymbol{a}_h = t_h \boldsymbol{b}_h.$$

Now, the columns of U must be unit vectors. And since $\forall h \in \mathbb{N}_m$,

$$= \|\boldsymbol{u}_{h}\|^{2} = \|\boldsymbol{X}\boldsymbol{D}_{v}^{1/2}t_{h}\boldsymbol{b}_{h}\|^{2}$$
$$= t_{h}^{2} \operatorname{tr}(\boldsymbol{X}\boldsymbol{D}_{v}^{1/2}\boldsymbol{b}_{h}\boldsymbol{b}_{h}^{\mathsf{T}}\boldsymbol{D}_{v}^{1/2}\boldsymbol{X}^{\mathsf{T}})$$
$$= t_{h}^{2} \operatorname{tr}(\boldsymbol{b}_{h}^{\mathsf{T}}\tilde{\boldsymbol{K}}\boldsymbol{b}_{h}) = t_{h}^{2}\lambda_{h},$$

then $t_h = \lambda_h^{-1/2}$. Therefore, an optimal value of u_h is obtained such that

$$\forall h \in N_m : \boldsymbol{u}_h = \lambda_h^{-1/2} \boldsymbol{X} \boldsymbol{D}_{\boldsymbol{v}}^{1/2} \boldsymbol{b}_h$$

and the conclusion follows.

To kernelize the method, input vectors are embedded in the inner product. Nevertheless, the input vectors remain explicitly in the result of Theorem 1. Fortunately in many applications, what are important are the distances from arbitrary vectors to the manifold, and not the basis vectors in U.

Corollary 2.

1

$$d(\mathbf{x}, \mathcal{S}_{ss}(\hat{\boldsymbol{U}}_{ss})) = K(\mathbf{x}, \mathbf{x}) - \left\| \operatorname{diag}(\lambda')^{-1/2} \boldsymbol{A} \boldsymbol{D}_{\boldsymbol{v}}^{1/2} \boldsymbol{k}(\mathbf{x}) \right\|^{2}.$$

The equation holds by substituting the result of Theorem 1 into Eq. $(A \cdot 1)$.

Thus, distances to the vector subspace spanned by the columns of the matrix \hat{U}_{ss} do not include any input vector explicitly, enabling us to kernelize the distance to the vector subspace even if the training data are weighted.

Note that the optimal solution is not unique since the set of eigenvectors of a symmetric matrix is not unique. For

instance, if u is an eigenvector corresponding to an eigenvalue λ , then so is the negative of u. Hence, in our problem, UR, where $R \in \mathbb{O}^{m \times m}$, is also an optimal solution whenever U is an optimal solution. The set of optimal solutions are further expanded when the *m*th largest eigenvalue of \tilde{K} is equal to its (m + 1)st largest eigenvalue.

We now build the relationship between Theorem 1 and some classical results. If inputs are given by vectors (e.g. inputs are not given by structured kernels), the non-kernelized method can be used to obtain the vector subspace linear to the input vectors. Furthermore, if $d \ll \ell$, the nonkernelized method is faster: The algorithm in Theorem 1 requires eigen-decomposition of an $\ell \times \ell$ symmetric matrix. The following result shows that the optimal basis vectors of U are also obtained by eigen-decomposition of a $d \times d$ matrix.

Corollary 3. The *m* columns of \hat{U}_{ss} defined in Theorem 1 are the *m* major eigenvectors of

$$\boldsymbol{M} = \sum_{i=1}^{\ell} v_i \boldsymbol{x}_i \boldsymbol{x}_i^{\mathsf{T}}.$$

Any set of the *m* major eigenvectors of *M* minimizes $J_{ss}(U; v)$.

Fitting of Affine Subspaces. In a similar manner as before, we give a method for fitting affine subspaces using weighted least squares. For any $\mathbf{x} \in \mathbb{R}^d$,

$$d(\mathbf{x}, S_{\mathrm{sa}}(U, \boldsymbol{\mu})) = \min_{\mathbf{y} \in S_{\mathrm{sa}}(U, \boldsymbol{\mu})} \|\mathbf{x} - \mathbf{y}\|^2$$
$$= \left\| (I_d - UU^{\mathsf{T}})(\mathbf{x} - \boldsymbol{\mu}) \right\|^2.$$
(A·4)

Given the weight $v_i \in \mathbb{R}_+$ for each data x_i , we wish to find the affine subspace minimizing

$$J_{\rm sa}(\boldsymbol{U},\boldsymbol{\mu};\boldsymbol{\nu}) \equiv \frac{\sum_{i=1}^{\ell} v_i d(\boldsymbol{x}, \mathcal{S}_{\rm sa}(\boldsymbol{U},\boldsymbol{\mu}))}{\sum_{i=1}^{\ell} v_i}.$$

