

Reduced Kernel Dictionary Learning

Denis C. ILIE-ABLACHIM

*Faculty of Automatic Control and Computers
University Politehnica of Bucharest
denis.ilie_ablachim@upb.ro*

Bogdan DUMITRESCU

*Faculty of Automatic Control and Computers
University Politehnica of Bucharest
bogdan.dumitrescu@upb.ro*

Abstract—In this paper we present new algorithms for training reduced-size nonlinear representations in the Kernel Dictionary Learning (KDL) problem. Standard KDL has the drawback of a large size of the kernel matrix when the data set is large. There are several ways of reducing the kernel size, notably Nyström sampling. We propose here a method more in the spirit of dictionary learning, where the kernel vectors are obtained with a trained sparse representation of the input signals. Moreover, we optimize directly the kernel vectors in the KDL process, using gradient descent steps. We show with three data sets that our algorithms are able to provide better representations, despite using a small number of kernel vectors, and also decrease the execution time with respect to KDL.

I. INTRODUCTION

Dictionary Learning (DL) is a representation learning method that aims to find a sparse representation for a set of signals, \mathbf{Y} , represented as a matrix with N columns (signals) of size m . The representation is achieved by computing a dictionary \mathbf{D} of size $m \times n$ and a sparse representation \mathbf{X} of size $n \times N$ such that a good approximation $\mathbf{Y} \approx \mathbf{D}\mathbf{X}$ is obtained. Most applications with dictionary learning are in image denoising, inpainting, signal reconstruction, clustering or classification.

In this paper we present new methods of dictionary learning that produce sparse representations in both linear and nonlinear spaces, starting from the Kernel Dictionary Learning (KDL) idea [1], [2]. In a first approach, this is done in two stages; a linear representation is built and the resulting optimized dictionary is used unchanged in the nonlinear space as kernel vectors during the training procedure. In the second approach, the kernel vectors are optimized alongside with the nonlinear representation. The main advantage of this method is the use of a reduced matrix, \mathbf{D} , containing the kernel vectors, which is also the dictionary in a standard DL problem. By the use of the reduced matrix, the built-in kernel is smaller and thus the problem complexity is reduced.

As KDL has a high complexity when N is large (and so the kernel matrix is large), solutions have been adopted from other problems where kernels appear. In Large-Scale Kernel Machines, various kernel enhancement or resizing strategies have been used, such as Nyström Sampling [3], [4] or Random Fourier Features (RFF) [5]. The first one computes a rank \hat{m}

approximation $\hat{\mathbf{K}}$ of the kernel matrix \mathbf{K} . The whole procedure consists in approximating the nonlinear mapping function $\varphi(\mathbf{Y})$ with a matrix $\hat{\mathbf{Y}}$, containing a compressed version of the original signals. Compared to the original problem, the new problem no longer requires a high computational cost. The Random Fourier Features method proposes to map the input signals to a randomized low-dimensional feature space. More exactly, the inner product, used in the kernel trick, is replaced with a randomized map $k(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x})^\top \varphi(\mathbf{y}) = z(\mathbf{x})^\top z(\mathbf{y})$, where $z: \mathbb{R}^m \mapsto \mathbb{R}^{\hat{m}}$ and $\hat{m} \ll m$. Both methods enable the use of fast linear methods, which will further use the resulting features. The use of reduced kernels dates from Support Vector Machine times, significant contributions being [6], [7], with applications in classification [8]. Other learning methods where a reduced kernel appears can be found in [9]. In these early works, kernel vectors are usually selected (randomly or with some heuristic) as a subset of input signals. Finally, there are KDL substitutes like [10], where the nonlinear transformation is performed by a neural encoder-decoder, with standard DL on the encoded signals.

The paper is organized as follows. In Section II we review the standard DL and KDL problems and the most common algorithm for solving them. In Section III we present our own contribution, named the Reduced Kernel Dictionary Learning (RKDL) problem, under three different scenarios, in which the kernel vectors are: a) the result of a DL problem, solved a priori; b) optimized together with the other KDL variables, using gradient descent, with an objective that is directly related to that of the KDL problem; c) optimized like before, but using a mixed objective that combines the nonlinear representation of the signals with the linear representation of the kernel vectors. The algorithms for b) and c) are new and their representation error is smaller. Section IV contains the experimental results, obtained on three public data sets: Digits, MNIST [11] and CIFAR-10 [12], under the three proposed scenarios.

