

# Improving KPCA Online Extraction by Orthonormalization in the Feature Space

João B. O. Souza Filho and Paulo S. R. Diniz, *Fellow, IEEE*

**Abstract**—Recently, some online kernel principal component analysis (KPCA) techniques based on the generalized Hebbian algorithm (GHA) were proposed for use in large data sets, defining kernel components using concise dictionaries automatically extracted from data. This brief proposes two new online KPCA extraction algorithms, exploiting orthogonalized versions of the GHA rule. In both the cases, the orthogonalization of kernel components is achieved by the inclusion of some low complexity additional steps to the kernel Hebbian algorithm, thus not substantially affecting the computational cost of the algorithm. Results show improved convergence speed and accuracy of components extracted by the proposed methods, as compared with the state-of-the-art online KPCA extraction algorithms.

**Index Terms**—Kernel methods, machine learning and generalized Hebbian algorithm (GHA), online kernel algorithms, principal component analysis (PCA).

## I. INTRODUCTION

Among the several kernel methods developed in the last decades, a remarkable technique is the Kernel PCA (KPCA), a powerful nonlinear generalization of the well-known PCA technique [3]. The KPCA is a useful tool for several applications, notably for image denoising [2] and novelty detection [4].

The KPCA was originally presented as a Gram matrix eigendecomposition problem [5], thus solvable by standard computational linear algebra techniques [6]. However, the adoption of this method on large-scale data sets, such as those involved in image processing applications, may result in high dimensional Gram matrices, turning the extraction process computationally infeasible.

Motivated by this fact, several iterative extraction algorithms were proposed [2], [7], [8]. One example is the kernel Hebbian algorithm (KHA) [2], which consists in a kernelization of the generalized Hebbian algorithm (GHA) [1], a widely used online PCA extraction algorithm. Although KHA does not involve Gram matrix estimation, it requires that all extraction data sets have to be known in advance.

Online KPCA extraction was only effectively achieved by some recently proposed algorithms [9]–[11], exploiting techniques, such as the approximated linear dependence (ALD) criterion [12] and coherence [13], to define concise dictionaries used in the expansion of kernel components. The online KHA (OKHA) [9] and subset KHA (SubKHA) [10] are useful online extensions of the KHA [2], essentially differing in how the dictionaries are built and updated. Moreover, as shown by Tanaka *et al.* [11], the extraction of kernel components using a subset of training samples can be stated as a generalized eigenvalue problem [11]. In his work, Tanaka also proposed a recursive least-square algorithm for KPCA extraction based on the kernelization of the iterative weighted rule technique [14].

The adoption of additional steps for the normalization or orthonormalization of component estimates has been shown to be effective for

some principal subspace extraction methods, such as Oja [15], [16] and PAST [17], as well as for kernel subspace component learning [18], improving the convergence behavior of such algorithms. In the case of GHA, the lack of orthogonality between component estimates often compromises its accuracy and convergence speed.

Inspired by these ideas, this brief proposes two new KPCA extraction algorithms. In both, we have included an additional orthogonalization or orthonormalization step implicitly applied to component estimates in the feature space after each KHA iteration. The proposed methods exploit the classical Gram–Schmidt algorithm [6] and possess low computational cost. Results show a significant improvement in both the convergence speed and accuracy of the extracted components. The supplementary material [19] includes the derivations of some paper equations.

## II. ORTHOGONAL GHA

A well-known rule for online principal component extraction is the GHA. Assuming the extraction of  $p$  components from a sequence of random vectors denoted by  $\mathbf{x}(k)$ , this algorithm can be described by the following set of equations:

$$y_j(k) = \mathbf{w}_j^T(k-1)\mathbf{x}(k) \quad (1)$$

$$\mathbf{w}_j(k) = \mathbf{w}_j(k-1) + \eta y_j(k) \mathbf{e}_j(k), \quad 1 \leq j \leq p \quad (2)$$

with

$$\mathbf{e}_j(k) = \mathbf{x}(k) - y_j(k)\mathbf{w}_j(k-1) - \mathbf{d}_j(k) \quad (3)$$

$$\mathbf{d}_j(k) = \begin{cases} 0, & j = 1 \\ \sum_{i=1}^{j-1} y_i(k)\mathbf{w}_i(k-1), & 2 \leq j \leq p \end{cases} \quad (4)$$

where the parameter  $\eta$  is the learning rate factor, which directly affects the convergence behavior of this algorithm.

Another relevant issue in GHA convergence is the orthonormality between component estimates. Even when assuming that two arbitrary component estimates  $\mathbf{w}_i(k-1)$  and  $\mathbf{w}_j(k-1)$  are orthonormal at the iteration  $(k-1)$ , the inner product between them after a GHA update is given by [19, Sec. I]

$$\mathbf{w}_i^T(k)\mathbf{w}_j(k) \approx \eta y_i(k)y_j(k), \quad 1 \leq i \neq j \leq p. \quad (5)$$

However, the convergence speed and accuracy of GHA may be improved if we periodically orthogonalize component estimates using the Gram–Schmidt orthogonalization (GSO) [20], [21], or the householder transformation [22], or employing the givens rotations [23]. The classical GSO (cGSO) is an attractive technique for this task, since it can be easily extended to work in the feature space, and has sufficient accuracy for the scope of applications commonly solved by the GHA. Besides, as we will show, the cGSO procedure can be incorporated to GHA equations with a mild increase in computational cost.

For this, initially consider the Amari rule [24] for the extraction of the first principal component as follows:

$$y_1(k) = \bar{\mathbf{w}}_1^T(k-1)\mathbf{x}(k) \quad (6)$$

$$\mathbf{w}_1(k) = \bar{\mathbf{w}}_1(k-1) + \eta y_1(k)\mathbf{x}(k) \quad (7)$$

$$\bar{\mathbf{w}}_1(k) = \frac{\mathbf{w}_1(k)}{\|\mathbf{w}_1(k)\|_2}. \quad (8)$$

Manuscript received December 23, 2015; revised October 28, 2016; accepted January 23, 2017.