Without loss of generality, we can assume that $\sum_{i=1}^{\ell} v_i = 1$ since the set of optimal solutions is invariant of the scalar product of vectors in \mathbb{R}^{ℓ} . Let us denote the training input vectors shifted with $Xv \in \mathbb{R}^d$ by

$$\forall i \in \mathbb{N}_{\ell} : \ \bar{x}_i \equiv x_i - Xv,$$

and the shifted design matrix by

 $\bar{X} \equiv \left[\bar{x}_1, \ldots, \bar{x}_\ell\right].$

Previously, we defined the functions $K(\cdot, \cdot)$ and $k(\cdot)$, and a matrix K to represent inner products of input vectors. For the shifted input vectors, we introduce two functions, \overline{K} : $\mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$ and $\overline{k} : \mathbb{R}^d \mapsto \mathbb{R}^\ell$, such that

$$\bar{k}(x,y) \equiv \langle x - Xv, y - Xv \rangle, \bar{k}(x) \equiv [\langle \bar{x}_1, x - Xv \rangle, \dots, \langle \bar{x}_\ell, x - Xv \rangle]^{\mathsf{T}}.$$

Then the shifted inner product matrix is given by $\bar{K} \in \mathbb{S}^{\ell}$,

where $\bar{K}_{i,j} = \langle \bar{x}_i, \bar{x}_j \rangle$. The following properties of the functions $\bar{K}(\cdot)$ and $\bar{k}(\cdot)$, and matrix \bar{K} , are useful for kernelization:

$$K(\mathbf{x}, \mathbf{y}) = K(\mathbf{x}, \mathbf{y}) - \langle \mathbf{k}(\mathbf{x}) + \mathbf{k}(\mathbf{y}), \mathbf{v} \rangle + \mathbf{v}^{\mathsf{T}} \mathbf{K} \mathbf{v},$$

$$\bar{\mathbf{k}}(\mathbf{x}) = \mathbf{k}(\mathbf{x}) - \langle \mathbf{k}(\mathbf{x}), \mathbf{v} \rangle \mathbf{1}_{\ell} - \bar{\mathbf{K}} \mathbf{v} + \langle \mathbf{v}, \bar{\mathbf{K}} \mathbf{v} \rangle \mathbf{1}_{\ell},$$

$$\bar{\mathbf{K}} = (\mathbf{I} - \mathbf{1}_{\ell} \mathbf{v}^{\mathsf{T}}) \mathbf{K} (\mathbf{I} - \mathbf{v} \mathbf{1}_{\ell}^{\mathsf{T}}).$$

Finally, we have the following weighted least squares estimation and the distance to the optimal affine manifold in an analytic and kernelizable form.

Theorem 4. Define a symmetric matrix $\bar{\mathbf{K}} \equiv \mathbf{D}_{v}^{1/2} \bar{\mathbf{K}} \mathbf{D}_{v}^{1/2}$. Suppose $r \geq m$, where $\bar{r} \equiv \operatorname{rank}(\bar{\mathbf{K}})$. Denote the m major eigenvectors of $\bar{\mathbf{K}}$ by $\bar{\lambda}_{1}, \ldots, \bar{\lambda}_{m} \in \mathbb{R}_{++}$ and the corresponding eigenvalues by $\bar{\mathbf{b}}_{1}, \ldots, \bar{\mathbf{b}}_{m} \in \mathbb{R}^{\ell}$. Let $\bar{\lambda} \equiv [\bar{\lambda}_{1}, \ldots, \bar{\lambda}_{m}]^{\mathsf{T}}$ and $\bar{\mathbf{B}} \equiv [\bar{\mathbf{b}}_{1}, \ldots, \bar{\mathbf{b}}_{m}]$, and define

$$\hat{U}_{\mathrm{sa}} \equiv X D_{\nu}^{1/2} \bar{B} \mathrm{diag}(\bar{\lambda})^{-1/2} \text{ and } \hat{\mu}_{\mathrm{sa}} \equiv X \nu.$$

Then

$$(\hat{U}_{\mathrm{sa}}, \hat{\mu}_{\mathrm{sa}}) \in \operatorname*{arg\,min}_{U \in \mathbb{O}^{d \times m}, \mu \in \mathbb{R}^d} J_{\mathrm{sa}}(U, \mu; \nu).$$

