A Fast Incremental Kernel Principal Component Analysis for Data Streams

Annie anak Joseph and Seiichi Ozawa

Abstract-Kernel Principal Component Analysis (KPCA) is widely used feature extraction as it have been proven that KPCA is powerful in many areas in pattern recognition. Considering that the conventional KPCA should decompose a kernel matrix of all training data, this would be an unrealistic assumption for data streams in real-world applications. Therefore, in this paper, we propose an online feature extraction called Chunk Incremental Kernel Principal Component Analysis (CIKPCA) that can handle data streams in an incremental mode. In the proposed method, the training data are assumed to be given in a chunk of multiple data at one time. In CIKPCA, an eigen-feature space is updated by solving the eigenvalue decomposition once whenever a chunk of data is given. However, if a chunk size is large, a kernel matrix to be decomposed is also large, resulting in high computational time. Considering that not all the data are useful for the eigen-feature space learning, the data in a chunk are first selected based on the importance. Several benchmark data sets in the UCI Machine Learning Repository are used to evaluate the performance of the proposed method. The experimental results show that our proposed method can accelerate the learning of the eigenfeature space compared to Takeuchi et al.'s IKPCA without reducing the recognition accuracy.

I. INTRODUCTION

Along with the rapid development of the technologies in real-world applications, data are often continuously generated in our life. This kind of data are called "data streams". Therefore, the importance of incremental learning is increasing due to the demand of online learning in real-time recognition is increased. An Incremental learning is not only important in classifier but also for feature extraction part. Principal Component Analysis (PCA) [1] has been one of the well-known feature extraction method in machine learning. Therefore, there have been proposed many extended works on PCA [2], [3], [4], [5], [6], [7]. One of them is Kernel Principal Component Analysis (KPCA) [8] which performs a nonlinear form of PCA. The main idea of KPCA is to map an input data into a higher dimensional feature space by using some specially designed mapping function ϕ . Then, principal components are computed in a very high or infinity dimensional feature space to approximate the data distribution by a low dimensional subspace called "eigenfeature space". However, the traditional batch KPCA has its own limitation where the computation time is increased with the increasing number of training data. This is because the size of a kernel matrix to be decomposed is increased for the larger size of the training data. Hence, the larger the

training data are, the higher the computation time are. In this case, an incremental learning algorithm, which can update an eigenspace model without keeping all the training data, is solicited.

Recently, Chin and Suter [5] proposed an incremental KPCA that can update an eigen-feature space incrementally by using a singular value decomposition. However, their proposed IKPCA requires large computation and memory costs for the high-dimensional data. On the other hand, Takeuchi et al. then proposed another incremental KPCA (IKPCA) [9] where they can approximate an eigen-feature space using linearly independent data for every incoming data. However, their proposed method still remain an open problem. Considering that the training data could appear in a chunk, IKPCA should be repeatedly applied to the individual data in a chunk. Therefore, the update of the eigen-feature space should be conducted for every incoming data. In order to overcome this problem, when a chunk of multiple data is given, these data should be learned at the same time. However, when a large chunk of data is given, the computational time of the independent data selection could dominate over those of the eigenvalue decomposition. Hence, the large computation and memory costs still remained an issue for the large chunk data [10].

To tackle this problem, we propose an extended approach to Takeuchi et al.'s IKPCA. Here, we assume that not all of the data are useful for the eigen-feature space learning. Therefore, when a chunk of multiple data is given, the data are first subdivided into small chunks and the useful data are selected from the small chunks. It is expected that the computation and memory costs could be reduced when the training data in a chunk are reduced. Then, the linearly independent data are selected from the reduced data. The eigenvectors are represented by the linear sum of the linearly independent data which are selected from the reduced data in a chunk. This is because the eigenvectors are not obtained directly in a high-dimensional feature space.

The rest of this paper is organized as follows: In Section II, we briefly explain KPCA and Takeuchi et al.'s IKPCA. Then, we proposed another incremental KPCA by extending Takeuchi et al.'s IKPCA called Chunk IKPCA (CIKPCA) in Section III. In Section IV, the performance of the proposed CIKPCA is evaluated by using several benchmark data sets from the UCI machine learning repository [11], and we give conclusions in Section V.

Annie anak Joseph and Seiichi Ozawa are with the Graduate School of Engineering, Kobe university, Kobe 657-8501, Japan (email: 097t805t@stu.kobe-u.ac.jp, ozawasei@kobe-u.ac.jp).

II. KERNEL PRINCIPAL COMPONENT ANALYSIS (KPCA)

Assume that N training data $\mathbf{X} = \{(\mathbf{x}_i, d_i)\}_{i=1}^N$ are given initially where \mathbf{x}_i and d_i are the *i*th input vector and its class label, respectively. These data are mapped into a feature space using a mapping function $\phi(\cdot)$. Here, $\phi(\mathbf{x})$ is the *l*dimensional data. A covariance matrix \mathbf{Q} of data in the feature space is given by

$$\boldsymbol{Q} = \frac{1}{N} \sum_{i=1}^{N} (\phi(\boldsymbol{x}_i) - \boldsymbol{c}) (\phi(\boldsymbol{x}_i) - \boldsymbol{c})^T$$
(1)

where c is a mean vector of the mapped data and it is defined as

$$\boldsymbol{c} = \frac{1}{N} \sum_{i=1}^{N} \phi(\boldsymbol{x}_i).$$
⁽²⁾

In KPCA, we have to find eigenvalues Λ and the associated eigenvectors Z. These are obtained by solving the following eigenvalue problem:

$$QZ = Z\Lambda.$$
 (3)

Since the dimensionality l could be infinity, the eigenvalue problem is not solved directly. Therefore, the so called kernel trick is applied to Eq. (3) by using the following kernel function:

$$k(\boldsymbol{x}, \boldsymbol{x}') = \phi(\boldsymbol{x})^T \phi(\boldsymbol{x}'). \tag{4}$$