The authors are with the Electrical Engineering Program, COPPE/POLI/Federal University of Rio de Janeiro, Rio de Janeiro 21941-972, Brazil (e-mail: jbfilho@poli.ufrj.br; diniz@smt.ufrj.br).

This paper has supplementary downloadable material available at <http://ieeexplore.ieee.org>, provided by the author.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNNLS.2017.2660441

As shown by Oja [25], if we assume the value of  $\eta$  small, the normalization step can be incorporated into (7), resulting in the following equations:

$$y_1(k) = \mathbf{w}_1^T(k-1)\mathbf{x}(k) \quad (9)$$

$$\mathbf{w}_1(k) = \mathbf{w}_1(k-1) + \eta y_1(k)[\mathbf{x}(k) - y_1(k)\mathbf{w}_1(k-1)]. \quad (10)$$

Note that a simple extension of Oja algorithm for the extraction of  $p$  ( $1 \leq j \leq p$ ) components is given by

$$y_j(k) = \bar{\mathbf{w}}_j^T(k-1)\mathbf{x}(k) \quad (11)$$

$$\mathbf{w}_j(k) = \bar{\mathbf{w}}_j(k-1) + \eta y_j(k)\mathbf{f}_j(k) \quad (12)$$

$$\tilde{\mathbf{w}}_j(k) = \mathbf{w}_j(k) - \mathbf{g}_j(k) \quad (13)$$

$$\bar{\mathbf{w}}_j(k) = \frac{\tilde{\mathbf{w}}_j(k)}{\|\tilde{\mathbf{w}}_j(k)\|_2} \quad (14)$$

where

$$\mathbf{f}_j(k) = \mathbf{x}(k) - y_j(k)\bar{\mathbf{w}}_j(k-1) \quad (15)$$

$$\mathbf{g}_j(k) = \begin{cases} 0, & j = 1 \\ \sum_{i=1}^{j-1} c_{ji}(k)\bar{\mathbf{w}}_i(k), & 2 \leq j \leq p \end{cases} \quad (16)$$

and

$$c_{ji}(k) = \mathbf{w}_j^T(k)\bar{\mathbf{w}}_i(k). \quad (17)$$

This extension simply employs multiple Oja rules coupled by a straightforward cGSO orthonormalization procedure applied to component estimates, according to (13), (14), (16), and (17).

Considering (12) and (13),  $j > i$ , as well as assuming component estimates orthonormal at iteration  $(k-1)$ , an approximated value for the projection coefficient defined by (17) is given by [19, Sec. II]

$$c_{ji}(k) \approx 2\eta y_i(k)y_j(k). \quad (18)$$

Assuming  $\bar{\mathbf{w}}_i(k) = \bar{\mathbf{w}}_i(k-1) + \Delta\bar{\mathbf{w}}_i(k)$ , where the vector  $\Delta\bar{\mathbf{w}}_i(k)$  represents an increment to the component estimate  $\bar{\mathbf{w}}_i(k-1)$  proportional to  $\eta$ , and considering (4) and (18), (13) can be rewritten as

$$\begin{aligned} \tilde{\mathbf{w}}_j(k) &\approx \mathbf{w}_j(k) - 2\eta y_j(k) \sum_{i=1}^{j-1} y_i(k)[\bar{\mathbf{w}}_i(k-1) + \Delta\bar{\mathbf{w}}_i(k)] \\ &\approx \mathbf{w}_j(k) - 2\eta y_j(k)\hat{\mathbf{d}}_j(k) + O(\eta^2) \\ &\approx \hat{\mathbf{w}}_j(k) \end{aligned} \quad (19)$$

with

$$\hat{\mathbf{w}}_j(k) = \mathbf{w}_j(k) - 2\eta y_j(k)\hat{\mathbf{d}}_j(k) \quad (20)$$

$$\hat{\mathbf{d}}_j(k) = \begin{cases} 0, & j = 1 \\ \sum_{i=1}^{j-1} \bar{\mathbf{w}}_i(k-1)y_i(k), & 2 \leq j \leq p. \end{cases} \quad (21)$$

Replacing (12) in (20), it results in the following equation:

$$\begin{aligned} \hat{\mathbf{w}}_j(k) &= \bar{\mathbf{w}}_j(k-1) + \eta y_j(k)\mathbf{f}_j(k) - 2\eta y_j(k)\hat{\mathbf{d}}_j(k) \\ &= \bar{\mathbf{w}}_j(k-1) + \eta y_j(k)[\mathbf{x}(k) - y_j(k)\bar{\mathbf{w}}_j(k-1) - \hat{\mathbf{d}}_j(k)] \\ &\quad - \eta y_j(k)\hat{\mathbf{d}}_j(k) \\ &= \bar{\mathbf{w}}_j(k-1) + \eta y_j(k)\hat{\mathbf{e}}_j(k) - \eta y_j(k)\hat{\mathbf{d}}_j(k). \end{aligned} \quad (22)$$

Thus, the proposed algorithm, here named as orthonormal GHA (ON-GHA), has the following equations:

$$y_j(k) = \bar{\mathbf{w}}_j^T(k-1)\mathbf{x}(k) \quad (23)$$

$$\hat{\mathbf{e}}_j(k) = \mathbf{x}(k) - y_j(k)\bar{\mathbf{w}}_j(k-1) - \hat{\mathbf{d}}_j(k) \quad (24)$$

$$\hat{\mathbf{w}}_j(k) = \bar{\mathbf{w}}_j(k-1) + \eta y_j(k)\hat{\mathbf{e}}_j(k) - \eta y_j(k)\hat{\mathbf{d}}_j(k) \quad (25)$$

$$\bar{\mathbf{w}}_j(k) = \frac{\hat{\mathbf{w}}_j(k)}{\|\hat{\mathbf{w}}_j(k)\|_2} \quad (26)$$

where the vector  $\hat{\mathbf{d}}_j(k)$  is given by (21). Note that the proposed ON-GHA has quite similar equations to GHA. Indeed, they essentially differ by the term  $-\eta y_j(k)\hat{\mathbf{d}}_j(k)$  in (25) and by the normalization step imposed by (26). If the normalization step is ignored, the ON-GHA equations become similar to the stochastic gradient ascent (SGA) algorithm proposed by Oja [26]. The orthonormality between component estimates from the ON-GHA corresponds to [19, Sec. III]

$$\bar{\mathbf{w}}_j^T(k)\bar{\mathbf{w}}_i(k) \approx 0, \quad 1 \leq i \neq j \leq p \quad (27)$$

as expected.