Proof. Substituting (A·4) into $J_{sa}(U, \mu; v)$ gives us

$$J_{\rm sa}(\boldsymbol{U},\boldsymbol{\mu};\boldsymbol{\nu}) = \frac{1}{2} \sum_{i=1}^{\ell} v_i \left\| (\boldsymbol{I}_d - \boldsymbol{U}\boldsymbol{U}^{\mathsf{T}})(\boldsymbol{x}_i - \boldsymbol{\mu}) \right\|^2, \qquad (A \cdot 5)$$

with the assumption that $\sum_{i=1}^{\ell} v_i = 1$. To minimize this, we set

$$\frac{\partial J_{\mathrm{sa}}(\boldsymbol{U},\boldsymbol{\mu})}{\partial \boldsymbol{\mu}} = (\boldsymbol{I}_d - \boldsymbol{U}\boldsymbol{U}^{\mathsf{T}}) \sum_{i=1}^{\ell} v_i(\boldsymbol{\mu} - \boldsymbol{x}_i) = 0,$$

and obtain the solutions

$$\forall \boldsymbol{\beta} \in \mathbb{R}^m, \qquad \boldsymbol{\mu} = \boldsymbol{X}\boldsymbol{\nu} + \boldsymbol{U}\boldsymbol{\beta}. \tag{A.6}$$

The result of Theorem 4 follows when $\beta = \mathbf{0}_m$. Substituting (A· 6) into (A· 5) yields

$$J_{\rm sa}(\boldsymbol{U}, \boldsymbol{X}\boldsymbol{v}; \boldsymbol{v}) = \frac{1}{2} \sum_{i=1}^{\ell} v_i \| (\boldsymbol{I}_d - \boldsymbol{U}\boldsymbol{U}^{\mathsf{T}}) \bar{\boldsymbol{x}}_i \|^2,$$

which is equal to $J_{ss}(U; v)$ when the \bar{x}_i 's are treated as input vectors. By applying Theorem 1, the theorem is established.

Corollary 5.

$$d(\mathbf{x}, \mathcal{S}_{sa}(\boldsymbol{U}_{sa}, \hat{\boldsymbol{\mu}}_{sa})) = \bar{K}(\mathbf{x}, \mathbf{x}) - \left\| \operatorname{diag}(\boldsymbol{\lambda})^{-1/2} \boldsymbol{B} \boldsymbol{D}_{\boldsymbol{\nu}}^{1/2} \bar{\boldsymbol{k}}(\mathbf{x}) \right\|^{2}.$$

The equation follows when the result of Theorem 4 is substituted in Eq. $(A \cdot 4)$.

For completeness, we give the non-kernelized form in the following corollary.

Corollary 6. The *m* columns in \hat{U}_{sa} defined in Theorem 4 are *m* major eigenvectors of

$$\bar{\boldsymbol{M}} = \sum_{i=1}^{c} v_i (\boldsymbol{x}_i - \hat{\boldsymbol{\mu}}_{\mathrm{sa}}) (\boldsymbol{x}_i - \hat{\boldsymbol{\mu}}_{\mathrm{sa}})^{\mathsf{T}}.$$

0

Any set of *m* major eigenvectors of \overline{M} minimizes $J_{sa}(U, \hat{\mu}_{sa}; v)$.



Naoya Ohta received the B.E. from Tokyo University of Agriculture and Technology in 1983, the M.E. from Tokyo Institute Technology in 1985, and the PhD degree from the University of Tokyo in 1998. From 1987, he worked for NEC Corporation as a researcher in the Central Research Laboratories. He was a research affiliate of the Media Laboratory of MIT from 1991 to 1992. In 1994, he moved to Gunma University and is currently a professor at the Graduate School of Science and Engineering, Gunma

University. His scientific interests include image processing, computer vision, and pattern recognition.



Raissa Relator received the B.S. and M.S. degrees in Mathematics from the University of the Philippines Diliman in 2005 and 2008, respectively. She is currently pursuing a Ph.D. in Computer Science at the Graduate School of Science and Engineering, Gunma University. Her research interests include pattern recognition, machine learning, and bioinformatics.



Tsuyoshi Kato received the B.E., M.E., and Ph.D. degree from Tohoku University, Sendai, Japan, in 1998, 2000, and 2003, respectively. From 2003 to 2005, he was with the National Institute of Advanced Industrial Science Technology (AIST) as a postdoctoral fellow in the Computational Biology Research Center (CBRC) in Tokyo. From 2005 to 2008, he was an assistant professor at the Graduate School of Frontier Sciences, University of Tokyo. From 2008 to 2010, he was an associate professor at the Center

for Informational Biology, Ochanomizu University. He then moved back to Graduate School of Frontier Sciences, University of Tokyo, and as of present, he is an associate professor at the Graduate School of Science and Engineering, Gunma University. His current scientific interests include pattern recognition, computer vision and bioinformatics. He is a member of IEICE Japan.



Takuma Tomarureceived the B.E. andM.E. from Gunma University, Gunma, Japan, in2009 and 2011, respectively. He was engaged inresearch on pattern recognition problems, especially sound recognition. He currently works forTaiyo Yuden Co., Ltd..