

Deep Sparse-coded Network (DSN)

Youngjune Gwon
Harvard University
gyj@eecs.harvard.edu

Miriam Cha
Harvard University
miriamcha@fas.harvard.edu

H. T. Kung
Harvard University
kung@harvard.edu

Abstract—We present Deep Sparse-coded Network (DSN), a deep architecture based on multilayer sparse coding. It has been considered difficult to learn a useful feature hierarchy by stacking sparse coding layers in a straightforward manner. The primary reason is the modeling assumption for sparse coding that takes in a dense input and yields a sparse output vector. Applying a sparse coding layer on the output of another tends to violate the modeling assumption. We overcome this shortcoming by interlacing nonlinear pooling units. Average- or max-pooled sparse codes are aggregated to form dense input vectors for the next sparse coding layer. Pooling achieves nonlinear activation analogous to neural networks while not introducing diminished gradient flows during the training. We introduce a novel backpropagation algorithm to finetune the proposed DSN beyond the pretraining via greedy layerwise sparse coding and dictionary learning. We build an experimental 4-layer DSN with the ℓ_1 -regularized LARS and the greedy- ℓ_0 OMP, and demonstrate superior performance over a similarly-configured stacked autoencoder (SAE) on CIFAR-10.

I. MOTIVATION

Representational power of single-layer feature learning is limited for tasks that involve large complex data objects such as a high-resolution image of human face. Best current practices in visual recognition use deep architectures based on autoencoder [1], restricted Boltzmann machine (RBM) [2], and convolutional neural network (CNN) [3]. A deep architecture stacks two or more layers of feature learning units in the hope of discovering hierarchical representations for data. In other words, deep architectures allow us to understand a feature at each layer using the features of the layer below. Such hierarchical decomposition is particularly useful when it is hard to resolve ambiguity of the low-level (or localized) features of data. Deep architectures also promote representational efficiency. We can achieve compaction of all characteristic features for the entire image, book, or lengthy multimedia clip to a single vector.

In an empirical analysis by Coates, Lee and Ng [4], a simple single-layer scheme, which encodes features with basis vectors learned from K-means clustering, is found on par (or sometimes, even superior) to RBM, deep neural network, and CNN for classification tasks on the CIFAR-10 and NORB datasets. We have also been able to draw a similar conclusion from our experiments with sparse coding. With these in mind, it is sound to build a deep architecture on sparse coding. Unfortunately, this takes more than just stacking sparse coding

units in layers. From sparse coding research on hierarchical feature learning [5], [6], we could deduce some explanations for the difficulty. First, sparse coding (in particular, the ℓ_1 -regularized LASSO or LARS) is computationally expensive for multilayering and associated optimizations. From our experience, it is cumbersome to simply connect multiple sparse coding units and run data as a feedforward network. Secondly, sparse coding makes an inherent assumption on the input being non-sparse. Therefore, a straightforward approach to take the output from one sparse coding unit for an input to another is flawed. Lastly, it is difficult to optimize all layers of sparse coding jointly. One consensual notion of deep learning suggests layer-by-layer unsupervised pretraining should be followed by supervised finetuning of the whole system, which is commonly done by backpropagation.

In this paper, we propose Deep Sparse-coded Network (DSN), a deep architecture for sparse coding as a principled extension from its single-layer counterpart. We consider both the ℓ_1 -regularized LASSO/LARS [7], [8] and greedy- ℓ_0 OMP [9] as a legitimate sparse coding method. Using max pooling as nonlinear activation analogous to neural networks, we avoid linear cascade of dictionaries and keep the effect of multilayering in tact. This architectural usage will remedy the problem of too many feature vectors by aggregating them to their maximum elements and help preserve translational invariance of higher-layer representations. Beyond the layer-by-layer pretraining, we propose a novel backpropagation algorithm that is specific to multilayer sparse coding interlaced by spatial max pooling. We have validated empirically that the proposed backpropagation algorithm can further optimize the classification performance by 4–6% gain in accuracy.

Rest of this paper is organized as follows. In Section 2, we provide a brief background on sparse coding. Section 3 will introduce DSN, explain its architectural principles, and discuss training algorithms. In Section 4, we present an empirical evaluation of DSN, and Section 5 concludes the paper.

