

# Task-Driven Dictionary Learning based on Convolutional Neural Network Features

Tom Tirer

*School of Electrical Engineering*  
*Tel Aviv University*  
 Tel Aviv, Israel  
 tirer.tom@gmail.com

Raja Giryes

*School of Electrical Engineering*  
*Tel Aviv University*  
 Tel Aviv, Israel  
 raja@tauex.tau.ac.il

**Abstract**—Modeling data as a linear combination of a few elements from a learned dictionary has been used extensively in the recent decade in many fields, such as machine learning and signal processing. The learning of the dictionary is usually performed in an unsupervised manner, which is most suitable for regression tasks. However, for other purposes, e.g. image classification, it is advantageous to learn a dictionary from the data in a supervised way. Such an approach has been referred to as *task-driven dictionary learning*. In this work, we integrate this approach with deep learning. We modify this strategy such that the dictionary is learned for features obtained by a convolutional neural network (CNN). The parameters of the CNN are learned simultaneously with the task-driven dictionary and with the classifier parameters.

## I. INTRODUCTION

In the recent decade many signal and image processing applications have benefited remarkably from adopting low dimensional models for the signals of interest. The most popular form of this low dimensional modeling assumes that signals admit a sparse representation in a given dictionary - this is referred to as the "sparsity model" [1]. A crucial factor in the success of sparse representation methods is the choice of the dictionary, which should provide a faithful encoding for a given signal, but also be able to discriminatively represent different signals.

Using off-the-shelf dictionaries, such as redundant wavelets or discrete cosine transform, may suffice for restoration problems of general images, but will not be effective enough for specific tasks, such as classification of images from a certain domain. Indeed, the dictionary learning (DL) approach, in which the dictionary is learned from the training data has led to state-of-the-art results in many practical applications, such as image reconstruction and image classification.

This work is supported by ERC StG SPADE.

There are two main strategies for dictionary learning: unsupervised DL and supervised DL. Unsupervised DL methods, which conveniently spare any need for data labeling, have been widely applied to image processing tasks, such as image denoising and compression [2], [3], [4]. Two popular unsupervised DL algorithms are K-SVD [5] and online DL [6]. However, without using additional label information of training data, unsupervised DL methods may not acquire discriminative representations for signals that belong to different classes. Therefore, while these methods are powerful for data reconstruction, they may not be a good choice for classification tasks. In contrast, supervised DL do use available class labels of training samples. This class discrimination information is exploited while learning the dictionary, and thus leads to better classification performance [7], [8].

In this work we build on the task-driven dictionary learning framework proposed in [7]. We use this framework for multiclass classification. However, we modify the original algorithm such that the dictionary is learned for features obtained by a convolutional neural network (CNN) [9]. The parameters of the CNN are learned simultaneously with the task-driven dictionary and with the classifier parameters. The proposed method outperforms the original algorithm that learns a linear transform of the input data. In addition, in the case of a small number of training examples, it can get better performance than the full CNN, from which we extract the features.

Our effort is related to other recent works that combines neural network training with low parsimonious data models. For example, in [10] it has been shown that a CNN may be interpreted as a convolutional sparse coding process, which has led to theoretical guarantees for the uniqueness of outputs of each layer of the

network. In [11] a novel deep network has been proposed that encodes features of hidden layers with nonnegative sparse representation. In [12] the network weights were decomposed as a sparse composition of a fixed set of filters. In [13], low-rank based loss has been proposed to improve the performance of neural networks.

## II. BACKGROUND

For the sake of completeness, we briefly describe the unsupervised DL and the (supervised) task-driven DL. Classical unsupervised dictionary learning techniques consider a finite training set of signals  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$  and minimize the empirical cost function

$$\min_{\mathbf{D}} \frac{1}{n} \sum_{i=1}^n \ell_u(\mathbf{x}_i, \mathbf{D}) \quad (1)$$

with respect to a dictionary  $\mathbf{D} \in \mathbb{R}^{m \times p}$ , where each column of  $\mathbf{D}$  represents a dictionary element.  $\ell_u$  (the subscript  $u$  denotes "unsupervised") is a loss-function such that  $\ell_u(\mathbf{x}, \mathbf{D})$  should be small if  $\mathbf{D}$  is "good" at representing the signal  $\mathbf{x}$  in a sparse fashion. Alternatively, it is more likely that one is interested in the minimization of the expected cost

$$\min_{\mathbf{D}} \mathbb{E}_{\mathbf{x}} \{\ell_u(\mathbf{x}, \mathbf{D})\}, \quad (2)$$

where the expectation is taken relative to the (unknown) probability distribution  $p(\mathbf{x})$  of the data. This kind of cost functions is usually minimized using stochastic gradient descent (SGD) methods.