II. KERNEL DICTIONARY LEARNING

The DL problem is formulated as following

$$\begin{aligned} \min_{\mathbf{D}, \mathbf{X}} \quad & \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_F^2 \\ \text{s.t.} \quad & \|\mathbf{x}_\ell\|_0 \leq s_x, \ell = 1 : N \\ & \|\mathbf{d}_j\| = 1, j = 1 : n, \end{aligned} \quad (1)$$

where $\|\cdot\|_0$ represents the 0-pseudo-norm, s_x is the sparsity level and \mathbf{d}_j is a column (named also atom) of \mathbf{D} .

This work was supported by grants of the Romanian Ministry of Research, Innovation and Digitization, CCCDI - UEFISCDI, project numbers PN-III-P2-2.1-PED-2019-3248 and PN-III-P4-PCE-2021-0154, within PNCDI III.

The standard dictionary learning problem can be solved by using simple strategies. In order to overcome the nonconvexity and the huge dimension of the problem, the most usual optimization procedure iterates two basic steps and is also known as DL by Alternate Optimization. In this way, the problem is divided in two subproblems: sparse coding and dictionary update. By alternating these two stages for a given number of iterations, good local solutions can be obtained. A simple iteration consists of computing the sparse representation \mathbf{X} , using Orthogonal Matching Pursuit (OMP) [13] while the dictionary \mathbf{D} is fixed, and then updating the dictionary atoms successively while the sparse representation is fixed. There are several methods of dictionary learning [14], but the one of interest to us is AK-SVD [15], [16] due to its low complexity and good performance.

In the DL problem, the input data are modeled through a linear representation, which in some cases may be seen as a limitation. In order to overcome this drawback, kernel representations can be used for a better quantification of similarities or differences between input vectors.

The kernel representation is an extension to nonlinearity. We do this by associating to a data vector $\mathbf{y} \in \mathbb{R}^m$ a feature vector $\varphi(\mathbf{y}) \in \mathbb{R}^{\tilde{m}}$, where $\varphi : \mathbb{R}^m \rightarrow \mathbb{R}^{\tilde{m}}$ is a nonlinear function. Typically, Mercer kernels are used, which can be expressed as a scalar product of feature vector functions $k(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x})^\top \varphi(\mathbf{y})$. In the last form, the scalar product can be replaced with a specific function definition, such as radial basis function (RBF) $k(\mathbf{x}, \mathbf{y}) = \exp(-\frac{\|\mathbf{x}-\mathbf{y}\|_2^2}{2\sigma^2})$ or polynomial kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} + \alpha)^\beta$.

The Kernel Dictionary Learning (KDL) [1], [2] problem is

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{Z}} \quad & \|\varphi(\mathbf{Y}) - \varphi(\mathbf{Y})\mathbf{A}\mathbf{Z}\|_F^2 \\ \text{s.t.} \quad & \|\mathbf{z}_\ell\|_0 \leq s_z, \ell = 1 : N \\ & \|\varphi(\mathbf{Y})\mathbf{a}_j\| = 1, j = 1 : n, \end{aligned} \quad (2)$$

where $\varphi(\mathbf{Y})$ represents the nonlinear extension of data and $\varphi(\mathbf{Y})\mathbf{A}$ is the kernel dictionary, where \mathbf{A} is the coefficients matrix of size $N \times p$. Depending on the used data set, the problem can be difficult to solve due the large kernel matrix $\mathbf{K}_{YY} = \varphi(\mathbf{Y})^\top \varphi(\mathbf{Y})$ that results from the trace form of the objective function. In this case, the problems with large data sets can involve large volume of memory and long execution times. The KDL problem can also be solved by alternate optimization. The sparse representation is computed according to the Kernel OMP algorithm [1]. The columns of the matrix \mathbf{A} are sequentially updated with an algorithm inspired by AK-SVD, while the representation matrix \mathbf{X} is fixed.

Both algorithms, AK-SVD and Kernel AK-SVD, alongside with the sparse representation computation are presented in [14].

III. REDUCED KERNEL DICTIONARY LEARNING

Nonlinear space can extend the horizon of data representation. However, KDL has disadvantages when the number of available signals is large. The size of the kernel increases in proportion to the size of the data. Thus, the problem becomes more complex from a numerical point of view.