Eigenvector matrix Z is represented as follows:

$$\boldsymbol{Z} = [\phi(\boldsymbol{x}_1), \cdots, \phi(\boldsymbol{x}_m)] \begin{bmatrix} \alpha_{1,1} & \cdots & \alpha_{1,m} \\ \vdots & \ddots & \vdots \\ \alpha_{m,1} & \cdots & \alpha_{m,m} \end{bmatrix} = \boldsymbol{\Phi}_m \boldsymbol{A}_m$$
(5)

where m is a set of linearly independent data that span a feature space. The kernel matrix can be represented by

$$\boldsymbol{H}_{ij} = \begin{bmatrix} k(\boldsymbol{x}_1, \boldsymbol{x}_1) & \cdots & k(\boldsymbol{x}_1, \boldsymbol{x}_j) \\ \vdots & \ddots & \vdots \\ k(\boldsymbol{x}_i, \boldsymbol{x}_1) & \cdots & k(\boldsymbol{x}_i, \boldsymbol{x}_j) \end{bmatrix}.$$
(6)

Finally, a kernel eigenvalue problem [12] is given by

$$\frac{1}{N}\boldsymbol{L}^{-1}\boldsymbol{H}_{Nm}^{T}(\boldsymbol{I}_{N}-\boldsymbol{1}_{N})\boldsymbol{H}_{Nm}(\boldsymbol{L}^{-1})^{T}(\boldsymbol{L}^{T}\boldsymbol{\alpha}_{i}) = \lambda_{i}(\boldsymbol{L}^{T}\boldsymbol{\alpha}_{i}).$$
(7)

Here, λ_i is an eigenvalue ($\lambda_i \geq 0$) and $\boldsymbol{L}^T \boldsymbol{\alpha}_i$ is the associated *i*th eigenvector spanning a feature space. \boldsymbol{I}_N is an $N \times N$ unit matrix and $\boldsymbol{1}_N$ is an $N \times N$ matrix with all elements are 1/N. \boldsymbol{L} is a lower triangular matrix and it is obtained by the Cholesky factorization, $\boldsymbol{H}_{mm} = \boldsymbol{L}\boldsymbol{L}^T$.

From the m eigenvectors, only the first d largest eigenvalues are selected to form a low dimensional eigen-feature space. The criterion to determine the augmentation of the eigen-feature space is by measuring the accumulation ratio below:

$$C(d) = \frac{\sum_{i=1}^{d} \lambda_i}{\sum_{i=1}^{m} \lambda_i}.$$
(8)

The accumulation ratio C(d) indicates how much the information remains in the feature space after the d principal components are selected.

In Takeuchi et al.'s incremental KPCA, the learning is carried out for every incoming data. In this case, a kernel eigenvalue problem needs to be solved one by one. Hence, the computation time are considered to be high. Therefore, more efficient algorithm is needed for data streams under realistic environments.

III. PROPOSED CHUNK IKPCA

The novelty of this paper lies in the reduction of computation and memory costs by introducing a mechanism of data selection in a chunk. Before an eigen-feature space augmentation, a chunk of data should be first subdivided into smaller chunks and the useful data are selected from the smaller chunks. In the following, the five main operations of the proposed method are described in details. In the final part of this section, the algorithm of the proposed method is summarized.

A. Data Selection

Before an incremental learning is carried out, batch KPCA is applied to the initial data $oldsymbol{X} = \{oldsymbol{x}_i\}_{i=1}^N.$ Then an eigenfeature space is constructed based on the initial data. The incremental learning is conducted to update the eigen-feature space when a chunk of new data $Y = \{y_i\}_{i=1}^{L}$ is given. The data in a chunk should be first selected based on the accumulation ratio before updating the eigen-feature space. Here, the accumulation ratio is used to judge whether a new eigenaxis should be added. It is the ratio of amount of information between eigenaxes and the original feature space. However, large computation and memory costs are needed to obtain an accumulation ratio if a chunk size is large. Takaomi et al. [10] investigated the influence of chunk size to the learning time. They found that more time is required for the large chunk data unless the large chunk data are divided into small chunks. Therefore, to overcome the problem of large computational time, when a new chunk of data is given, the data are first subdivided into smaller chunks. For the notational simplicity, let us denote a divided chunk as $y = \{y_i\}_{i=1}^l$. Then, the useful data are selected from y based on the accumulation ratio. Here, the accumulation ratio is updated with the data in a smaller chunk as follows:

$$C'(d) = \frac{\sum_{i=1}^{d} \lambda'_i}{\sum_{i=1}^{m} \lambda'_i},\tag{9}$$

where

$$\sum_{i=1}^{d} \lambda'_{i} = \frac{N}{N+\tilde{l}} \sum_{i=1}^{d} \lambda_{i} + \frac{N\tilde{l}}{(N+\tilde{l})^{2}} \sum_{i=1}^{d} \left\{ \boldsymbol{\alpha}_{i}^{T} (\tilde{\boldsymbol{c}} - \tilde{\boldsymbol{c}}_{y}) \right\}^{2} + \frac{1}{N+\tilde{l}} \sum_{i=1}^{d} \sum_{j=1}^{\tilde{l}} \left\{ \boldsymbol{\alpha}_{i}^{T} (\boldsymbol{K}_{m}(\boldsymbol{y}_{j}) - \tilde{\boldsymbol{c}}) \right\}^{2}$$
(10)

$$\sum_{i=1}^{m} \lambda'_{i} = \frac{N}{N+\tilde{l}} \sum_{i=1}^{m} \lambda_{i} + \frac{N\tilde{l}}{(N+\tilde{l})^{2}} \boldsymbol{c}^{T} \boldsymbol{c}$$
$$- \frac{\tilde{l}}{(N+\tilde{l})^{2}} \sum_{i=1}^{\tilde{l}} \sum_{j=1}^{\tilde{l}} k(\boldsymbol{y}_{i}, \boldsymbol{y}_{j})$$
$$- \frac{2N}{(N+\tilde{l})^{2}} \sum_{i=1}^{\tilde{l}} \boldsymbol{\beta}_{i}^{T} \tilde{\boldsymbol{c}}$$
$$+ \frac{1}{N+\tilde{l}} \sum_{j=1}^{\tilde{l}} k(\boldsymbol{y}_{j}, \boldsymbol{y}_{j}).$$
(11)