We define  $\ell_u(\mathbf{x}, \mathbf{D})$  as the optimal value of the elastic-net sparse coding problem [14]

$$\ell_u(\mathbf{x}, \mathbf{D}) \triangleq \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda_1 \|\boldsymbol{\alpha}\|_1 + \frac{\lambda_2}{2} \|\boldsymbol{\alpha}\|_2^2, \quad (3)$$

where  $\lambda_1$  and  $\lambda_2$  are regularization parameters. Setting  $\lambda_2 = 0$  yields the  $\ell_1$  sparse decomposition problem (also known as basis pursuit denoising [15] or Lasso [16]), while using  $\lambda_2 > 0$  results in a more stable solution. To prevent the  $\ell_2$ -norm of  $\mathbf{D}$  from being arbitrarily large, which would lead to arbitrarily small values of  $\boldsymbol{\alpha}$ , it is common to constrain  $\mathbf{D}$  to the following convex set of matrices

$$\mathcal{D} \triangleq \{\mathbf{D} \in \mathbb{R}^{m \times p} \text{ s.t. } \forall j \in \{1, \dots, p\} \|\mathbf{d}_j\|_2 \leq 1\}. \quad (4)$$

where  $\{\mathbf{d}_j\}$  denote the columns of  $\mathbf{D}$ .

We turn now to describe the task-driven DL formulation. Given a dictionary  $\mathbf{D}$ , a vector  $\mathbf{x}$  can be represented as a sparse vector  $\boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D})$  defined as the solution of (3). We assume that each signal  $\mathbf{x} \in \mathcal{X}$  is associated

with a variable  $\mathbf{y} \in \mathcal{Y}$ , which we want to predict from  $\mathbf{x}$  (e.g., a label in classification tasks). We can now use the sparse vector  $\boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D})$  as a feature representation of a signal  $\mathbf{x}$  in the following minimization formulation:

$$\min_{\mathbf{W} \in \mathcal{W}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \{\ell_s(\mathbf{y}, \mathbf{W}, \boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D}))\} + \frac{\nu}{2} \|\mathbf{W}\|_F^2, \quad (5)$$

where  $\mathbf{W}$  are model parameters which we want to learn,  $\mathcal{W}$  is a convex set, and  $\nu$  is a regularization parameter. In this equation,  $\ell_s$  (the subscript  $s$  denotes "supervised") is a convex loss function that measures how well one can predict  $\mathbf{y}$  by observing  $\boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D})$  given the model parameters  $\mathbf{W}$ . For instance, it can be the square, logistic, or hinge loss. The expectation is taken with respect to the unknown probability distribution  $p(\mathbf{x}, \mathbf{y})$  of the data.

Finally, the task-driven dictionary learning formulation consists of jointly learning  $\mathbf{W}$  and  $\mathbf{D}$  by solving

$$\min_{\mathbf{D} \in \mathcal{D}, \mathbf{W} \in \mathcal{W}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \{\ell_s(\mathbf{y}, \mathbf{W}, \boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D}))\} + \frac{\nu}{2} \|\mathbf{W}\|_F^2. \quad (6)$$

Note that the function in the expectation is non-convex and non-smooth. The non-convexity is typical to DL problems, and handled by alternating minimization: performing sparse coding with a fixed dictionary, and updating the dictionary using the obtained representations. The non-smoothness is mitigated by the expectation operator together with properties of  $\boldsymbol{\alpha}^*(\mathbf{x}_i, \mathbf{D})$ , obtained by the elastic-net formulation with  $\lambda_2 > 0$ , as detailed in [7].