In order to overcome this limitations we use a reduced space on which the kernel matrix is built. This strategy is implemented by using as kernel vectors not the full set of signals, \mathbf{Y} , but a smaller set of vectors, \mathbf{D} , trained with DL as a dictionary for linear representations, thus replacing (2) with

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{Z}} \quad & \|\varphi(\mathbf{Y}) - \varphi(\mathbf{D})\mathbf{A}\mathbf{Z}\|_F^2 \\ \text{s.t.} \quad & \|\mathbf{z}_\ell\|_0 \leq s_z, \ell = 1 : N \\ & \|\varphi(\mathbf{D})\mathbf{a}_j\| = 1, j = 1 : n. \end{aligned} \quad (3)$$

This problem has two advantages that can be used as needed. If the number of training signals is very large, a dictionary with a much smaller number of atoms can be used. On the other hand, we have problems where the number of training signals is small or the given signals are not representative enough for the representation problem. In this case it is recommended to use large dictionaries. In our work the case of interest is the first one.

A. Standard Reduced Kernel Dictionary Learning

A first approach to (3) consists of solving it in two steps. In the first step we design \mathbf{D} by solving a DL problem and thus obtaining a trained dictionary. We call this method RKDL-D, due the use of matrix \mathbf{D} ; it was introduced in [17]. Here is a brief reminder of the atom update step of RKDL-D. (Sparse representation can be easily derived from Kernel OMP.)

Expressing the objective of (3) in its trace form and isolating the current atom \mathbf{a}_j , we can write

$$\text{Tr} \left[\begin{aligned} & \left(\varphi^\top(\mathbf{Y}) - \sum_{i \neq j} \mathbf{z}_i \mathbf{a}_i^\top \varphi^\top(\mathbf{D}) - \mathbf{z}_j \mathbf{a}_j^\top \varphi^\top(\mathbf{D}) \right) \\ & \left(\varphi(\mathbf{Y}) - \varphi(\mathbf{D}) \sum_{i \neq j} \mathbf{a}_i \mathbf{z}_i^\top - \varphi(\mathbf{D}) \mathbf{a}_j \mathbf{z}_j^\top \right) \end{aligned} \right].$$

We compute the partial derivatives with respect to the current atom \mathbf{a}_j

$$\frac{\partial \text{Tr}(\cdot)}{\partial \mathbf{a}_j} = 2\|\mathbf{z}_j\|^2 \mathbf{K}_{DD} \mathbf{a}_j + 2\mathbf{K}_{DD} \mathbf{R} \mathbf{z}_j - 2\mathbf{K}_{YD} \mathbf{z}_j$$

and the current sparse representation vector \mathbf{z}_j

$$\frac{\partial \text{Tr}(\cdot)}{\partial \mathbf{z}_j} = 2\mathbf{a}_j^\top \mathbf{K}_{DD} \mathbf{a}_j \mathbf{z}_j + 2\mathbf{R}^\top \mathbf{K}_{DD} \mathbf{a}_j - 2\mathbf{K}_{YD} \mathbf{a}_j.$$

where we have used the notations $\mathbf{K}_{YD} = K(\mathbf{Y}, \mathbf{D}) = \varphi(\mathbf{Y})^\top \varphi(\mathbf{D})$, $\mathbf{K}_{DD} = K(\mathbf{D}, \mathbf{D})$ and $\mathbf{R} = \sum_{i \neq j} \mathbf{a}_i \mathbf{z}_i^\top$.

Setting these derivatives to zero, we obtain the optimal atom (for fixed representation)

$$\mathbf{a}_j = (\|\mathbf{z}_j\|_2^2 \mathbf{K}_{DD})^{-1} (\mathbf{K}_{YD} + \mathbf{K}_{DD} \mathbf{R}) \mathbf{z}_j \quad (4)$$

and optimal representation (for fixed atom)

$$\mathbf{z}_j = (\mathbf{K}_{YD} - \mathbf{R}^\top \mathbf{K}_{DD}) \mathbf{a}_j. \quad (5)$$

The RKDL-D procedure of updating atoms is summarized in Algorithm 1. For brevity, we did not use special notations, but only the signals from \mathbf{Y} where \mathbf{a}_j appears in the representation are involved in the computation.