Here, \tilde{c} , $\tilde{c_y}$, $K_m(y_i)$ and β_i are defined by

•

$$\boldsymbol{\Phi}_{m}^{T}\boldsymbol{c} = \frac{1}{N} \left[\sum_{i=1}^{N} k(\hat{\boldsymbol{x}}_{1}, \boldsymbol{x}_{i}), \cdots, \sum_{i=1}^{N} k(\hat{\boldsymbol{x}}_{m}, \boldsymbol{x}_{i}) \right]^{T}$$
$$\equiv \tilde{\boldsymbol{c}}$$
(12)

$$\boldsymbol{\Phi}_{m}^{T}\boldsymbol{c}_{y} = \frac{1}{\tilde{l}} \left[\sum_{i=1}^{\tilde{l}} k(\hat{\boldsymbol{x}}_{1}, \boldsymbol{y}_{i}), \ \cdots, \sum_{i=1}^{\tilde{l}} k(\hat{\boldsymbol{x}}_{m}, \boldsymbol{y}_{i}) \right]^{T}$$
$$\equiv \tilde{\boldsymbol{c}}_{y}$$
(13)

$$\begin{split} \boldsymbol{\Phi}_{m}^{T} \phi(\boldsymbol{y}_{j}) &= \left[k(\hat{\boldsymbol{x}}_{1}, \boldsymbol{y}_{j}), \cdots, k(\hat{\boldsymbol{x}}_{m}, \boldsymbol{y}_{j}) \right]^{T} \\ &\equiv \boldsymbol{K}_{m}(\boldsymbol{y}_{j}) \end{split} \tag{14}$$

$$\boldsymbol{\beta}_i \equiv \boldsymbol{H}_{mm}^{-1} \boldsymbol{K}_m(\boldsymbol{y}_i). \tag{15}$$

When the updated accumulation ratio in Eq. (9) is lower than a threshold θ , the useful data should be included in the set of linearly independent data Φ_m . In order to construct a minimum dimensional eigen-feature space, the useful data are selected based on the highest variance $\sigma^2(\phi(\boldsymbol{y}_i))$ from \boldsymbol{y} . Then, $\phi(y_i)$ with the highest $\sigma^2(\phi(\boldsymbol{y}_i))$ is redefined as the (m+1)th independent data $\phi(\hat{\boldsymbol{x}}_{m+1})$ and added to Φ_m as shown below.

$$\boldsymbol{\Phi}_{m+1} = \begin{bmatrix} \boldsymbol{\Phi}_m & \phi(\hat{\boldsymbol{x}}_{m+1}) \end{bmatrix}.$$
(16)

Here, $\sigma^2(\phi(\boldsymbol{y}_i))$ is calculated as follows:

$$\sigma^{2}(\phi(\boldsymbol{y}_{i})) \approx \frac{1}{\|\boldsymbol{h}_{i}\|^{2}(N+\tilde{l})} \sum_{j=1}^{\tilde{l}} \left\{ k(\boldsymbol{y}_{i}, \boldsymbol{y}_{j}) - \frac{N}{N+\tilde{l}} \beta_{i}^{T} \tilde{c} - \frac{1}{N+\tilde{l}} \sum_{k=1}^{\tilde{l}} k(\boldsymbol{y}_{i}, \boldsymbol{y}_{k}) - \sum_{k=1}^{d} \left(\boldsymbol{K}_{m}^{T}(\boldsymbol{y}_{i}) \boldsymbol{\alpha}_{k} \right) \boldsymbol{\alpha}_{k}^{T} \times \left\{ \boldsymbol{K}_{m}(\boldsymbol{y}_{j}) - \frac{N\tilde{\boldsymbol{c}} + \tilde{l}\tilde{\boldsymbol{c}}_{y}}{N+\tilde{l}} \right\} \right\}^{2}$$
(17)

and \tilde{c} , \tilde{c}_y , and $K_m(y_i)$ are updated by

$$\tilde{\boldsymbol{c}} = \boldsymbol{\Phi}_{m+1}^T \boldsymbol{c} = \begin{bmatrix} \boldsymbol{\Phi}_m^T \boldsymbol{c} \\ \beta \left(\hat{\boldsymbol{x}}_{m+1} \right)^T \boldsymbol{\Phi}_m^T \boldsymbol{c} \end{bmatrix}$$
(18)

$$\tilde{\boldsymbol{c}}_{y} = \boldsymbol{\Phi}_{m+1}^{T} \boldsymbol{c}_{y} = \begin{bmatrix} \boldsymbol{\Phi}_{m}^{T} \boldsymbol{c}_{y} \\ \frac{1}{\tilde{l}} \sum_{i=1}^{\tilde{l}} k(\hat{\boldsymbol{x}}_{m+1}, \boldsymbol{y}_{i}) \end{bmatrix}$$
(19)

$$\boldsymbol{K}_{m+1}(\boldsymbol{y}_i) = \begin{bmatrix} \boldsymbol{K}(\boldsymbol{y}_i) \\ k(\hat{\boldsymbol{x}}_{m+1}, \boldsymbol{y}_i) \end{bmatrix}.$$
 (20)