### III. ONLINE ORTHONORMAL KERNEL HEBBIAN ALGORITHM

To kernelize the ON-GHA, we will consider a concise and constructive dictionary  $\mathbf{D}_m$  with  $m$  elements, composed of some dynamically selected input samples as follows:

$$\mathbf{D}_m = [\mathbf{x}(a) \quad \dots \quad \mathbf{x}(q)], \quad a \neq q. \quad (28)$$

The mapping of dictionary members into the feature space can be expressed as

$$\Psi_m = [\Phi(\mathbf{x}(a)) \quad \dots \quad \Phi(\mathbf{x}(q))] \quad (29)$$

where the arbitrary function  $\Phi(\mathbf{x}(a))$  denotes the mapping of an input sample  $\mathbf{x}(a)$  into this space. As in the OKHA algorithm, we will explore the ALD criterion to evaluate dictionary augmentation at each algorithm iteration. For this, initially consider the following approximate mapping for  $\Phi(\mathbf{x}(k))$ :

$$\Phi(\mathbf{x}(k)) \approx \Psi_m \boldsymbol{\beta}_{k,m}. \quad (30)$$

This mapping corresponds to a linear combination of the columns of  $\Psi_m$  according to the components of the vector  $\boldsymbol{\beta}_{k,m}$ . Thus, an optimal value for  $\boldsymbol{\beta}_{k,m}$  is given by [12]

$$\boldsymbol{\beta}_{k,m} = \mathbf{K}_m^{-1} \boldsymbol{\kappa}_{\mathbf{x}(k)}^m \quad (31)$$

where the matrix  $\mathbf{K}_m$  and the vector  $\boldsymbol{\kappa}_{\mathbf{x}(k)}^m$  are defined as follows:

$$\mathbf{K}_m = \begin{bmatrix} k(\mathbf{x}(a), \mathbf{x}(a)) & \dots & k(\mathbf{x}(a), \mathbf{x}(q)) \\ & \ddots & \\ k(\mathbf{x}(q), \mathbf{x}(a)) & \dots & k(\mathbf{x}(q), \mathbf{x}(q)) \end{bmatrix} \quad (32)$$

$$\boldsymbol{\kappa}_{\mathbf{x}(k)}^m = [k(\mathbf{x}(k), \mathbf{x}(a)); \quad \dots \quad k(\mathbf{x}(k), \mathbf{x}(q))]. \quad (33)$$

The square error produced in the approximative mapping from (30) is given by

$$\epsilon_k^2 = \|\Phi(\mathbf{x}(k)) - \Psi_m \boldsymbol{\beta}_{k,m}\|_2^2 \quad (34)$$

Thus, if this error is greater than a threshold, denoted as the constant  $\nu$  in [9], the dictionary should be updated. Using (31), (34) can be rewritten in the following form [9]:

$$\epsilon_k^2 = k(\mathbf{x}(k), \mathbf{x}(k)) - \boldsymbol{\beta}_{k,m}^T \boldsymbol{\kappa}_{\mathbf{x}(k)}^m. \quad (35)$$

Now, assume that kernel principal component estimates can be expressed as a linear combination of dictionary members mapped into the feature space as follows ( $1 \leq j \leq p$ ):

$$\hat{\mathbf{w}}_j(k) = \Psi_m \hat{\boldsymbol{\alpha}}_j(k) \quad (36)$$

$$\bar{\mathbf{w}}_j(k) = \Psi_m \bar{\boldsymbol{\alpha}}_j(k). \quad (37)$$

Using (30), (31), and (37), (23) can be rewritten considering the ON-GHA operating in the feature space as follows<sup>1</sup>:

$$\begin{aligned} y_j^{fs}(k) &= \bar{\mathbf{w}}_j^T(k-1)\Phi(\mathbf{x}(k)) \\ &= \bar{\alpha}_j^T(k-1)\Psi_m^T\Psi_m\beta_{k,m} \\ &= \bar{\alpha}_j^T(k-1)\mathbf{K}_m\mathbf{K}_m^{-1}\kappa_{\mathbf{x}(k)}^m \\ &= \bar{\alpha}_j^T(k-1)\kappa_{\mathbf{x}(k)}^m. \end{aligned} \quad (38)$$

From (21), we have

$$\hat{\mathbf{d}}_j^{fs}(k) = \begin{cases} 0, & j = 1 \\ \Psi_m\hat{\mathbf{h}}_j(k), & 2 \leq j \leq p \end{cases} \quad (39)$$

where

$$\hat{\mathbf{h}}_j(k) = \sum_{i=1}^{j-1} \bar{\alpha}_i(k-1)y_i^{fs}(k). \quad (40)$$

Considering (24), it follows that:

$$\begin{aligned} \hat{\mathbf{e}}_j^{fs}(k) &= \Psi_m[\beta_{k,m} - y_j^{fs}(k)\bar{\alpha}_j(k-1) - \hat{\mathbf{h}}_j(k)] \\ &= \Psi_m\hat{\mathbf{r}}_j(k). \end{aligned} \quad (41)$$