II. SPARSE CODING BACKGROUND

Originally used to explain neuronal activations [10], sparse coding is an unsupervised method to learn an efficient representation of data using a small number of basis vectors. It has been used to discover high-level features present in data

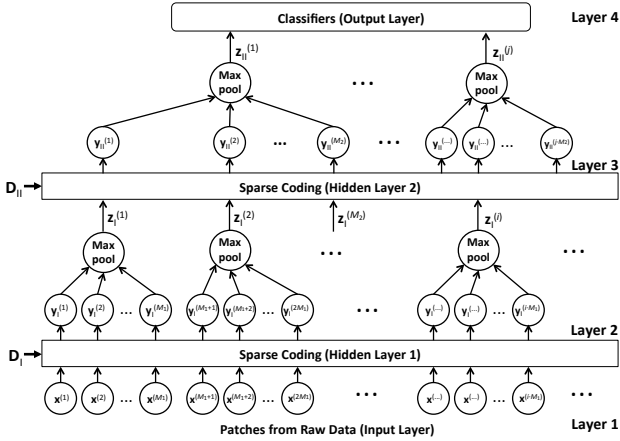


Fig. 1. Deep Sparse-coded Network (DSN) with four layers

from unlabeled examples. Given an example $\mathbf{x} \in \mathbb{R}^N$, sparse coding searches for a representation $\mathbf{y} \in \mathbb{R}^K$ (*i.e.*, the feature vector for \mathbf{x}) while simultaneously updating the dictionary $\mathbf{D} \in \mathbb{R}^{N \times K}$ of K basis vectors by

$$\min_{\mathbf{D}, \mathbf{y}} \|\mathbf{x} - \mathbf{D}\mathbf{y}\|_2^2 + \lambda \|\mathbf{y}\|_1 \quad \text{s.t.} \quad \|\mathbf{d}_i\|_2 \leq 1, \forall i \quad (1)$$

where \mathbf{d}_i is i th dictionary atom in \mathbf{D} , and λ is a regularization parameter that penalizes over the ℓ_1 -norm, which induces a sparse solution. With $K > N$, sparse coding typically trains an overcomplete dictionary.

A more direct way to control sparsity is to regularize on the ℓ_0 pseudo-norm $\|\mathbf{y}\|_0$, describing the number of nonzero elements in \mathbf{y} . However, it is known to be intractable to compute the sparsest ℓ_0 solution in general. The approach in Eq. (1) is called least absolute shrinkage and selection operator (LASSO) [7], a convex relaxation of the ℓ_0 sparse coding that induces sparse \mathbf{y} 's. We use least angle regression (LARS) [8] to solve the LASSO problem. We also consider orthogonal matching pursuit (OMP) [9], a greedy- ℓ_0 sparse coding algorithm that computes an at-most S -sparse \mathbf{y} extremely fast by

$$\min_{\mathbf{D}, \mathbf{y}} \|\mathbf{x} - \mathbf{D}\mathbf{y}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{y}\|_0 \leq S. \quad (2)$$

III. DEEP SPARSE-CODED NETWORK (DSN)

A. Notation

We denote input vector \mathbf{x} , its sparse code \mathbf{y} , and the pooled sparse code \mathbf{z} . An i th input vector (patch) is designated as $\mathbf{x}^{(i)}$. Sparse coding layers are the “hidden” layers of DSN. We use subscripted Roman numerals to indicate sparse coding layers. For example, \mathbf{D}_I means the first hidden layer’s dictionary. Note that the first hidden layer is the overall layer 2. Accordingly, \mathbf{y}_1 is the sparse code computed at the first hidden layer with the input \mathbf{x} and \mathbf{D}_I , and \mathbf{z}_1 the pooled sparse code over multiple \mathbf{y}_1 's.