Besides, the coefficient matrix A_d is updated by

$$\mathbf{A}_{d+1} = \begin{bmatrix} \mathbf{A}_{d} & -\frac{1}{||\mathbf{h}_{i}||} \sum_{j=1}^{d} \left(\mathbf{K}_{m}^{T}(\mathbf{y}_{i}) \boldsymbol{\alpha}_{j} \right) \boldsymbol{\alpha}_{j} \\ \mathbf{0} & \frac{1}{||\mathbf{h}_{i}||} \end{bmatrix}.$$
(21)

This procedure is repeated until the accumulation ratio becomes higher than a threshold θ . For the notational simplicity, the selected data are denoted as $\tilde{Y} = \{y_i\}_{i=1}^{\tilde{L}}$. The data are arranged into smaller chunks if the size of useful data \tilde{Y} are large. This is to ensure the constant update speed and memory usages. Then, the necessity of the eigen-feature augmentation is checked after the process of data selection.

B. Feature Space Augmentation

In the initial learning phase, the number of eigen-axes are determined by finding minimum d such that the accumulation ratio in Eq. (9) exceeds a threshold. The eigenvectors $\mathbf{Z} = \{z_1, \ldots, z_d\} \in \Re^{l \times d}$ are represented as the linear combination of m independent data Φ_m

$$\boldsymbol{z}_{i} = [\phi(\hat{\boldsymbol{x}}_{1}), \cdots, \phi(\hat{\boldsymbol{x}}_{m})] \begin{bmatrix} \alpha_{1i} \\ \vdots \\ \alpha_{mi} \end{bmatrix} = \boldsymbol{\Phi}_{m} \boldsymbol{\alpha}_{i}.$$
(22)

Using the selected data \tilde{Y} , the accumulation ratio in Eq. (8) is updated based on Eqs. (9)-(15). If the selected data includes almost all energy in the current eigen-feature space, the dimensionality does not need to be expanded. However, if the eigen-feature space includes certain energy in the complementary eigen-feature space, the dimensional augmentation is needed, or essential information would be lost. Therefore, if the accumulation ratio is lower than θ , it means that selected data may contain some linearly independent data. Then, linearly independent data should be searched and added to the set of linearly independent data Φ_m .

C. Selection of Eigen-axis

Assume that the updated accumulation ratio in Eq. (9) is lower than a θ , it means that linearly independent should be obtained from $\tilde{\mathbf{Y}}$. The data with the highest variance $\sigma^2(\phi(\mathbf{y}_i))$ is found from the selected data. The variance of data projected to the eigenvector is calculated by using Eq. (17) and $\tilde{\mathbf{c}}$, $\tilde{\mathbf{c}}_y$, and $\mathbf{K}_m(\mathbf{y}_i)$ are updated by Eqs. (18)-(20). Here, $\phi(\mathbf{y}_i)$ is defined as the m + 1th linearly independent data and added to a set of Φ_m as shown in Eq. (16). Besides, the coefficient matrix is updated by using Eq. (21). Here, we rewrite m + 1 as m' and d + 1 as d' for the sake of the notational simplicity. By adding $\sigma^2(\phi(\mathbf{y}_i))$ to the numerator of Eq. (9), the accumulation ratio C'(d) is updated. If the updated accumulation ratio C'(d) is still smaller than a threshold θ , another data should be selected. Assume that

)

we select (m' + 1)-th data. A candidate axis is obtained as follows:

$$\hat{\boldsymbol{z}}_{i}^{\prime} = \frac{\boldsymbol{h}_{i}^{\prime}}{||\boldsymbol{h}_{i}^{\prime}||} \tag{23}$$

where

In Eq. (23), $||h'_i||$ can be obtained by calculating its squared norm as follows:

$$||\boldsymbol{h}'_{i}||^{2} = k(\boldsymbol{y}_{i}, \boldsymbol{y}_{i}) - 2 \sum_{j=1}^{d'} \left(\boldsymbol{K}_{m'}^{T}(\boldsymbol{y}_{i})\boldsymbol{\alpha}_{j}\right)^{2} + \sum_{j=1}^{d'} \sum_{k=1}^{d'} \left(\boldsymbol{K}_{m'}^{T}(\boldsymbol{y}_{i})\boldsymbol{\alpha}_{k}\right) \left(\boldsymbol{K}_{m'}^{T}(\boldsymbol{y}_{i})\boldsymbol{\alpha}_{k}\right) \times \boldsymbol{\alpha}_{j}^{T} \boldsymbol{H}_{(m')(m')} \boldsymbol{\alpha}_{k}.$$
(25)

The variance $\sigma^2(\phi(\boldsymbol{y}_i))$ can be calculated as follows:

$$\sigma'^{2}(\phi(\boldsymbol{y}_{i})) = \frac{1}{||\boldsymbol{h}'_{i}||^{2}(N+\tilde{L})} + \sum_{j=1}^{\tilde{L}} \left\{ k(\boldsymbol{y}_{i}, \boldsymbol{y}_{j}) - \frac{N}{N+\tilde{L}} \boldsymbol{\beta}_{i}^{T} \tilde{c} - \frac{1}{N+\tilde{L}} \sum_{k=1}^{\tilde{L}} k(\boldsymbol{y}_{i}, \boldsymbol{y}_{k}) - \sum_{k=1}^{d'} \left(\boldsymbol{K}_{m'}^{T}(\boldsymbol{y}_{i}) \boldsymbol{\alpha}_{k} \right) \boldsymbol{\alpha}_{k}^{T} \boldsymbol{K}_{m'}(\boldsymbol{y}_{j}) + \frac{N}{N+\tilde{L}} \sum_{k=1}^{d'} \left(\boldsymbol{K}_{m'}^{T}(\boldsymbol{y}_{i}) \boldsymbol{\alpha}_{k} \right) \boldsymbol{\alpha}_{k}^{T} \tilde{c} + \frac{\tilde{L}}{N+\tilde{L}} \sum_{k=1}^{d'} \left(\boldsymbol{K}_{m'}^{T}(\boldsymbol{y}_{i}) \boldsymbol{\alpha}_{k} \right) \times \boldsymbol{\alpha}_{k}^{T} \tilde{c}_{j} \right\}^{2}$$
(26)