Similarly, using (36), (37), (40), and (41), (25) can be expressed as

$$\begin{aligned} \Psi_m\hat{\alpha}_j(k) &= \Psi_m\bar{\alpha}_j(k-1) + \eta y_j^{fs}(k)\Psi_m[\hat{\mathbf{r}}_j(k) - \hat{\mathbf{h}}_j(k)] \\ \hat{\alpha}_j(k) &= \bar{\alpha}_j(k-1) + \eta y_j^{fs}(k)[\hat{\mathbf{r}}_j(k) - \hat{\mathbf{h}}_j(k)]. \end{aligned} \quad (42)$$

Finally, (26) can be rewritten as follows:

$$\begin{aligned} \Psi_m\bar{\alpha}_j(k) &= \frac{\Psi_m\hat{\alpha}_j(k)}{\sqrt{\hat{\alpha}_j^T(k)\Psi_m^T\Psi_m\hat{\alpha}_j(k)}} \\ \bar{\alpha}_j(k) &= \frac{\hat{\alpha}_j(k)}{\sqrt{\hat{\alpha}_j^T(k)\mathbf{K}_m\hat{\alpha}_j(k)}} = \frac{\hat{\alpha}_j(k)}{\sqrt{\hat{\alpha}_j^T(k)\hat{\mathbf{p}}_j(k)}} \end{aligned} \quad (43)$$

where

$$\hat{\mathbf{p}}_j(k) = \mathbf{K}_m\hat{\alpha}_j(k). \quad (44)$$

Thus, considering matrices  $\hat{\mathbf{A}}$  and  $\bar{\mathbf{A}}$  with columns defined by the vectors  $\hat{\alpha}_j$  and  $\bar{\alpha}_j$  ( $1 \leq j \leq p$ ), respectively, the online orthonormal KHA (ON-KHA) can be written in the following matrix form:

$$\mathbf{y}(k) = \bar{\mathbf{A}}^T(k-1)\kappa_{\mathbf{x}(k)}^m \quad (45)$$

$$\hat{\mathbf{A}}(k) = \bar{\mathbf{A}}(k-1) + \eta[\beta_{k,m}\mathbf{y}^T(k) - \bar{\mathbf{A}}(k-1)\mathbf{M}(k)] \quad (46)$$

$$\hat{\mathbf{P}}(k) = [\hat{\mathbf{p}}_1(k) \ \cdots \ \hat{\mathbf{p}}_p(k)] = \mathbf{K}_m\hat{\mathbf{A}}(k) \quad (47)$$

$$\mathbf{n}(k) = \left[ 1/\sqrt{\hat{\alpha}_1^T(k)\hat{\mathbf{p}}_1(k)}; \ \cdots \ 1/\sqrt{\hat{\alpha}_p^T(k)\hat{\mathbf{p}}_p(k)} \right] \quad (48)$$

$$\bar{\mathbf{A}}(k) = \hat{\mathbf{A}}(k)\text{diag}[\mathbf{n}(k)] \quad (49)$$

with

$$\mathbf{M}(k) = 2\text{ut}[\mathbf{y}(k)\mathbf{y}^T(k)] - \text{diag}[\mathbf{y}(k)\mathbf{y}^T(k)] \quad (50)$$

where the operator  $\text{ut}[\mathbf{A}]$  returns the upper triangular part of the matrix  $\mathbf{A}$ , while  $\text{diag}[\mathbf{n}]$  produces a diagonal matrix, whose elements are defined by the components of the vector  $\mathbf{n}$ .

A less computational complex version of ON-KHA may be derived if the normalization involved in cGSO is ignored. In this case, we have the online orthogonal KHA (OO-KHA), which may be

considered as a kernelized version of the SGA [26] algorithm. The OO-KHA equations are given by

$$\mathbf{y}(k) = \mathbf{A}^T(k-1)\kappa_{\mathbf{x}(k)}^m \quad (51)$$

$$\mathbf{A}(k) = \mathbf{A}(k-1) + \eta[\beta_{k,m}\mathbf{y}^T(k) - \mathbf{A}(k-1)\mathbf{M}(k)] \quad (52)$$

where the matrix  $\mathbf{M}(k)$  is given by (50). Note that the iterative equations of KHA, OKHA, and SubKHA may be written in the same form of (52), but employing a matrix  $\mathbf{M}$  given by

$$\mathbf{M}(k) = \text{ut}[\mathbf{y}(k)\mathbf{y}^T(k)]. \quad (53)$$

Thus, comparing (50) and (53), we conclude that the proposed orthogonalization scheme does not significantly increase the computational efforts involved in KHA.

#### A. Dictionary Augmentation

In case of member inclusion, the new dictionary mapped into the feature space will be given by

$$\Psi_{m+1} = [\Psi_m \ \Phi(\mathbf{x})]. \quad (54)$$

Thus, the mapping of the current input sample according to (30) will have a value of  $\beta_{k,m+1}$  corresponding to

$$\beta_{k,m+1} = \begin{bmatrix} \mathbf{o}_{m \times 1} \\ 1 \end{bmatrix} \quad (55)$$

where all  $m$  components of the column vector  $\mathbf{o}_{m \times 1}$  are zero.

Assume that component estimates before and after member inclusion are given by  $\mathbf{W}^m$  and  $\mathbf{W}^{m+1}$ , respectively. Since dictionary member inclusion should not modify these estimates, we have

$$\mathbf{W}^{m+1} = \mathbf{W}^m \quad (56)$$

$$\Psi_{m+1}\mathbf{A}^{m+1} = \Psi_m\mathbf{A}^m \quad (57)$$

$$\mathbf{A}^{m+1} = \begin{bmatrix} \mathbf{A}^m \\ \mathbf{o}_{1 \times p} \end{bmatrix} \quad (58)$$

where  $\mathbf{A}^m$  and  $\mathbf{A}^{m+1}$  will correspond to the matrix  $\bar{\mathbf{A}}$  (ON-KHA) or  $\mathbf{A}$  (OO-KHA) before and after member inclusion, according to the algorithm considered, respectively.