Algorithm 1: RKDL-D – update step

Data: complementary kernel matrix $\mathbf{K}_{DD} \in \mathbb{R}^{p \times p}$
 partial kernel matrix $\mathbf{K}_{YD} \in \mathbb{R}^{N \times p}$
 current kernel dictionary $\mathbf{A} \in \mathbb{R}^{N \times n}$
 representation matrix $\mathbf{Z} \in \mathbb{R}^{n \times N}$

Result: updated kernel dictionary \mathbf{A} , representation \mathbf{Z}

- 1 Compute sum $\mathbf{S} = \sum_{i=1}^n \mathbf{Z}_i^\top \mathbf{a}_i^\top$
 - 2 **for** $j = 1$ **to** n **do**
 - 3 Modify sum: $\mathbf{R} = \mathbf{S} - \mathbf{Z}_j \mathbf{a}_j^\top$
 - 4 Update atom:
 $\mathbf{a}_j = (\|\mathbf{z}\|_2^2 \mathbf{K}_{DD})^{-1} (\mathbf{K}_{YD}^\top + \mathbf{K}_{DD} \mathbf{R}) \mathbf{Z}_j$
 - 5 Normalize atom: $\mathbf{a}_j \leftarrow \mathbf{a}_j / (\mathbf{a}_j^\top \mathbf{K}_{DD} \mathbf{a}_j)^{\frac{1}{2}}$
 - 6 Update representation: $\mathbf{Z}_j^\top \leftarrow (\mathbf{K}_{YD} - \mathbf{R} \mathbf{K}_{DD}) \mathbf{a}_j$
 - 7 Recompute error: $\mathbf{S} = \mathbf{R} + \mathbf{Z}_j \mathbf{a}_j^\top$
-

B. Optimized Reduced Kernel Dictionary Learning

The improvement that we propose here is to update the dictionary \mathbf{D} , containing the kernel vectors, during the non-linear optimization procedure. We keep the idea of alternate optimization. The matrices \mathbf{Z} , \mathbf{A} and \mathbf{D} are updated successively. As above, for \mathbf{Z} we use Kernel OMP and the atoms of \mathbf{A} are updated as described by Algorithm 1. Updating the dictionary \mathbf{D} must be done with a different procedure, detailed below. Since the dictionary \mathbf{D} is updated together with the nonlinear representation, we call the resulting method Optimized Reduced Kernel Dictionary Learning (ORKDL-D).

In order to solve the optimization problem for \mathbf{D} , we update each column \mathbf{d}_j independently, by the use of the trace form of the objective function of (3):

$$\text{Tr} [\mathbf{K}_{YY} - 2\mathbf{K}_{YD} \mathbf{A} \mathbf{Z} + \mathbf{Z}^\top \mathbf{A}^\top \mathbf{K}_{DD} \mathbf{A} \mathbf{Z}].$$

We compute the partial derivatives with respect to the i th element of the current column \mathbf{d}_j , for both nonlinear terms

$$\begin{aligned} \frac{\partial \text{Tr}[\mathbf{K}_{YD} \mathbf{A} \mathbf{Z}]}{\partial \mathbf{d}_{ij}} &= \text{Tr} \left[\left(\frac{\partial \text{Tr}[\mathbf{K}_{YD} \mathbf{A} \mathbf{Z}]}{\partial \mathbf{K}_{YD}} \right)^\top \cdot \frac{\partial \mathbf{K}_{YD}}{\partial \mathbf{d}_{ij}} \right] \\ &= \text{Tr} \left[\mathbf{A} \mathbf{Z} \cdot \frac{\partial \mathbf{K}_{YD}}{\partial \mathbf{d}_{ij}} \right] \end{aligned}$$

and

$$\begin{aligned} \frac{\partial \text{Tr}[\mathbf{Z}^\top \mathbf{A}^\top \mathbf{K}_{DD} \mathbf{A} \mathbf{Z}]}{\partial \mathbf{d}_{ij}} &= \text{Tr} \left[\left(\frac{\partial \text{Tr}[\mathbf{Z}^\top \mathbf{A}^\top \mathbf{K}_{DD} \mathbf{A} \mathbf{Z}]}{\partial \mathbf{K}_{DD}} \right)^\top \cdot \frac{\partial \mathbf{K}_{DD}}{\partial \mathbf{d}_{ij}} \right] \\ &= \text{Tr} \left[\mathbf{A} \mathbf{Z} \mathbf{Z}^\top \mathbf{A}^\top \cdot \frac{\partial \mathbf{K}_{DD}}{\partial \mathbf{d}_{ij}} \right]. \end{aligned}$$