where \tilde{c} and \tilde{c}_y can be calculated recursively as follows:

$$\tilde{\boldsymbol{c}} = \begin{bmatrix} \tilde{\boldsymbol{c}} \\ \left(\boldsymbol{H}_{mm}^{-1} \boldsymbol{K}_{m}(\hat{\boldsymbol{y}}) \right)^{T} \tilde{\boldsymbol{c}} \end{bmatrix}$$
(27)

$$\tilde{\boldsymbol{c}}_{y} = \begin{bmatrix} \tilde{\boldsymbol{c}}_{y} \\ \frac{1}{\tilde{L}} \sum_{i=1}^{\tilde{L}} k(\hat{\boldsymbol{x}}_{m'}, \boldsymbol{y}_{i}) \end{bmatrix}$$
(28)

$$\boldsymbol{K}_{m'}(y_i) = \begin{bmatrix} \boldsymbol{K}_m(y_{(i)}) \\ k(\hat{\boldsymbol{x}}_{m'}, y_{(i)}) \end{bmatrix}.$$
 (29)

Here, $\phi(\boldsymbol{y}_i)$ with the highest $\sigma'^2(\phi(\hat{\boldsymbol{x}}_{m'+1}))$ is selected as $\phi(\hat{\boldsymbol{x}}_{m'+1})$ and added to $\Phi_{m'+1}$ as follows:

$$\boldsymbol{\Phi}_{m'+1} = \begin{bmatrix} \boldsymbol{\Phi}_{m'} & \phi(\hat{\boldsymbol{x}}_{m'+1}) \end{bmatrix}$$
(30)

where $\phi(\hat{x}_{m'+1})$ is a new linearly independent data. By adding $\sigma'^2(\phi(\hat{x}_{m'+1}))$ to the numerator of C(d'), the new

accumulation ratio C(d' + 1) is obtained. The coefficient matrix is updated as follows:

$$\boldsymbol{A}_{d'+1} = \begin{bmatrix} \boldsymbol{A}_{d'} & -\frac{1}{||\boldsymbol{h}_i||} \sum_{j=1}^{d'} \left(\boldsymbol{K}_{m'}^T(\boldsymbol{y}_i) \boldsymbol{\alpha}_j \right) \boldsymbol{\alpha}_j \\ \boldsymbol{0} & \frac{1}{||\boldsymbol{h}_i||} \end{bmatrix}.$$
(31)

This procedure is repeated until the accumulation ratio becomes higher than a threshold θ .

D. Update of Feature Space

The update of feature space depends on two cases. The first case is when the new eigenvectors are added and the second case is when the eigenvector is not added into the feature space. When l new eigenvectors are added, the eigenvector matrix is updated by

$$\boldsymbol{Z}'_{d+l} = \begin{bmatrix} \boldsymbol{Z}_d & \hat{\boldsymbol{Z}}_l \end{bmatrix} \boldsymbol{R}$$
(32)

and \hat{Z}_l can be represented as follows:

$$\hat{\boldsymbol{Z}}_{l} = \begin{bmatrix} \boldsymbol{z}_{d+1} & \cdots & \boldsymbol{z}_{d+l} \end{bmatrix}$$
(33)

where $\mathbf{R} \in R^{(d+l) \times (d+l)}$ is a rotation matrix.

In order to obtain the rotation matrix \mathbf{R} , we solve the following intermediate kernel eigenvalue problem.

$$\frac{1}{N+\tilde{L}} \left(N \begin{bmatrix} \mathbf{\Lambda} & \mathbf{0}_{d \times l} \\ \mathbf{0}_{d \times l} & \mathbf{0}_{l \times l} \end{bmatrix} + \sum_{i=1}^{\tilde{L}} \boldsymbol{f}_{i} \boldsymbol{f}_{i}^{T} + \frac{N\tilde{L}}{N+\tilde{L}} \boldsymbol{p} \boldsymbol{p}^{T} \right) \boldsymbol{R} = \boldsymbol{R} \boldsymbol{\Lambda}_{d+l}^{\prime}$$
(34)

where both f and p are given by

$$\boldsymbol{f}_{i} = \begin{bmatrix} \boldsymbol{\alpha}_{1}^{T} \\ \vdots \\ \boldsymbol{\alpha}_{d+l}^{T} \end{bmatrix} \begin{pmatrix} \boldsymbol{K}_{m+l}(\boldsymbol{y}_{i}) - \tilde{\boldsymbol{c}}_{y}^{(m+l)} \end{pmatrix}$$
(35)

$$\boldsymbol{p} = \begin{bmatrix} \boldsymbol{\alpha}_{1}^{T} \\ \vdots \\ \boldsymbol{\alpha}_{d+l}^{T} \end{bmatrix} \left(\tilde{\boldsymbol{c}}^{(m+l)} - \tilde{\boldsymbol{c}}_{y}^{(m+l)} \right).$$
(36)

On the other hand, when no l new eigenvector is added, the updated eigenvectors are given as follows:

$$\boldsymbol{Z}'_{d} = \boldsymbol{Z}_{d}\boldsymbol{R} \tag{37}$$

and the intermediate eigenvalue is given by

$$\frac{1}{N+\tilde{L}}\left(N\boldsymbol{\Lambda}_{d}+\sum_{i=1}^{\tilde{L}}\boldsymbol{f}_{i}\boldsymbol{f}_{i}^{T}+\frac{N\tilde{L}}{N+\tilde{L}}\boldsymbol{p}\boldsymbol{p}^{T}\right)\boldsymbol{R}=\boldsymbol{R}\boldsymbol{\Lambda}'_{d}.$$
(38)