The formulation (6) is further extended in [7] by including a linear transform of the input data, represented by a matrix  $\mathbf{Z}$ , which allows to increase or decrease the dimension of the features  $\mathbf{x}$ . Thus, we have

$$\min_{\mathbf{D} \in \mathcal{D}, \mathbf{W} \in \mathcal{W}, \mathbf{Z} \in \mathcal{Z}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \{\ell_s(\mathbf{y}, \mathbf{W}, \boldsymbol{\alpha}^*(\mathbf{Z}\mathbf{x}, \mathbf{D}))\} + \frac{\nu_1}{2} \|\mathbf{W}\|_F^2 + \frac{\nu_2}{2} \|\mathbf{Z}\|_F^2. \quad (7)$$

where  $\nu_1$  and  $\nu_2$  are two regularization parameters, and  $\mathcal{Z}$  is a convex set.

The optimization of (7) can be carried out using the projected stochastic gradient descent scheme, presented in Algorithm 1.

## III. THE PROPOSED METHOD

Inspired by the excellent performance of recent deep-learning-based methods at different tasks [17], [18], [19], [20], we turn to extend the task-driven DL such that the dictionary is learned for features obtained by a convolutional neural network (CNN) [9], instead of a simple linear transform of the input data. In more details, instead

**Algorithm 1:** SGD algorithm for task-driven DL

---

**Input:**  $p(\mathbf{x}, \mathbf{y})$  (a way to draw i.i.d samples),  
 $\lambda_1, \lambda_2, \nu_1, \nu_2$  (regularization parameters),  
 $\mathbf{D} \in \mathcal{D}, \mathbf{W} \in \mathcal{W}, \mathbf{Z} \in \mathcal{Z}$  (initializations),  $T$   
(number of iterations),  $t_0, \rho$  (learning rate  
parameters).  
**for**  $t = 1$  **to**  $T$  **do**  
  Draw  $(\mathbf{x}, \mathbf{y})$  from  $p(\mathbf{x}, \mathbf{y})$ ;  
  Sparse coding: compute  $\alpha^*(\mathbf{Z}\mathbf{x}, \mathbf{D})$  using  
  elastic-net;  
   $\Lambda = \{j \in \{1, \dots, p\} : \alpha^*[j] \neq 0\}$ ;  
   $\beta_\Lambda^* = (\mathbf{D}_\Lambda^T \mathbf{D}_\Lambda + \lambda_2 \mathbf{I})^{-1} \nabla_{\alpha_\Lambda} \ell_s(\mathbf{y}, \mathbf{W}, \alpha^*)$ ;  
   $\beta_{\Lambda^c}^* = \mathbf{0}$ ;  
   $\rho_t = \min(\rho, \rho \frac{t_0}{t})$ ;  
   $\mathbf{W} = \Pi_{\mathcal{W}}[\mathbf{W}$   
   $\quad - \rho_t (\nabla_{\mathbf{W}} \ell_s(\mathbf{y}, \mathbf{W}, \alpha^*) + \nu_1 \mathbf{W})]$ ;  
   $\mathbf{Z} = \Pi_{\mathcal{Z}}[\mathbf{Z} - \rho_t (\mathbf{D} \beta^* \mathbf{x}^T + \nu_2 \mathbf{Z})]$ ;  
   $\mathbf{D} = \Pi_{\mathcal{D}}[\mathbf{D}$   
   $\quad - \rho_t (-\mathbf{D} \beta^* \alpha^{*T} + (\mathbf{Z}\mathbf{x} - \mathbf{D} \alpha^*) \beta^{*T})]$ ;  
**end**  
**Output:**  $\mathbf{D} \in \mathcal{D}, \mathbf{W} \in \mathcal{W}, \mathbf{Z} \in \mathcal{Z}$ .