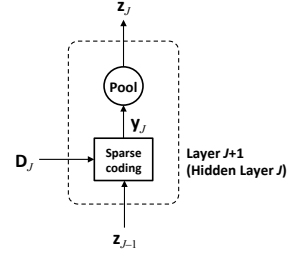


Fig. 2. DSN layering module

B. Architectural Overview

Deep Sparse-coded Network (DSN) is a feedforward network built on multilayer sparse coding. In Figure 1, we present an exemplar 4-layer DSN. This is a deep architecture since there are two hidden layers of sparse coding, each of which can learn corresponding level’s feature representations and train own dictionary of basis vectors. Similar to neural network, layers 1 and 4 are the input and output layers. The input layer takes in vectorized patches drawn from the raw data, which will be sparse coded and max pooled, propagating up the layers. The output layer consists of classifiers or regressors specific to application needs.

Figure 2 depicts a stackable layering module to build DSN. Sparse coding and pooling units together constitute the module. The J th hidden layer (for $J \geq II$) takes in pooled sparse codes \mathbf{z}_{J-1} 's from the previous hidden layer and produces \mathbf{y}_J using dictionary \mathbf{D}_J . Max pooling \mathbf{y}_J 's yields pooled sparse code \mathbf{z}_J that are passed as the input for hidden layer $J + 1$.

C. Algorithms

Hinton *et al.* [11] suggested *pretrain* a deep architecture with layer-by-layer unsupervised learning and finetune via *backpropagation*, a supervised algorithm popularized by neural network. We explain training algorithms for DSN using the architecture in Figure 1.

1) *Pretraining via layer-by-layer sparse coding and dictionary learning*: DSN takes in *spatially* contiguous patches from an image or *temporally* consecutive patches from time-series data to make the overall feature learning meaningful. Optimally, patches are preprocessed by normalization and whitening. The input layer is organized as pooling groups of M_1 patches: $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(M_1)}\}$, $\{\mathbf{x}^{(M_1+1)}, \mathbf{x}^{(M_1+2)}, \dots, \mathbf{x}^{(2M_1)}\}, \dots$. Sparse coding and dictionary learning at hidden layer 1 compute sparse codes $\mathbf{y}_1^{(i)}$'s

while learning \mathbf{D}_I jointly

$$\begin{aligned} \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M_1)}\} &\xrightarrow{\mathbf{D}_I} \{\mathbf{y}_I^{(1)}, \dots, \mathbf{y}_I^{(M_1)}\} \\ \{\mathbf{x}^{(M_1+1)}, \dots, \mathbf{x}^{(2M_1)}\} &\xrightarrow{\mathbf{D}_I} \{\mathbf{y}_I^{(M_1+1)}, \dots, \mathbf{y}_I^{(2M_1)}\} \\ &\vdots \end{aligned}$$

Max pooling at hidden layer 1 aggregates multiple sparse codes

$$\begin{aligned} \{\mathbf{y}_I^{(1)}, \dots, \mathbf{y}_I^{(M_1)}\} &\xrightarrow{\max \text{ pool}} \mathbf{z}_I^{(1)} \\ &\vdots \end{aligned}$$

Hidden layer 1 passes the pooled sparse codes $\{\mathbf{z}_I^{(1)}, \mathbf{z}_I^{(2)}, \dots\}$ to hidden layer 2. Sparse coding and dictionary learning continue at hidden layer 2 using \mathbf{z}_I 's as input

$$\begin{aligned} \{\mathbf{z}_I^{(1)}, \dots, \mathbf{z}_I^{(M_2)}\} &\xrightarrow{\mathbf{D}_{II}} \{\mathbf{y}_{II}^{(1)}, \dots, \mathbf{y}_{II}^{(M_2)}\} \\ &\vdots \end{aligned}$$

Pooling groups at hidden layer 2 consist of M_2 pooled sparse codes from hidden layer 1. Max pooling by M_2 yields

$$\begin{aligned} \{\mathbf{y}_{II}^{(1)}, \dots, \mathbf{y}_{II}^{(M_2)}\} &\xrightarrow{\max \text{ pool}} \mathbf{z}_{II}^{(1)} \\ &\vdots \end{aligned}$$

Pretraining completes by producing dictionaries $\{\mathbf{D}_I \in \mathbb{R}^{N \times K_1}, \mathbf{D}_{II} \in \mathbb{R}^{K_1 \times K_2}\}$ and the highest hidden layer's pooled sparse codes $\{\mathbf{z}_{II}^{(1)}, \mathbf{z}_{II}^{(2)}, \dots\}$ with each $\mathbf{z}_{II}^{(j)} \in \mathbb{R}^{K_2}$.