The two partial derivatives of kernel matrices with respect to the current atom are sparse matrices with non-zero columns or rows only where their index is equal to the index of the current atom. The two matrices are computed as follows:

$$\frac{\partial \mathbf{K}_{YD}}{\partial \mathbf{d}_{ij}} = \begin{bmatrix} 0 & \dots & \frac{\partial k(\mathbf{y}_1, \mathbf{d}_j)}{\partial \mathbf{d}_{ij}} & \dots & 0 \\ \vdots & & \frac{\partial k(\mathbf{y}_2, \mathbf{d}_j)}{\partial \mathbf{d}_{ij}} & & \vdots \\ \vdots & & \vdots & & \vdots \\ 0 & \dots & \frac{\partial k(\mathbf{y}_N, \mathbf{d}_j)}{\partial \mathbf{d}_{ij}} & \dots & 0 \end{bmatrix} \quad (6)$$

and

$$\frac{\partial \mathbf{K}_{DD}}{\partial \mathbf{d}_{ij}} = \begin{bmatrix} 0 & \dots & \frac{\partial k(\mathbf{d}_1, \mathbf{d}_j)}{\partial \mathbf{d}_{ij}} & \dots & 0 \\ 0 & \dots & \frac{\partial k(\mathbf{d}_2, \mathbf{d}_j)}{\partial \mathbf{d}_{ij}} & \dots & 0 \\ \vdots & & \vdots & & \vdots \\ \frac{\partial k(\mathbf{d}_j, \mathbf{d}_1)}{\partial \mathbf{d}_{ij}} & \dots & \frac{\partial k(\mathbf{d}_j, \mathbf{d}_j)}{\partial \mathbf{d}_{ij}} & \dots & \frac{\partial k(\mathbf{d}_j, \mathbf{d}_m)}{\partial \mathbf{d}_{ij}} \\ \vdots & & \vdots & & \vdots \\ 0 & \dots & \frac{\partial k(\mathbf{d}_n, \mathbf{d}_j)}{\partial \mathbf{d}_{ij}} & \dots & 0 \end{bmatrix}. \quad (7)$$

Depending on the kernel function of interest, the partial derivatives are computed via

$$\frac{\partial k(\mathbf{x}, \mathbf{y})}{\partial \mathbf{x}} = -\exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{2\sigma^2}\right) \frac{(\mathbf{x} - \mathbf{y})}{\sigma^2}$$

for the radial basis function kernel $k(\mathbf{x}, \mathbf{y}) = \exp(-\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{\sigma^2})$, and

$$\frac{\partial k(\mathbf{x}, \mathbf{y})}{\partial \mathbf{x}} = \beta (\mathbf{x}^\top \mathbf{y} + \alpha)^{\beta-1} \mathbf{y}$$

for the polynomial kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} + \alpha)^\beta$.

Since explicit solutions to the optimization problem on \mathbf{d}_j do not seem to exist, we update each column \mathbf{d}_j through gradient descent procedure applied on each element. At gradient descent iteration ℓ , the next element value is computed via

$$\mathbf{d}_{ij}^{(\ell+1)} = \mathbf{d}_{ij}^{(\ell)} - \gamma \mathbf{g}_{ij}^{(\ell)},$$

where $\gamma \in \mathbb{R}_+$ represents the chosen learning rate and $\mathbf{g}_{ij}^{(\ell)}$ is the gradient

$$\mathbf{g}_{ij}^{(\ell)} = \text{Tr} \left[\mathbf{A} \mathbf{Z} \left(\mathbf{Z}^\top \mathbf{A}^\top \cdot \frac{\partial \mathbf{K}_{DD}}{\partial \mathbf{d}_{ij}} - 2 \cdot \frac{\partial \mathbf{K}_{YD}}{\partial \mathbf{d}_{ij}} \right) \right],$$

where the matrices (6) and (7) are computed for the dictionary at iteration ℓ . The update step of the ORKDL-D algorithm is obtained by introducing a few steps of gradient descent as described above for updating the dictionary \mathbf{D} , after the update of \mathbf{A} described by Algorithm 1.