Algorithm 1 Proposed Method

- **Require:** Initial training data $\mathbf{X} = \{x_i\}_{i=1}^N$. 1: Obtain a threshold of accumulation ratio θ based on the cross validation.
- 2: Perform KPCA for X and obtain an initial eigen-feature space $\Omega =$ $\{ \boldsymbol{\Phi}_m, \boldsymbol{A}_d, \boldsymbol{\Lambda}_d, N \}.$
- 3: Obtain the linearly independent data such that the accumulation ratio C(d) is higher than a θ .
- 4: loop
- 5: **Input:** A new chunk of training data $Y = \{y_i\}_{i=1}^{L}$.
- Subdivide Y into smaller chunks $y = \{y_i\}_{i=1}^l$ 6:
- Perform the data selection based on the accumulation ratio and obtain 7: the selected data $\tilde{\boldsymbol{Y}} = \{\boldsymbol{y}_i\}_{i=1}^L$.
- Update the C(d) by Eq. (9). 8:
- if $C(d) \geq \theta$ then Q٠ 10:
- Obtain \mathbf{R} by Eq. (38). Update the coefficient matrix: $A'_d = A_d R$. 11:
- 12: else
- 13: while $C(d) < \theta$ do
- 14: for i=1: \tilde{L} do
- 15: Calculate variance $\sigma^2(\phi(\boldsymbol{y}_i))$ by Eq. (17).
- 16: end for 17: Obtain linearly independent data with the m + 1 largest variances, $\sigma^2(\phi(\boldsymbol{y}_i))$.
- 18: Update the coefficient matrix A_{d+1} by Eq. (21).
- Calculate the accumulation ratio C(d + 1) by adding 19: $\sigma^2(\phi(\boldsymbol{y}_i))$ to the numerator of Eq. (9).
- 20: end while
- 21: Obtain *R* by solving Eq. (34).
- 22: Update the coefficient Matrix: $A'_{d+l} = A_{d+l}R$. 23: end if
- Update $\|\boldsymbol{c}\|^2$ and $\tilde{\boldsymbol{c}}'$ by Eqs.(39)-(40). 24.
- 25: end loop

E. Mean Update

 $\|c^2\|$ are updated by

$$\|\boldsymbol{c}'\|^{2} = \frac{1}{(N+\tilde{L})^{2}} \left(N^{2} \|\boldsymbol{c}\|^{2} + 2N \sum_{i=1}^{L} \boldsymbol{\beta}_{i}^{T} \tilde{\boldsymbol{c}} + \sum_{i=1}^{\tilde{L}} \sum_{j=1}^{\tilde{L}} k(\boldsymbol{y}_{i}, \boldsymbol{y}_{j}) \right)$$
(39)

$$\tilde{\boldsymbol{c}}' = \frac{1}{N + \tilde{L}} (N \tilde{\boldsymbol{c}} + \tilde{L} \tilde{\boldsymbol{c}}_y).$$
(40)

When the new selected data are given, $\|c'\|^2$ and \tilde{c}' are replaced with $\|\boldsymbol{c}\|^2$ and $\tilde{\boldsymbol{c}}$, respectively.

The brief learning procedures are shown in Algorithm 1.

IV. EXPERIMENTAL RESULTS

A. Experiment Setup

In order to evaluate the effectiveness of our proposed method, we conduct several experiments based on the five data sets which are selected from the UCI Machine Learning Repository [11]. The evaluated data sets are listed in Table I.

For Poker Hand data, 2000 data are randomly selected from the original 1025010 data as train data and another 1000 data are selected as the test data. On the other hand, for Optical-digit data, 1500 data are selected from 3823 data as the train data and all 1797 test data are used to measure the performance. The initial data consists of 100 training data and remaining data are used for the incremental learning.

In our experiments, the results are averaged over 30 trials because incremental learning depends on the sequence of the training data. The kernel function used in the experiments is Gaussian kernel $k(\boldsymbol{x}, \boldsymbol{y}) = \exp(-\gamma || \boldsymbol{x} - \boldsymbol{y} ||^2)$. The parameter γ and θ are selected based on the 5-fold cross validation in the initial learning.

The computer used here has the following specifications:

- 1) Intel(R) Core(TM)2 Duo (3.16 GHz) CPU.
- 2) 2GB Random acess memory (RAM).
- 3) Windows 7 (32 bit) OS.

All the programs are develop with Matlab (R2007b).

B. Performance Evaluation of the CIKPCA

The objective of this experiment is to measure the performance of the proposed method. In this experiment, the recognition accuracy, the computational time, and the eigenspace dimension of the traditional KPCA, IKPCA and CIKPCA are evaluated based on the five UCI Machine Learning Repository data sets.

Table II shows the comparison of the recognition accuracy [%] with the standard deviation, the computational time (sec.) with the standard deviation and the dimension of eigenspace with the standard deviation for the traditional KPCA, IKPCA and CIKPCA. The chunk size is set at L = 100.

Table II (a) shows the average and the standard deviation of the recognition accuracy that evaluated over 30 trials. The recognition accuracy of the proposed method is compared with the conventional KPCA and Takeuchi et al.'s IKPCA. From the results, we observed that the recognition accuracy based on the five data sets are almost similar for all the models. The recognition accuracy are almost similar with the other two models although the training data in CIKPCA are reduced. Therefore, it suggests that essential information is not lost actually.