---

of training the dictionary for sparse representations of the linearly transformed  $\mathbf{x}$ :

$$\tilde{\mathbf{x}}_{linear} \triangleq \mathbf{Z}\mathbf{x}, \quad (8)$$

we train it for sparse representations of the features

$$\tilde{\mathbf{x}}_{CNN} \triangleq \mathbf{f}_{CNN}(\mathbf{x}; \mathbf{z}), \quad (9)$$

where  $\mathbf{f}_{CNN}$  is a non-linear function of  $\mathbf{x}$ , obtained by common building blocks of CNNs, such as convolutional layers and max-pooling layers.  $\mathbf{z}$  is an  $N_z \times 1$  vector that denotes all the parameters of the network, e.g. the weights of the convolutional layers. The proposed architecture is demonstrated in Fig. 1. Note that the parameters of the CNN layers are not fixed. They are learned simultaneously with the task-driven dictionary  $\mathbf{D}$  and with the model parameters  $\mathbf{W}$ . Thus, the formulation in (7) is extended to

$$\min_{\mathbf{D} \in \mathcal{D}, \mathbf{W} \in \mathcal{W}, \mathbf{z} \in \mathcal{Z}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left\{ \ell_s(\mathbf{y}, \mathbf{W}, \alpha^*(\tilde{\mathbf{x}}(\mathbf{x}; \mathbf{z}), \mathbf{D})) \right\} + \frac{\nu_1}{2} \|\mathbf{W}\|_F^2 + \frac{\nu_2}{2} \|\mathbf{z}\|_2^2. \quad (10)$$

Note the general dependency of  $\alpha^*$  on  $\mathbf{z}$  through  $\tilde{\mathbf{x}}$ , which equals to  $\tilde{\mathbf{x}}_{linear}$  when linear transformation is used and to  $\tilde{\mathbf{x}}_{CNN}$  when a CNN is employed.

We restrict ourselves to standard static neural networks, which have bounded results given that their inputs are bounded [21]. Therefore, the function  $\alpha^*(\tilde{\mathbf{x}}_{CNN}, \mathbf{D})$

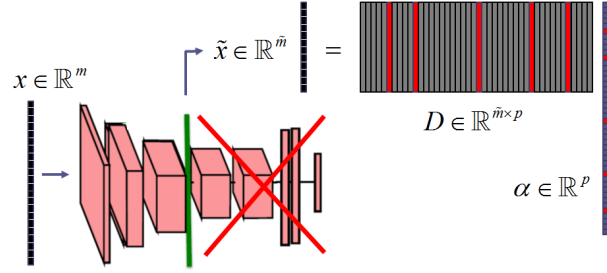


Fig. 1: Proposed architecture for joint dictionary and CNN learning.

is still uniformly Lipschitz, and the differentiability proofs in [7] are still valid. However, the update rule for  $\mathbf{z}$  has to be modified to our general case

$$\mathbf{z} = \Pi_{\mathcal{Z}}[\mathbf{z} - \rho_t (\nabla_{\mathbf{z}} \ell_s + \nu_2 \mathbf{z})]. \quad (11)$$

Following the computation technique of [22], we have

$$\nabla_{\mathbf{z}} \ell_s = \frac{\partial \tilde{\mathbf{x}}}{\partial \mathbf{z}} \mathbf{D} \beta^*, \quad (12)$$

where we used the chain rule, and  $\beta^*$  is defined exactly as in Algorithm 1. As a sanity check, in the case  $\tilde{\mathbf{x}} = \mathbf{Z}\mathbf{x}$  (where  $\mathbf{Z} \in \mathbb{R}^{\tilde{m} \times m}$ ), we have

$$\frac{\partial(\mathbf{Z}\mathbf{x})}{\partial \text{vec}(\mathbf{Z})} = \mathbf{x} \otimes \mathbf{I}_{\tilde{m}}, \quad (13)$$

where  $\otimes$  is the Kronecker product and  $\mathbf{I}_{\tilde{m}}$  is the identity matrix of size  $\tilde{m} \times \tilde{m}$ . Therefore, we have

$$\nabla_{\mathbf{z}} \ell_s = (\mathbf{x} \otimes \mathbf{I}_{\tilde{m}}) \mathbf{D} \beta^* = (\mathbf{x} \otimes \mathbf{D}) \beta^* = \text{vec}(\mathbf{D} \beta^* \mathbf{x}^T), \quad (14)$$

which is the vectorization of  $\partial \ell_s / \partial \mathbf{Z} = \mathbf{D} \beta^* \mathbf{x}^T$  that appears in [7]. Exploiting the fact that the optimization algorithm is based on stochastic gradient descent, we use backpropagation [23] to calculate  $\nabla_{\mathbf{z}} \ell_s$ , with the vector  $\mathbf{D} \beta^*$  being used as the derivatives of the network w.r.t. the feature layer's output.