C. Mixed Optimized Reduced Kernel Dictionary Learning

The previous subsection presents an update procedure for the kernel vectors \mathbf{D} . A possible drawback is that direct optimization of the objective of (3) may lead to a dictionary \mathbf{D} with weaker representation power for the entire set of signals, \mathbf{Y} . We present here a further improvement. To keep the representation significance of \mathbf{D} , which is a natural property to require in the context of kernel methods, we introduce the standard DL objective (1) into the ORKDL-D problem, obtaining

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{Z}, \mathbf{D}, \mathbf{X}} \quad & \|\varphi(\mathbf{Y}) - \varphi(\mathbf{D}) \mathbf{A} \mathbf{Z}\|_F^2 + \lambda \|\mathbf{Y} - \mathbf{D} \mathbf{X}\|_F^2 \\ \text{s.t.} \quad & \|\mathbf{z}_\ell\|_0 \leq s_z, \ell = 1 : N \\ & \|\varphi(\mathbf{D}) \mathbf{a}_j\| = 1, j = 1 : n \\ & \|\mathbf{x}_\ell\|_0 \leq s_x, \ell = 1 : N \\ & \|\mathbf{d}_j\| = 1, j = 1 : n, \end{aligned} \quad (8)$$

where $\lambda \in \mathbb{R}_+$ is a penalty constant. Since a mixed objective is proposed, we name this method Mixed Optimized Reduced Kernel Dictionary Learning (MORKDL-D).

By following two directions of optimization, we update the nonlinear representation while conserving the representation power of the linear space. The current optimization problem is solved similarly with the previous problems. The update of a column of \mathbf{D} is slightly different. The gradient of the current atom uses also the partial derivative of the linear term and is

$$\tilde{\mathbf{g}}_{ij} = \mathbf{g}_{ij} - \mathbf{e}_{ij}$$

where \mathbf{e}_{ij} represents the i th element of the vector $(\mathbf{Y} - \mathbf{D}\mathbf{X}) \cdot \mathbf{X}_j^\top$; here, \mathbf{X}_j^\top is the j th row of \mathbf{X} . On the other hand, the linear sparse representation matrix \mathbf{X} is computed with the OMP algorithm. As before, the nonlinear sparse representation matrix \mathbf{Z} is computed using the Kernel OMP algorithm, since \mathbf{Z} and \mathbf{X} can be optimized independently when \mathbf{A} and \mathbf{D} are fixed.

IV. EXPERIMENTS

In this section we present the main results obtained with the proposed methods. We used three different data sets: Digits, MNIST and CIFAR-10. For each experiment we trained a kernel dictionary for representing the whole data set or a subset extracted by its corresponding label. For example, for the Digits data set we used all the signals during the training procedure, while for the MNIST and CIFAR-10 data sets we used signals specific to a selected label. For the MNIST data set we used label 5, while for the CIFAR-10 data set we used the first label.

All the algorithms were implemented in Matlab and Python and were run on a Desktop PC with Ubuntu 20.04 as operating system. The PC has a processor with 16 cores, with a base frequency of 2.90 GHz (Max Turbo Frequency 4.80 GHz), and 80 GB RAM memory. During the experiments, we computed the objective function of (3) (we report the values $\|\varphi(\mathbf{Y}) - \varphi(\mathbf{D})\mathbf{AZ}\|_F / \sqrt{mN}$ of the error per signal element) at each iteration and measured the overall execution time. Note that for MORKDL-D we report the values of the same error, not of the objective of (8). Of course, for KDL we compute the objective of (2), which corresponds to the standard case $\mathbf{D} = \mathbf{Y}$. The error and execution time were computed by the mean on ten different rounds. As kernel function we used the radial basis function $k(\mathbf{x}, \mathbf{y}) = \exp(-\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{\sigma^2})$. The kernel parameters were chosen according to a grid search, and for all data sets we chose $\sigma = 10$. All the algorithms trained a kernel dictionary \mathbf{A} of size $p = 20$, with sparsity level $s_z = 4$, using 10 iterations. For the reduced methods we initially trained a dictionary \mathbf{D} of size $n = 50$, having a sparsity level of $s_x = 5$, with 10 iterations of the AK-SVD algorithm. The dictionary \mathbf{D} dimensions ensures a reduced kernel by having a smaller number of atoms compared to the number of data signals. For example, the Digits data set contains $N = 5000$ signals of size $m = 784$, while the MNIST data set consists of approximately $N = 6000$ signals of the same size, of the same class. From the CIFAR-10 data set we used $N = 5000$ signals, specific to the interest class, of size $m = 1024$.

In the ORKDL-D and MORKDL-D algorithms, the dictionary \mathbf{D} is updated with three gradient descent iterations.

For each of the three data sets we chose a learning rate that ensures a smooth decrease of the objective function. We took $\gamma = 5 \cdot 10^{-4}$ for the Digits and MNIST data sets and $\gamma = 6 \cdot 10^{-4}$ for the CIFAR-10 data set. For all our experiments we used $\lambda = 1$.