Table II (b) shows the computation time for KPCA, IKPCA and CIKPCA. It is observed that the computational time of CIKPCA is obviously faster than KPCA and IKPCA for all the data sets. In KPCA, the kernel matrix of all training data is computed, leading to the large memory and computation time. The learning of IKPCA is not affected by the chunk sizes. Even though the data are received in a chunk, IKPCA would learn the data one by one and the eigenvalue decomposition is solved for every incoming data. On the other hand, the computational time of the kernel matrix decomposition of the selected data is reduced by reducing the amount of data in a chunk. Furthermore, the eigenvalue decomposition is only solved once when a chunk

TABLE I

EVALUATED DATA SETS.

Data Sets	# Attrib.	# Classes	# Train	# Test
Ozone	72	2	1000	848
Thyroid	21	3	3772	3428
Poker Hand	10	10	2000	1000
Optical-digit	64	10	1500	1797
Adult	14	2	22611	22611

PERFORMANCE COMPARISON OF (A) RECOGNITION ACCURACY [%], (B) COMPUTATIONAL TIME (SEC.) AND (C) EIGENSPACE DIMENSION.

			((a)							
	Data S	Sets	KPCA		IKPCA	CIK	PCA				
	Ozo	Ozone		95.0±0.01 9		95.1:	±0.01				
	Thyre	Thyroid		9	1.2 ± 0.02	91.4:	$91.4 {\pm} 0.02$				
	Poker I	Poker Hand		4	4.6 ± 0.02	44.7:	44.7±0.02				
	Optical	-digit	88.5±0.01	8	38.5 ± 0.01 88		88.1 ± 0.01				
	Adu	lt	73.0±0.26	7	2.4 ± 0.25	72.6:	72.6 ± 0.25				
(b)											
D	ata Sets	KPCA			IKPCA		CIKPCA				
	Ozone	117.05±120.60)	2.91±7.79		0.53±0.17				
Т	Thyroid	792.44±1.88			2.13±1.63		$0.86 {\pm} 0.46$				
Pol	ker Hand	133.30±17.43			1.18 ± 1.00		0.46 ± 0.35				
Op	tical-digit	1293.52±677.33		3	70.21±241.64		1.53 ± 0.50				
	Adult	$186,418.4\pm 64,463$		3.0	95.6±61.8		8 38.8±23.8				
(c)											
[Data Se	ts	KPCA		IKPCA		CIKPCA				
[Ozone	;	10.5±12.5		17.9±30.9		8.7±6.4				
	Thyroi	Thyroid 11.2		15.0 ± 9.1		16.3±3.8		1			
	Poker Hand		$8.8 {\pm} 0.9$	8.9±2.1		8.5±1.1					
ĺ	Optical-digit		38.4±24.1	111.77±126.5		33.9±10.7					
[Adult		11.1±4.1	15.9 ± 9.1		13	13.2±5.3				

of selected data is given. Therefore, the computation time for CIKPCA is reduced.

To understand the results of the computational time in Table II (b), let us study the computational complexity of IKPCA and CIKPCA. The highest computation time in IKPCA comes from the calculation of $||h||^2$. Its computational complexity is $O(m^2d^2)$ where m and d are the number of linearly independent data and the dimensions of an eigenfeature space, respectively.

On the other hand, the computation in CIKPCA is dominated by the calculation of the denominator of the accumulation ratio shown in Eq.(11) or $||h||^2$ in Eq. (25). Its computation complexity is $O(m^3L + nL^2)$ or $O(m^2d^2L)$ where *n* and *L* are the number of input space and the size of data chunk, respectively. To calculate the denominator, we need to obtain $\beta^T(\Phi_m^T c)$ where $\beta^T = H_{mm}^{-1} K_m(x)$ (see Eq. (15)). Here, the computational of H_{mm}^{-1} is $O(m^3)$. In addition, the computation complexity to calculate the third term of the right hand side in Eq.(11) is $O(nL^2)$. Therefore, the computational time is dominated in the calculation of the third term in Eq.(11) when large chunk of *L* data is received. Otherwise, the computational complexity is dominated in the calculation of $||h||^2$ when the moderate size of chunk data is received.

However, in our proposed method, a chunk of data is divided into small chunks and only some training data are selected when a large chunk of data is received. Therefore, the size of data chunk is usually small (\tilde{l} for small chunks and \tilde{L} for selected data in a chunk). In this case, our proposed method would not encounter the problem of large data chunk. The computational complexity of Eq.(11) and $||h||^2$ is reduced to $O(m^3 + n)$ and $O(m^2d^2)$, respectively.

In addition, Table II (c) shows the eigen-space dimension-

ality for the three models. It is observed that the dimension of the eigenspace for CIKPCA are smaller than KPCA for almost all the data sets. Thus, it suggests that although the dimension of the eigen-feature space is small, but it contains almost the same energy as KPCA for the classification. The computational time are reduced with only some useful features.

It is concluded that the computational time of CIKPCA are obviously faster than KPCA and IKPCA while the recognition accuracies are almost similar for all the models. Therefore CIKPCA is more efficient than KPCA and IKPCA.

C. Efficiency of IKPCA with Different Number of Chunk Sizes

As mentioned before, Takaomi et al. [10] found that more time is needed for the large chunk data. Therefore, the objective of this experiment is to measure the influence of the data division and the data selection when a chunk of data is received. The computational complexity of the proposed method is often dominated in the calculation of $||h||^2$. Its computational complexity is $O(m^2d^2)$. Thus, in order to measure the important of the data division and the data selection, the computation time of CIKPCA and IKPCA are compared with regards to the different number of chunk sizes. Fig. 1 shows the influence of the chunk sizes on the learning time for the three data sets, Ozone, Thyroid, and Poker Hand data. Since the data in a chunk have been selected, it is expected that the eigen-feature space is constructed with only some useful data. Therefore, the learning time are expected to be shortened. In this experiment, we evaluate the learning time of IKPCA and CIKPCA for the chunks sizes at 100, 300, 500, 800 and 1000.