#### IV. APPLICATION TO DIGITS CLASSIFICATION

We consider here a classification task using the MNIST handwritten digits dataset [24], which contains 70,000  $28 \times 28$  images, 60,000 of them for training and 10,000 for testing.

We use a single dictionary and a multiclass loss function. This strategy has less computational cost compared to other strategies such as one-versus-all and one-versus-one. Thus, in our setting,  $\mathcal{Y} = \{1, \dots, 10\}$ . Given a vector  $\mathbf{x}$ , we want to learn (jointly with  $\mathbf{D}$  and  $\mathbf{z}$ ) the parameters  $\mathbf{W} \in \mathbb{R}^{p \times 10}$  of a linear model to predict  $\mathbf{y} \in \mathcal{Y}$ , using the sparse representation  $\alpha^*(\tilde{\mathbf{x}}, \mathbf{D})$  as

features. Each column of  $\mathbf{W}$  is associated with a specific class, i.e.  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_{10}]$ , and the maximal entry of  $\mathbf{W}^T \boldsymbol{\alpha}^*$  determines the predicted class. The loss function that we use is the categorical cross-entropy loss (also known as softmax loss), which is given by

$$\ell_s(y, \mathbf{W}, \boldsymbol{\alpha}^*) = -\log \left( \frac{e^{w_y^T \boldsymbol{\alpha}^*}}{\sum_{\tilde{y} \in \mathcal{Y}} e^{w_{\tilde{y}}^T \boldsymbol{\alpha}^*}} \right). \quad (15)$$

Note that the derivatives required for Algorithm 1 can be computed easily. Let

$$P_y \triangleq \frac{e^{w_y^T \boldsymbol{\alpha}^*}}{\sum_{\tilde{y} \in \mathcal{Y}} e^{w_{\tilde{y}}^T \boldsymbol{\alpha}^*}}, \quad (16)$$

then, it can be shown that

$$\frac{\partial \ell_s(y, \mathbf{W}, \boldsymbol{\alpha}^*)}{\partial \alpha^*[i]} = \sum_{\tilde{y} \in \mathcal{Y}} (-w_y[i] + w_{\tilde{y}}[i]) P_{\tilde{y}}, \quad (17)$$

$$\frac{\partial \ell_s(y, \mathbf{W}, \boldsymbol{\alpha}^*)}{\partial w_{\tilde{y}}^*[i]} = \alpha^*[i] (P_{\tilde{y}} - \delta[\tilde{y} - y]), \quad (18)$$

where  $\delta[\cdot]$  denotes the Kronecker delta.

We compare the performance of three methods:

- full-CNN: an 8-layers CNN from [25].
- TDDL-LIN: task-driven DL with  $\tilde{\mathbf{x}} = \mathbf{Z}\mathbf{x}$ .
- TDDL-CNN: task-driven with  $\tilde{\mathbf{x}} = \mathbf{f}_{CNN}(\mathbf{x}; \mathbf{z})$ .

For the full-CNN, we use the exact 8-layers CNN architecture and training scheme (including the optimization algorithm) supplied in the MatConvNet toolbox [25] for the MNIST dataset.

For the task-driven methods, similar to [7], each image is normalized to zero mean and unit norm. Also, it is worthwhile to decrease the dimension of the input vectors in order to reduce the runtime. For TDDL-LIN we use a  $196 \times 784$  matrix  $\mathbf{Z}$ , which is initialized with entries independently drawn from the standard normal distribution  $\mathcal{N}(0, 1)$ . For TDDL-CNN we keep only the first 4 layers of the full CNN, and reduce the number of filters in the third layer. The first layer is a convolutional layer with 20 filters of size  $5 \times 5 \times 1$  and stride 1. The second layer is a  $2 \times 2$  max-pooling with stride 2. The third layer is a convolutional layer with 12 filters of size  $5 \times 5 \times 20$  and stride 1. The fourth layer is a  $2 \times 2$  max-pooling with stride 2. For an input image of size  $28 \times 28$ , this results with a  $192 \times 1$  feature vector. We initialize the weights with standard normal weights, and the biases with zeros. Note that the linear transform of TDDL-LIN has 153,664 learnable parameters, while the non-linear transform of TDDL-CNN has only 6,532 learnable parameters.