Finally, the inclusion of a dictionary member modifies the matrices  $\mathbf{K}_m^{-1}$  and  $\mathbf{K}_m$  used in (31) and (47). Low computational complex formulas to calculate these matrices are the following [9]:

$$\mathbf{K}_{m+1}^{-1} = \begin{bmatrix} \mathbf{K}_m^{-1} & \mathbf{o}_{p \times 1} \\ \mathbf{o}_{1 \times p} & 0 \end{bmatrix} + \frac{1}{\epsilon_k^2} \begin{bmatrix} -\beta_{k,m} \\ 1 \end{bmatrix} \begin{bmatrix} -\beta_{k,m}^T & 1 \end{bmatrix} \quad (59)$$

$$\mathbf{K}_{m+1} = \begin{bmatrix} \mathbf{K}_m & \kappa_{\mathbf{x}(k)}^m \\ \kappa_{\mathbf{x}(k)}^m & k(\mathbf{x}(k), \mathbf{x}(k)) \end{bmatrix} \quad (60)$$

where the vector  $\beta_{k,m}$  and the variable  $\epsilon_k^2$  are given by (31) and (35), respectively.

#### B. Computational Cost

The ON-KHA algorithm is summarized in Table I. The algorithm is similar to OKHA in the phase of dictionary evaluation, as both follows the ALD criterion. This evaluation is conducted in Steps 4–7 and has a computational cost per iteration of  $O(m^2)$ . If the dictionary is updated, the matrix  $\mathbf{K}_m^{-1}$  is incrementally produced in Step 10 by a simple  $O((m+1)^2)$  procedure. The update of kernel component estimates, which occurs in Steps 16–19, is an  $O((m/2)p^2)$  procedure, similar to the OKHA and SubKHA algorithms. The normalization phase adds  $O(m^2p)$  operations, mainly due to Step 20. The learning rate follows a heuristically motivated exponential decay formula stated in Step 17. Clearly, the overall ON-KHA computational cost is  $O(mp(m + (p/2) + (m/p)))$ . For the OO-KHA, the steps 20 to 22 are eliminated, resulting in additional  $O((p(p-1))/2)$  operations when compared with OKHA and SubKHA.

<sup>1</sup>Note that the orthonormalization of the component estimates of ON-GHA will occur in the feature space.

TABLE I  
ON-KHA ALGORITHM

```

1: Define a subspace distance threshold limit  $\nu$ , an initial learning rate value  $\eta(0)$ ,
   and a learning decay factor  $\gamma_\eta$ .
2: Initialize:  $m = 1$ ,  $\mathbf{K}_1 = k(\mathbf{x}(1), \mathbf{x}(1))$ ,  $\mathbf{K}_1^{-1} = 1/k(\mathbf{x}(1), \mathbf{x}(1))$ ,  $\mathbf{D}_1 =$ 
    $[\mathbf{x}(1)]$  and  $\bar{\mathbf{A}}(1) = [a_1 \dots a_p]$ , where  $a_1 \dots a_p$  are small random values.
3: for all  $k \geq 2$  do
4:    $\kappa_{\mathbf{x}(k)}^m = [k(\mathbf{d}_1, \mathbf{x}(k)) \dots k(\mathbf{d}_m, \mathbf{x}(k))]$ , where  $\mathbf{d}_i$  is the  $i$ -th column of
      $\mathbf{D}_m$ .
5:    $\beta_{k,m} = \mathbf{K}_m^{-1} \kappa_{\mathbf{x}(k)}^m$ 
6:    $\epsilon_k^2 = k(\mathbf{x}(k), \mathbf{x}(k)) - \beta_{k,m}^T \kappa_{\mathbf{x}(k)}^m$ 
7:   if  $\epsilon_k^2 \geq \nu$  then
8:      $m = m + 1$ 
9:      $\mathbf{K}_m = \begin{bmatrix} \mathbf{K}_{m-1} & \kappa_{\mathbf{x}(k)}^{m-1} \\ \kappa_{\mathbf{x}(k)}^{m-1 T} & k(\mathbf{x}(k), \mathbf{x}(k)) \end{bmatrix}$ 
10:     $\mathbf{K}_m^{-1} = \begin{bmatrix} \mathbf{K}_{m-1}^{-1} & \mathbf{0}_{(m-1) \times 1} \\ \mathbf{0}_{1 \times (m-1)} & 0 \end{bmatrix} + \frac{1}{\epsilon_k^2} \begin{bmatrix} -\beta_{k,m-1} \\ -\beta_{k,m-1}^T & 1 \end{bmatrix}$ 
11:     $\beta_k = [\mathbf{0}_{1 \times (m-1)}; 1]$ 
12:     $\mathbf{D}_m = [\mathbf{D}_{m-1} \ \mathbf{x}(k)]$ 
13:     $\bar{\mathbf{A}}(k-1) = \begin{bmatrix} \bar{\mathbf{A}}(k-1) \\ \mathbf{0}_{1 \times p} \end{bmatrix}$ 
14:     $\kappa_{\mathbf{x}(k)}^m = \begin{bmatrix} \kappa_{\mathbf{x}(k)}^{m-1} \\ k(\mathbf{x}(k), \mathbf{x}(k)) \end{bmatrix}$ 
15:  end if
16:   $\mathbf{y}(k) = \bar{\mathbf{A}}(k-1) \kappa_{\mathbf{x}(k)}^m$ 
17:   $\eta(k) = \gamma_\eta \eta(k-1)$ 
18:   $\mathbf{M}(k) = 2\text{ut}[\mathbf{y}(k) \mathbf{y}^T(k)] - \text{diag}[\mathbf{y}(k) \mathbf{y}^T(k)]$ 
19:   $\hat{\mathbf{A}}(k) = \bar{\mathbf{A}}(k-1) + \eta(k) \left( \beta_{k,m}^T \mathbf{y}(k) - \bar{\mathbf{A}}(k-1) \mathbf{M}(k) \right)$ 
20:   $\hat{\mathbf{P}}(k) = [\hat{\mathbf{p}}_1(k) \dots \hat{\mathbf{p}}_p(k)] = \mathbf{K}_m \hat{\mathbf{A}}(k)$ 
21:   $\mathbf{n}(k) = [1/\sqrt{\hat{\alpha}_1^T(k) \hat{\mathbf{p}}_1(k)}; \dots; 1/\sqrt{\hat{\alpha}_p^T(k) \hat{\mathbf{p}}_p(k)}]$ 
22:   $\bar{\mathbf{A}}(k) = \hat{\mathbf{A}}(k) \text{diag}[\mathbf{n}(k)]$ 
23: end for