Max pooling is crucial for our DSN architecture. It reduces the total number of features by aggregating sparse codes to their max elements. More importantly, max pooling serves as a nonlinear activation function in neural network. Without *nonlinear* pooling, multilayering has no effect: $\mathbf{x} = \mathbf{D}_I \mathbf{y}_I$ and $\mathbf{y}_I = \mathbf{D}_{II} \mathbf{y}_{II}$ implies $\mathbf{x} = \mathbf{D}_I \mathbf{D}_{II} \mathbf{y}_{II} \approx \mathbf{D} \mathbf{y}_{II}$ because linear cascade of dictionaries is simply $\mathbf{D} \approx \mathbf{D}_I \mathbf{D}_{II}$ regardless of total number of layers.

2) *Training classifiers at output layer:* DSN learns each layer's dictionary greedily during pretraining. The resulting highest hidden layer output \mathbf{z}_{II} is already a powerful feature for classification tasks. Suppose DSN output layer predicts a class label $\hat{l} = h_{\mathbf{w}}(\phi)$, where $h_{\mathbf{w}}(\cdot)$ is a standard linear classifier or logistic/softmax regression that takes a feature encoding ϕ as input. Note that ϕ is encoded on \mathbf{z}_{II} , but depends on DSN setup. For instance, we may have $\phi = [\mathbf{z}_{II}^{(1)}; \mathbf{z}_{II}^{(2)}; \mathbf{z}_{II}^{(3)}; \mathbf{z}_{II}^{(4)}]$ if the highest hidden layer yields four pooled sparse codes per training example.

For simplicity, assume $\phi = \mathbf{z}_{II}$. A DSN classifier then computes $\hat{l} = h_{\mathbf{w}}(\mathbf{z}_{II}) = \mathbf{w}^\top \mathbf{z}_{II} + w_0$. We train the classifier weight $\mathbf{w} = [w_1 \dots w_{K_2}]^\top$ and the bias w_0 using labeled examples $\{(\mathbf{X}_1, l_1), \dots, (\mathbf{X}_m, l_m)\}$ in a supervised process by filling the input layer with patches from each example—the

i th example \mathbf{X}_i consists of $\{\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}, \dots\}$, where $\mathbf{x}_i^{(k)}$ is the k th patch—and working up the layers to compute \mathbf{z}_{II} 's that are used to train the DSN classifiers.

3) *Backpropagation:* By now, we have the DSN output layer with trained classifiers, and this is a good working pipeline for discriminative tasks. However, we might further improve the performance of DSN by optimizing the whole network in a supervised setting. Is backpropagation possible for DSN?

DSN backpropagation is quite different from conventional neural network or deep learning architectures. We explain our backpropagation idea again using the example DSN in Figure 1. The complete feedforward path of DSN is summarized by

$$\mathbf{x} \xrightarrow{\mathbf{D}_I} \mathbf{y}_I \xrightarrow{\max \text{ pool}} \mathbf{z}_I \xrightarrow{\mathbf{D}_{II}} \mathbf{y}_{II} \xrightarrow{\max \text{ pool}} \mathbf{z}_{II} \xrightarrow{\text{classify}} \hat{l}$$

We define the loss or cost function for the DSN classification

$$J(\mathbf{z}_{II}) = \frac{1}{2} \|\hat{l} - l\|^2 = \frac{1}{2} \|h_{\mathbf{w}}(\mathbf{z}_{II}) - l\|^2 \quad (3)$$