TABLE I: Accuracy rates of different classification methods

Train set size	full-CNN	TDDL-LIN	TDDL-CNN
60K	<b>0.990</b>	0.971	0.984
1K	0.917	0.898	<b>0.928</b>

For both TDDL-based methods, we use an  $\tilde{m} \times 400$  dictionary  $\mathbf{D}$ , which is initialized in the same manner. From each class, we take 40 input features (i.e. 40  $\tilde{\mathbf{x}}$  vectors, using the fixed initial value of  $\mathbf{z}$  or  $\mathbf{Z}$ ) to initialize an online dictionary learning (ODL) algorithm applied only on training samples from the same class. The resulted 40 atoms from each class are then collected into  $\mathbf{D}$ . This strategy has led to better performance and more discriminative sparse representations. With this initial dictionary  $\mathbf{D}$  at hand, we obtain an initial value of  $\mathbf{W}$  by optimizing (5) (again,  $\mathbf{z}$  is fixed), and get a triplet of parameters  $(\mathbf{D}, \mathbf{W}, \mathbf{z})$ , which are used to initialize the task-driven DL algorithm.

We use  $\lambda_1 = 0.8$  and  $\lambda_1 = 4$  for TDDL-LIN and TDDL-CNN, respectively. These choices lead to an average sparsity of 25-35 for  $\boldsymbol{\alpha}^*$  at the end of the task-driven DL (note that after ODL alone, the sparsity is about 10, but it grows during the TDDL, as the classification error reduces, making the representation more discriminative). The rest of the parameters are  $\lambda_2 = 0$  (though the proof for differentiability [7] requires  $\lambda_2 > 0$ , the algorithm still converged with  $\lambda_2 = 0$ ),  $\nu_1 = 0.1$ ,  $\nu_2 = 0$ ,  $t_0 = 4,000$  and  $\rho = 0.05$ . The mini-batches size is 100.

We perform two experiments. In the first one, we use all the training data as is (i.e. 60,000 training samples) with  $T = 40,000$ . The second experiment reflects an extreme case, in which only 100 training samples are given per class (i.e. 1,000 training samples). The latter experiment is repeated 20 times with  $T = 1,000$ , each time for a different subset of 100 training samples. Table I summarizes the accuracy rates of the examined methods on the test set.

Fig. 2 presents the model parameters  $\mathbf{W}$  ( $w_y$  for each  $y \in \mathcal{Y}$ ) that were learned for TDDL-CNN in the first experiment. Note that the classifier emphasizes different indices of  $\boldsymbol{\alpha}^*$  (i.e., using different dictionary atoms) for each class, a fact which implies that the dictionary is discriminative.

From the results in Table I, it can be concluded that TDDL-CNN outperforms TDDL-LIN. It can be also observed that TDDL-CNN outperforms the full-CNN when the labeled data is limited. This demonstrates that it may be possible to improve CNN performance by adding a data-model prior in the learning process.

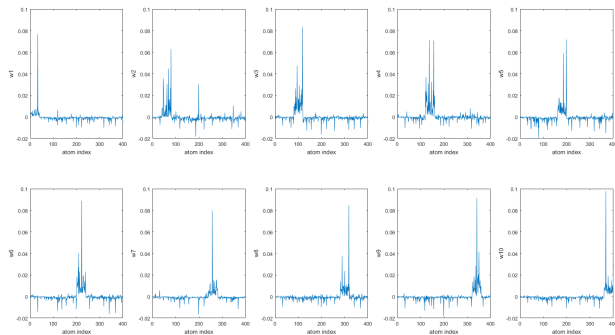


Fig. 2: The learned  $w_y, y \in \mathcal{J}$ , for TDDL-CNN (first exp.).