```

#### IV. KPCA QUALITY EXTRACTION ASSESSMENT

To evaluate the convergence speed and accuracy of the proposed algorithms, we considered the average value of the cosine of the angle between components produced by a reference extraction ( $\mathbf{w}_{i,\text{ref}}$ ) and by the algorithms upon evaluation ( $\mathbf{w}_{i,\text{alg}}$ ) as follows:

$$\text{AC} = \frac{1}{p} \sum_{i=1}^p \frac{\mathbf{w}_{i,\text{alg}}^T \mathbf{w}_{i,\text{ref}}}{\|\mathbf{w}_{i,\text{alg}}\|_2 \|\mathbf{w}_{i,\text{ref}}\|_2}. \quad (61)$$

Here, the reference components were produced by standard KPCA. Clearly, this index attains a maximum value of 1 when both components only differ in their norms. This value, here denoted as the average cosine (AC) index, can be expressed in the following form [19, Sec. IV]:

$$\text{AC} = \frac{1}{p} \sum_{i=1}^p \frac{\alpha_{i,\text{alg}}^T [\mathbf{K}_m \ \mathbf{K}_{\text{dic} \times \text{dic}}^{\text{ref}}] \alpha_{i,\text{ref}}}{\sqrt{\alpha_{i,\text{alg}}^T \mathbf{K}_m \alpha_{i,\text{alg}}} \sqrt{\alpha_{i,\text{ref}}^T \mathbf{K}^{\text{ref}} \alpha_{i,\text{ref}}}} \quad (62)$$

where the matrix  $\mathbf{K}_m$  is defined by (32), and the Gram matrix employed in the reference extraction is denoted by  $\mathbf{K}^{\text{ref}}$ . The matrix  $\mathbf{K}_{\text{dic} \times \text{dic}}^{\text{ref}}$  is obtained by the following decomposition of the matrix  $\mathbf{K}^{\text{ref}}$ :

$$\mathbf{K}^{\text{ref}} = \begin{bmatrix} \mathbf{K}_m & \mathbf{K}_{\text{dic} \times \text{dic}}^{\text{ref}} \\ \mathbf{K}_{\text{dic} \times \text{dic}}^{\text{ref}} & \mathbf{K}_{\text{dic} \times \text{dic}}^{\text{ref}} \end{bmatrix}. \quad (63)$$

#### V. RESULTS

The proposed algorithms were experimentally evaluated using the USPS database of handwritten digits [27]. This set is composed of grayscale handwritten digit images of size  $16 \times 16$ , which were coded using integer numbers in the range of 0–255. As we reproduced one experiment conducted in [28], our evaluation data set involved only the first 100 images related to the digits 1 to 3 from USPS, and

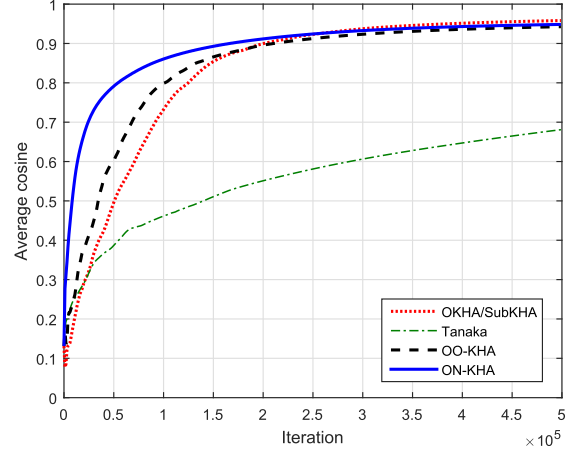


Fig. 1. Mean curve of the AC index obtained by each method considering the full sample dictionary (see text).

we considered the extraction of 16 kernel components using a radial basis function kernel with  $\sigma = 8$ .

The accuracy and convergence speed of the proposed ON-KHA and OO-KHA were experimentally compared with SubKHA, OKHA, and Tanaka algorithms using the AC index as the figure of merit. Two evaluation scenarios were considered for kernel component expansion: using all data set members (full sample dictionary), and based on a small subset of them (small sample dictionary). In the last case, the parameters of each method were tuned to result in about 49 members (16.3% of the total).

We produced 25 simulations to each method, considering that the elements of the matrix  $\mathbf{A}$  follow a random zero-mean Gaussian distribution with variance equal to 0.01. To define the inputs of the algorithms, each one of the 300 selected images was converted into a column vector having 256 dimensions, and was normalized to be in 0 to 1 range. These vectors were randomly presented to the algorithms. The initial learning rate value ( $\eta$ ) and the learning rate decay factor ( $\gamma_\eta$ ) were chosen as 0.05 and 0.999995, respectively. Simulations were conducted for 500 000 iterations.

All algorithms, except the SubKHA and the Tanaka, employed a value of  $\nu$  equal to  $10^{-3}$  and 0.25 for full and small sample dictionaries, respectively. Similarly, for the SubKHA, we considered the maximum number of dictionary members given by 300 and 49, and the constant  $\alpha$  was chosen equal to 1. For the Tanaka algorithm, we utilized: decreasing weighting factors starting from 0.9 with a regular step of 0.05; a factor  $\beta$  equal to 0.9999; the coherence threshold equal to 1 and 0.742 (full and small sample dictionaries, respectively); and the matrix  $\mathbf{Q}$  with random elements following a zero-mean Gaussian distribution whose variance equals 0.1.