For a better visualization of the results, we show the evolution of the nonlinear error (objective function of (3)) during the training procedure for each algorithm using the three data sets Digits (Figure 1), MNIST (Figure 2) and CIFAR-10 (Figure 3). As we can see, the improvements are visible from iteration two of training. More results are given in two tables, containing the value of the error at the last iteration (Table I) and the execution time (Table II). All the reduced methods achieve a smaller error compared to the KDL method. At first glance we notice that very large kernel spaces are not necessary to obtain good results. The introduction of a small dimensional space is enough to obtain satisfactory results. For situations where an improvement of the nonlinear representation space is needed, several iterations of gradient descent can be run as needed. According to the results obtained on the three data sets it can be seen that the biggest advantage of the proposed methods is the reduction of the execution time. For example, the non-optimized reduced form obtains an execution time at least six times shorter than the KDL method. In the problems where we try to train the kernel space, an additional time will be introduced (for dictionary \mathbf{D} update), but the execution times are still shorter than those of KDL. For reduced kernel methods the execution time is reduced by approximately 20-45% with respect to the KDL problem. All the implementations are available at <https://github.com/denisilie94/rkdl>.

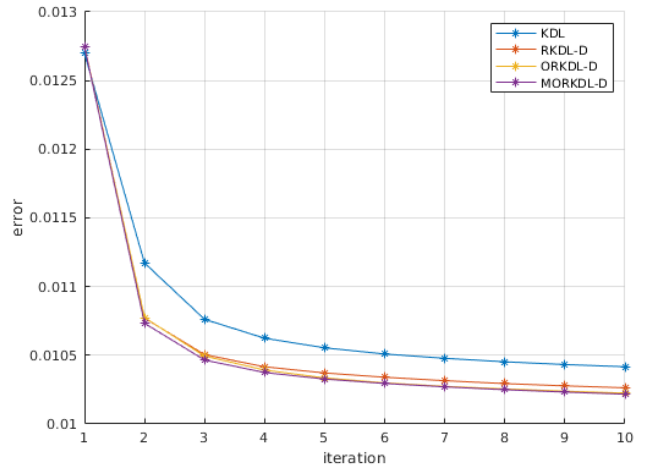


Fig. 1. Digits - representation error per iteration

V. CONCLUSIONS

In this paper we have presented a new approach to the Kernel Dictionary Learning problem by introducing a reduced kernel and thus obtaining a new algorithm, named RKDL. This method is suitable for problems with large data sets, where standard KDL requires large memory sizes and long execution