## V. CONCLUSION

In this work we have proposed a CNN based task-driven dictionary learning, where the parameters of the CNN are optimized simultaneously with the dictionary and a given classifier. Our method clearly outperforms task-driven dictionary learning that is restricted to a linear transform of the input data. Moreover, our work shows the potential of improving the performance of deep neural networks in the presence of a small number of training examples using the sparsity prior. It also leads to the following question: Is it possible to use the unlabeled data to improve the dictionary and further increase the performance gap between TDDL-CNN and full-CNN? One approach is using the semi-supervised formulation in [7] combined with a CNN. Yet, our efforts in this direction have led to a small improvement (e.g. the accuracy of TDDL-CNN for MNIST with 1K training samples increased from 0.928 to 0.935). Presumably, the limited improvement follows from the fact that a plain unsupervised dictionary learning makes the dictionary less discriminative. Thus, we leave the semi-supervised learning task to a future work. Such a future work may examine also other deep architectures and loss functions.

## REFERENCES

- [1] A. M. Bruckstein, D. L. Donoho, and M. Elad, "From sparse solutions of systems of equations to sparse modeling of signals and images," *SIAM review*, vol. 51, no. 1, pp. 34–81, 2009.
- [2] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Trans. Imag. Proc.*, vol. 15, no. 12, pp. 3736–3745, 2006.
- [3] O. Bryt and M. Elad, "Compression of facial images using the K-SVD algorithm," *Journal of Visual Communication and Image Representation*, vol. 19, no. 4, pp. 270–282, 2008.
- [4] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman, "Non-local sparse models for image restoration," in *ICCV*, 2009.
- [5] M. Aharon, M. Elad, and A. Bruckstein, "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Transactions on signal processing*, vol. 54, no. 11, pp. 4311–4322, 2006.
- [6] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, "Online dictionary learning for sparse coding," in *ICML*, 2009.
- [7] J. Mairal, F. Bach, and J. Ponce, "Task-driven dictionary learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 4, pp. 791–804, 2012.
- [8] Z. Jiang, Z. Lin, and L. S. Davis, "Label consistent K-SVD: Learning a discriminative dictionary for recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 11, pp. 2651–2664, 2013.
- [9] Y. LeCun, Y. Bengio, et al., "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, vol. 3361, no. 10, p. 1995, 1995.
- [10] V. Pappayan, Y. Romano, and M. Elad, "Convolutional neural networks analyzed via convolutional sparse coding," *Journal of Machine Learning Research*, vol. 18, no. 83, pp. 1–52, 2017.
- [11] X. Sun, N. M. Nasrabadi, and T. D. Tran, "Supervised deep sparse coding networks," *arXiv:1701.08349*, 2017.
- [12] Q. Qiu, X. Cheng, R. Calderbank, and G. Sapiro, "Dcfnet: Deep neural network with decomposed convolutional filters," *arXiv:1802.04145*, 2018.
- [13] J. Lezama, Q. Qiu, P. Muse, and G. Sapiro, "OLÉ: Orthogonal low-rank embedding, a plug and play geometric loss for deep learning," in *CVPR*, 2018.
- [14] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 67, no. 2, pp. 301–320, 2005.
- [15] S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM review*, vol. 43, no. 1, pp. 129–159, 2001.
- [16] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 267–288, 1996.
- [17] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012.
- [18] G. E. Dahl, D. Yu, L. Deng, and A. Acero, "Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition," *IEEE Transactions on audio, speech, and language processing*, vol. 20, no. 1, pp. 30–42, 2012.
- [19] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv:1409.1556*, 2014.
- [20] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [21] J. Bruna, A. Szlam, and Y. LeCun, "Signal recovery from pooling representations," in *ICML*, 2014.
- [22] J. Yang, K. Yu, and T. Huang, "Supervised translation-invariant sparse coding," in *CVPR*, 2010.
- [23] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by error propagation," tech. rep., California Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- [24] Y. LeCun, C. Cortes, and C. Burges, "MNIST handwritten digit database," *AT&T Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, vol. 2, 2010.
- [25] A. Vedaldi and K. Lenc, "MatConvNet: Convolutional neural networks for Matlab," in *ACM international conference on Multimedia*, pp. 689–692, 2015.