The mean curves of the AC index produced by each method in case of full and small dictionaries are shown in Figs. 1 and 2. In the case of full dictionary, the OKHA and SubKHA curves are superimposed, as expected. The Tanaka algorithm was the less accurate and slowest converging method in this comparison. The proposed OO-KHA performed better than the OKHA and SubKHA, reaching the performance of the ON-KHA in about 200 000 iterations. Assuming an AC of 0.8, the OKHA and SubKHA took about 125 000 iterations to reach this value, whereas the ON-KHA and OO-KHA achieved it in about 55 000 and 100 000 iterations, respectively. Thus, the proposed methods were  $\approx 2.27$  (ON-KHA) and  $\approx 1.25$  (OO-KHA) times faster in this experiment.

Concerning the small sample case, a similar behavior to full-dictionary simulations was observed, as expected. However, the use of

TABLE II  
MEAN AND STANDARD DEVIATION VALUES FOR THE VALUE OF AC CONSIDERING DIFFERENT ITERATIONS,  
ALGORITHMS, AND DICTIONARIES SIZES (SEE TEXT)

Full sample dictionary							
Iteration	OKHA	Tanaka	SubKHA	OO-KHA	ON-KHA	$\chi^2(4)$	$p$ value
50,000	0.4887 $\pm$ 0.0572	0.3841 $\pm$ 0.0719	0.4887 $\pm$ 0.0572	0.5974 $\pm$ 0.0513	<b>0.7891 <math>\pm</math> 0.0537</b>	101.7	$p < 0.001$
150,000	<b>0.8525 <math>\pm</math> 0.0579</b>	0.5097 $\pm$ 0.0594	<b>0.8525 <math>\pm</math> 0.0579</b>	<b>0.8650 <math>\pm</math> 0.0777</b>	<b>0.8923 <math>\pm</math> 0.0566</b>	63.17	$p < 0.001$
500,000	<b>0.9582 <math>\pm</math> 0.0306</b>	0.6809 $\pm$ 0.0869	<b>0.9582 <math>\pm</math> 0.0306</b>	<b>0.9426 <math>\pm</math> 0.0571</b>	<b>0.9484 <math>\pm</math> 0.0432</b>	59.24	$p < 0.001$
Small sample dictionary							
50,000	<b>0.4220 <math>\pm</math> 0.041</b>	0.3294 $\pm$ 0.0765	0.2637 $\pm$ 0.0185	0.4905 $\pm$ 0.0486	<b>0.6327 <math>\pm</math> 0.0499</b>	106.71	$p < 0.001$
150,000	<b>0.6320 <math>\pm</math> 0.0644</b>	0.4249 $\pm$ 0.0727	0.3820 $\pm$ 0.0446	<b>0.6718 <math>\pm</math> 0.0582</b>	<b>0.6800 <math>\pm</math> 0.0465</b>	91.63	$p < 0.001$
500,000	<b>0.6917 <math>\pm</math> 0.0531</b>	0.5518 $\pm$ 0.0830	0.4229 $\pm$ 0.0439	<b>0.7021 <math>\pm</math> 0.0510</b>	<b>0.6965 <math>\pm</math> 0.0499</b>	83.46	$p < 0.001$

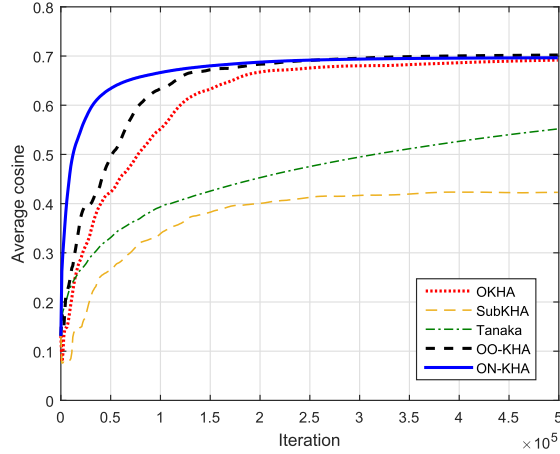


Fig. 2. Mean curve of the AC index obtained by each method considering the small sample dictionary (see text).

only  $\approx 16\%$  of data set samples to expand kernel components resulted in lower values of the AC index in all algorithms, as expected. Another remarkable result was the slow convergence observed for the SubKHA, caused by a frequent exchange of dictionary members. In this case, the SubKHA performed worse than the Tanaka algorithm. Again, the proposed ON-KHA and OO-KHA converged faster than the remaining algorithms in about 100 000 and 150 000 iterations, respectively.

We used the nonparametric test of Kruskal–Wallis (KW) [29] to statistically evaluate the values of the AC index attained by each method at the following iterations: 50 000, 150 000, and 500 000, considering both small and full sample dictionaries. Post-hoc analysis was performed using the Tuckey’s HSD [29] procedure. We considered the confidence level equal to 0.10. Table II summarizes the results, where boldface letters indicate statistically similar values. Since the values of  $p$  reported in Table II are lower than the confidence level for all simulations, we should refuse the hypothesis that the methods performed similarly in both evaluation scenarios. At iteration 50 000, according to the post-hoc tests, the proposed ON-KHA was the best performing method, achieving the highest value for the AC index for both dictionary sizes, and a value of  $p$  lower than 0.0747 in the worst case comparison, i.e., the values of  $p$  attained in all pairwise comparisons satisfy  $p < 0.0747$ . At iteration 150 000, the ON-KHA still showed a higher value of the AC index than the proposed OO-KHA, but both performed similar to most of the alternative methods, concerning small and full-dictionary sizes ( $p < 0.4707$ ). However, the ON-KHA outperformed Tanaka algorithm ( $p < 0.001$ ) in full-dictionary simulations, and the Tanaka ( $p < 0.001$ ) and the SubKHA ( $p < 0.001$ ) in

small-dictionary experiments. At iteration 500 000, the proposed methods and OKHA practically showed the same value for the AC index for both dictionary sizes ( $p > 0.9932$ ).