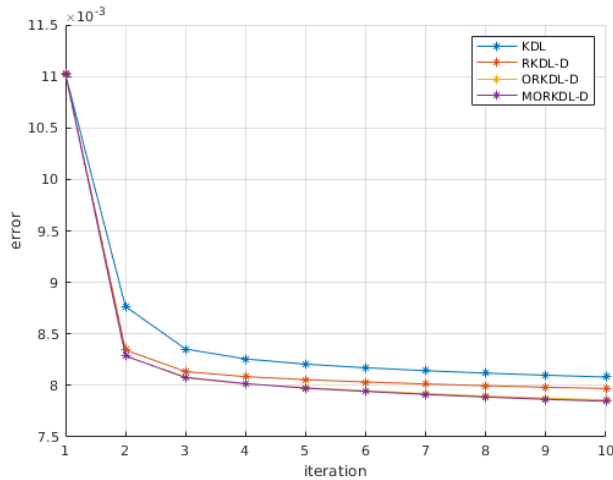


Fig. 2. MNIST - representation error per iteration

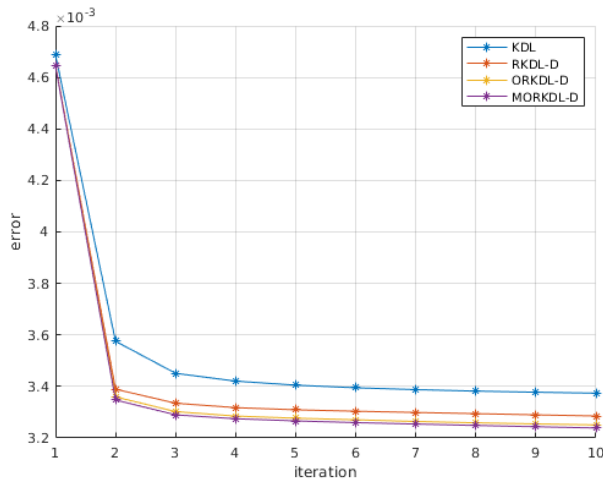


Fig. 3. CIFAR-10 - representation error per iteration

TABLE I
LAST ERROR VALUE

Algorithm \ Data set	Digits	MNIST	CIFAR-10
KDL	$1.041 \cdot 10^{-2}$	$8.081 \cdot 10^{-3}$	$3.373 \cdot 10^{-3}$
RKDL-D	$1.026 \cdot 10^{-2}$	$7.970 \cdot 10^{-3}$	$3.285 \cdot 10^{-3}$
ORKDL-D	$1.022 \cdot 10^{-2}$	$7.859 \cdot 10^{-3}$	$3.250 \cdot 10^{-3}$
MORKDL-D	$1.021 \cdot 10^{-2}$	$7.846 \cdot 10^{-3}$	$3.238 \cdot 10^{-3}$

TABLE II
EXECUTION TIME IN SECONDS

Algorithm \ Data set	Digits	MNIST	CIFAR-10
KDL	61.7	92.5	65.1
RKDL-D	9.2	10.9	9.6
ORKDL-D	39.2	48.4	51.6
MORKDL-D	38.9	49	51.4

times. Moreover, we have demonstrated that most of the time a large representation space for the kernel matrix is not always needed to obtain satisfactory results. The reduced form of the kernel is enough to get similar results or even better. Since the kernel matrix is smaller, the required execution time is also much shorter.

The RKDL algorithm was presented under three different forms: standard RKDL-D, optimized RKDL-D and mixed optimized RKDL-D, the latter two including optimization of the kernel vectors. All of them obtain competitive results with the standard KDL problem.

REFERENCES

- [1] H. Van Nguyen, V.M. Patel, N.M. Nasrabadi, and R. Chellappa. Design of non-linear kernel dictionaries for object recognition. *IEEE Transactions on Image Processing*, 22(12):5123–5135, 2013.
- [2] J.J. Thiagarajan, K.N. Ramamurthy, and A. Spanias. Multiple kernel sparse representations for supervised and unsupervised learning. *IEEE Transactions on Image Processing*, 23(7):2905–2915, 2014.
- [3] A. Golts and M. Elad. Linearized kernel dictionary learning. *IEEE Journal of Selected Topics in Signal Processing*, 10(4):726–739, 2016.
- [4] S. Kumar, M. Mohri, and A. Talwalkar. Sampling methods for the Nyström method. *The Journal of Machine Learning Research*, 13(1):981–1006, 2012.
- [5] A. Rahimi and B. Recht. Random features for large-scale kernel machines. *Advances in neural information processing systems*, 20, 2007.
- [6] Y.J. Lee and O.L. Mangasarian. RSVM: Reduced support vector machines. In *Proceedings of the 2001 SIAM International Conference on Data Mining*, pages 1–17. SIAM, 2001.
- [7] K.M. Lin and C.J. Lin. A study on reduced support vector machines. *IEEE Transactions on Neural Networks*, 14(6):1449–1459, 2003.
- [8] G.M. Fung, O.L. Mangasarian, and A.J. Smola. Minimal kernel classifiers. *Journal of Machine Learning Research*, 3(Nov):303–321, 2002.
- [9] W. Deng, Q. Zheng, and K. Zhang. Reduced kernel extreme learning machine. In *Proceedings of the 8th international conference on computer recognition systems CORES 2013*, pages 63–69. Springer, 2013.
- [10] J. Hu and Y.P. Tan. Nonlinear dictionary learning with application to image classification. *Pattern Recognition*, 75:282–291, 2018.
- [11] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [12] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [13] Y.C. Pati, R. Rezaifar, and P.S. Krishnaprasad. Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition. In *Proceedings of 27th Asilomar conference on signals, systems and computers*, pages 40–44. IEEE, 1993.
- [14] B. Dumitrescu and P. Irofti. *Dictionary learning algorithms and applications*. Springer, 2018.
- [15] M. Aharon, M. Elad, and A. Bruckstein. K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, 54(11):4311–4322, 2006.
- [16] R. Rubinstein, M. Zibulevsky, and M. Elad. Efficient implementation of the K-SVD algorithm using batch orthogonal matching pursuit. Technical report, Computer Science Department, Technion, 2008.
- [17] D.C. Ilie-Ablachim and B. Dumitrescu. Anomaly detection with selective dictionary learning. In *International Conference on Control, Decision and Information Technologies*. IEEE, 2022.