Fig. 1. Learning time (sec.) vs. different number of chunk sizes for (a) Ozone, (b) Thyroid and (c) Poker Hand data.

As seen in Fig. 1, we observed that the computation time for CIKPCA are significantly shorter than IKPCA with regards to the different number of chunk sizes for the three data sets. In CIKPCA, the data are first selected when a chunk of data is given; therefore, the training data in a chunk are reduced. The kernel matrix is computed based on the smaller size of data leading to the significant reduction in the computation time. Furthermore, the computation time are decreased as the chunk size becomes large. This is due to the reduction of the total number of eigenvalue decomposition when the chunk size becomes large. For IKPCA, the training data is learned one by one although the training data are given in a chunk. Hence, the learning time is not affected by the chunk sizes.

Since the computation time are reduced significantly in the proposed method, we further evaluate the efficiency of CIKPCA by evaluating the influence of chunk sizes on the recognition accuracy.

Fig. 2 shows the recognition accuracy of IKPCA and CIKPCA with regards to the different number of chunk sizes for the three data sets. As seen in Fig. 2, the recognition accuracies are almost similar between IKPCA and CIKPCA for the three data sets. The learning accuracies of the proposed method are almost the same with IKPCA with regards to the different number of chunk sizes although the amount of training data are reduced. It suggests that the essential information is not lost and the selected data contains some

useful information for the eigen-feature augmentation.

Therefore, it is concluded that the proposed method is more efficient than Takeuchi et al.'s IKPCA where the computation time are reduced significantly without reducing the recognition accuracy.

V. CONCLUSIONS

As a conclusion, fast CIKPCA is proposed. To reduce the computational time of the proposed method, the data given in a chunk are first divided into small chunks and data are selected from the smaller chunks based on the accumulation ratio. The reason of data selection is to retain only some useful data for the eigen-feature space learning. In addition, the computation time are reduced since the computation time of kernel matrix decomposition are reduced.

Then, compact feature space is constructed by the linear combination of the linearly independent data that are selected from the reduced data and the eigenvalue problem is solved only once for the selected data.

In this paper, CIKPCA is compared with KPCA and Takeuchi et al.'s IKPCA. Experimental results demonstrated that CIKPCA requires shorter learning time than IKPCA and KPCA with similar recognition accuracy for all the models. To further evaluate the effectiveness of CIKPCA, the computational time of CIKPCA is compared with IKPCA with regards to the different number of chunk sizes. The results of all the data sets revealed that CIKPCA requires



Fig. 2. Recognition Accuracy (%.) vs. different number of chunk sizes for (a) Ozone, (b) Thyroid and (c) Poker Hand data.

shorter time compared to IKPCA for all the chunk sizes. The learning time are reduced when the chunk sizes becomes large. In addition, the total number of eigenvalue decomposition are reduced when the chunk size becomes large. In general, the proposed method could accelerate the learning time without reducing the recognition accuracy. Therefore, we conclude that division of data into smaller chunks and selection of training data is important in order to accelerate the learning time. In the future, we will compare the proposed method with Chunk Incremental Principal Component Analysis (CIPCA).

REFERENCES

- [1] I. Jolliffe, Principal component analysis, Springer, New York, 1986.
- [2] F. Han, and H. Liu, "High Dimensional Semiparametric Scale-invariant Principal Component Analysis," *Advances in Neural Information Processing Systems (NIPS)*, pp. 171-179, 2012.
- [3] P. Honeine, "Online Kernel Principal Component Analysis: A reduced-Order Model," *IEEE Trans. on Pattern Analysis and Machine Intelli*gence, vol. 34, no. 9, pp. 1814-1826, 2012.
- [4] H. Zhao, P. Chi, and J. T. Kwok, "A Novel Incremental Principal Component Analysis and Its Application for Face Recognition," *IEEE Trans. on Systems, Man and Cybernetics, Part B*, vol. 36, no. 4, pp. 873-886, 2006.
- [5] T.-J. Chin, and D. Suter, "Incremental Kernel Principal Component Analysis," *IEEE Trans. on Image Processing*, vol. 16, no. 6, pp. 1662-1674, 2007.
- [6] S. Ozawa, S. Pang, and N. Kasabov, "Incremental Learning of Chunk Data for On-line Pattern Classification Systems," *IEEE Trans. on Neural Networks*, vol. 19, no. 6, pp. 1061-1074, 2008.

- [7] S. Ozawa, S. L. Toh, S. Abe, S. Pang, and N. Kasabov, "Incremental Learning of Feature Space and Classifier for Face Recognition," *Neural Networks*, vol. 18, nos. 5-6, pp. 575-578, 2005.
- [8] B. Schölkopf, A. Smola, K. Muller, "Kernel Principal Component Analysis," In:W. Gerstner, M. Hasler, A. Germond, J.-D. Nicoud (eds.) ICANN 1997. LNCS, Springer-Heideberg, no. 1327, pp. 583-588, 1997.
- [9] Y. Takeuchi, S. Ozawa, S. Abe, "An Efficient Incremental Kernel Principal Component Analysis for Online Feature Extraction," *Proc.* of Int. Joint Conf. Neural Networks 2007, pp. 2346-2351, 2007.
- [10] T. Tokumoto and S. Ozawa, "Property of Learning Chunk Data Using Incremental Kernel Principal Component Analysis," *Proc. IEEE Workshop on Evolving and Adaptive Intelligent Systems (EAIS2012, Madrid)*, pp. 7-10, 2012.
- [11] S. Asunction, D. J. Newman, "UCI Machine Learning Repository," UC. Irvine, School of Info. and Comp. Sci., 2007.
- [12] S. Abe, Support Vector Machines for Pattern Classification, Second Edition, Springer-Verlag, London, 2010.