## VI. CONCLUSION

In this brief, we have proposed two new online KPCA extraction algorithms, exploiting an implicit orthogonalization (OO-KHA) or orthonormalization (ON-KHA) of kernel component estimates in the feature space. Both methods involve simple equations, which eases their embedded implementation and possesses a similar computational cost to the state-of-the-art online KPCA algorithms. Simulations have shown impressive gains in the accuracy of the components extracted and in the convergence speed.

## REFERENCES

- [1] T. D. Sanger, “Optimal unsupervised learning in a single-layer linear feedforward neural network,” *Neural Netw.*, vol. 2, no. 6, pp. 459–473, 1989.
- [2] K. I. Kim, M. Franz, and B. Scholkopf, “Iterative kernel principal component analysis for image modeling,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 9, pp. 1351–1366, Sep. 2005.
- [3] I. T. Jolliffe, *Principal Component Analysis*, 2nd ed. New York, NY, USA: Springer, 2002.
- [4] H. Hoffmann, “Kernel PCA for novelty detection,” *Pattern Recognit.*, vol. 40, no. 3, pp. 863–874, 2007.
- [5] B. Schölkopf, A. Smola, and K.-R. Müller, “Nonlinear component analysis as a kernel eigenvalue problem,” *Neural Comput.*, vol. 10, no. 5, pp. 1299–1319, Jul. 1998.
- [6] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 3rd ed. Baltimore, MD, USA: The Johns Hopkins Univ. Press, 1996.
- [7] T.-J. Chin and D. Suter, “Incremental kernel principal component analysis,” *IEEE Trans. Image Process.*, vol. 16, no. 6, pp. 1662–1674, Jun. 2007.
- [8] W. Liao, A. Pizurica, W. Philips, and Y. Pi, “A fast iterative kernel PCA feature extraction for hyperspectral images,” in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, Sep. 2010, pp. 1317–1320.
- [9] P. Honeine, “Online kernel principal component analysis: A reduced-order model,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 9, pp. 1814–1826, Sep. 2012.
- [10] Y. Washizawa, “Adaptive subset kernel principal component analysis for time-varying patterns,” *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 23, no. 12, pp. 1961–1973, Dec. 2012.
- [11] T. Tanaka, Y. Washizawa, and A. Kuh, “Adaptive kernel principal components tracking,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Mar. 2012, pp. 1905–1908.
- [12] Y. Engel, S. Mannor, and R. Meir, “The kernel recursive least-squares algorithm,” *IEEE Trans. Signal Process.*, vol. 52, no. 8, pp. 2275–2285, Aug. 2004.
- [13] C. Richard, J. C. M. Bermudez, and P. Honeine, “Online prediction of time series data with kernels,” *IEEE Trans. Signal Process.*, vol. 57, no. 3, pp. 1058–1067, Mar. 2009.
- [14] J. Yang, Y. Zhao, and H. Xi, “Weighted rule based adaptive algorithm for simultaneously extracting generalized eigenvectors,” *IEEE Trans. Neural Netw.*, vol. 22, no. 5, pp. 800–806, May 2011.
- [15] K. Abed-Meraim, S. Attallah, A. Chkeif, and Y. Hua, “Orthogonal Oja algorithm,” *IEEE Signal Process. Lett.*, vol. 7, no. 5, pp. 116–119, May 2000.

- [16] S. Attallah and K. Abed-Meraim, "Fast algorithms for subspace tracking," *IEEE Signal Process. Lett.*, vol. 8, no. 7, pp. 203–206, Jul. 2001.
- [17] K. Abed-Meraim, A. Chkeif, and Y. Hua, "Fast orthonormal PAST algorithm," *IEEE Signal Process. Lett.*, vol. 7, no. 3, pp. 60–62, Mar. 2000.
- [18] T. Tanaka, "Dictionary-based online kernel principal subspace analysis with double orthogonality preservation," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Apr. 2015, pp. 4045–4049.
- [19] *Supplementary Material for the paper Improving KPCA Online Extraction by Orthonormalization in the Feature Space*. [Online]. Available: <http://dx.doi.org/10.1109/TNNLS.2017.2660441>
- [20] A. Dax, "A modified Gram–Schmidt algorithm with iterative orthogonalization and column pivoting," *Linear Algebra Appl.*, vol. 310, nos. 1–3, pp. 25–42, May 2000.
- [21] W. Hoffmann, "Iterative algorithms for Gram–Schmidt orthogonalization," *Computing*, vol. 41, no. 4, pp. 335–348, 1989.
- [22] Å. Björck, "Numerics of Gram–Schmidt orthogonalization," *Linear Algebra Appl.*, vols. 197–198, pp. 297–316, Jan./Feb. 1994.
- [23] J. W. Daniel, W. B. Gragg, L. Kaufman, and G. W. Stewart, "Reorthogonalization and stable algorithms for updating the Gram–Schmidt  $QR$  factorization," *Math. Comput.*, vol. 30, no. 136, pp. 772–795, Oct. 1976.
- [24] A. Cichocki and S.-I. Amari, *Adaptive Blind Signal and Image Processing: Learning Algorithms and Applications*. Hoboken, NJ, USA: Wiley, 2002.
- [25] E. Oja, "Simplified neuron model as a principal component analyzer," *J. Math. Biol.*, vol. 15, no. 3, pp. 267–273, 1982.
- [26] E. Oja, "Principal components, minor components, and linear neural networks," *Neural Netw.*, vol. 5, no. 6, pp. 927–935, Nov./Dec. 1992.
- [27] *USPS Handwritten Digits*, accessed on Dec. 12, 2014. [Online]. Available: <http://www.cs.nyu.edu/~roweis/data.html>
- [28] S. Günter, N. N. Schraudolph, and S. V. N. Vishwanathan, "Fast iterative kernel principal component analysis," *J. Mach. Learn. Res.*, vol. 8, pp. 1893–1918, Aug. 2007.
- [29] D. J. Sheskin, *Handbook of Parametric and Nonparametric Statistical Procedures*, 4th ed. London, U.K.: Chapman & Hall, 2007.