Our objective now is to propagate the loss value down the reverse path and adjust sparse codes. Fixing classifier weights \mathbf{w} , we back-estimate optimal \mathbf{z}_{II}^* that minimizes $J(\mathbf{z}_{II})$. To do so, we perform gradient descent learning with $J(\mathbf{z}_{II})$ that adjusts each element of vector \mathbf{z}_{II} by

$$z_{II,k} := z_{II,k} - \alpha \frac{\partial J(\mathbf{z}_{II})}{\partial z_{II,k}} \quad (4)$$

where α is the learning rate, $\mathbf{z}_{II} = [z_{II,1} \ z_{II,2} \ \dots \ z_{II,K_2}]^\top$, and K_2 is the number of basis vectors in dictionary \mathbf{D}_{II} for hidden layer 2. Since an optimal \mathbf{z}_{II}^* is estimated by correcting \mathbf{z}_{II} , the partial derivative is with respect to each element $z_{II,k}$

$$\frac{\partial J(\mathbf{z}_{II})}{\partial z_{II,k}} = [h_{\mathbf{w}}(\mathbf{z}_{II}) - l] \frac{\partial h(\mathbf{z}_{II})}{\partial z_{II,k}} = [h_{\mathbf{w}}(\mathbf{z}_{II}) - l] w_k$$

Here, note our linear classifier $h_{\mathbf{w}}(\mathbf{z}_{II}) = w_0 + w_1 \cdot z_{II,1} + \dots + w_{K_2} \cdot z_{II,K_2}$. Therefore, the following gradient descent rule adjusts \mathbf{z}_{II} to \mathbf{z}_{II}^*

$$z_{II,k} := z_{II,k} + \alpha [l - h_{\mathbf{w}}(\mathbf{z}_{II})] w_k \quad (5)$$

This update rule is intuitive because it down-propagates the error $[l - h_{\mathbf{w}}(\mathbf{z}_{II})]$ proportionately to the contribution from each $z_{II,k}$ and adjusts accordingly.

Using the corrected \mathbf{z}_{II}^* , we can correct the unpooled original \mathbf{y}_{II} 's to optimal \mathbf{y}_{II}^* 's by a procedure called *putback* illustrated in Figure 3. At hidden layer 2, we have performed max pooling by M_2 . For putback, we need to keep the original M_2 \mathbf{y}_{II} 's that have resulted \mathbf{z}_{II} in memory so that corrected values at \mathbf{z}_{II}^* are *put back* to corresponding locations at the original sparse codes \mathbf{y}_{II} 's and yield error-adjusted \mathbf{y}_{II}^* 's.

With \mathbf{y}_{II}^* , going down a layer is straightforward. By sparse coding relation, we just compute $\mathbf{z}_I^* = \mathbf{D}_{II} \mathbf{y}_{II}^*$. Next, we do

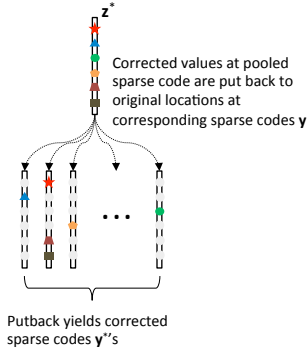


Fig. 3. Putback corrects sparse codes \mathbf{y} from \mathbf{z}^*

another putback at hidden layer 1. Using \mathbf{z}_I^* , we obtain \mathbf{y}_I^* 's. Each pooling group at hidden layer 1 originally has M_1 \mathbf{y}_I 's that need to be saved in memory.

With \mathbf{y}_I^* 's, we should now correct \mathbf{D}_I , not \mathbf{x} , because it does not make sense to correct given data input. Hence, down-propagation of the error for DSN stops here, and we up-propagate corrected sparse codes to finetune the dictionaries and classifier weights. We consider two methods to adjust basis vectors of the dictionaries. First method is to use rank-1 update in batches. The idea is to compute the residual matrix that isolates a contribution by each basis vector only. We update each basis vector iteratively with the first principle component computed via singular value decomposition of the residual (*i.e.*, the first column of matrix \mathbf{U}). This method is also used in the inner-loop of K-SVD [12].

Our second method uses online gradient descent. We define the loss function with respect to \mathbf{D}_I

$$J(\mathbf{D}_I) = \frac{1}{2} \|\mathbf{D}_I \mathbf{y}_I^* - \mathbf{x}\|_2^2 \quad (6)$$

Adjusting \mathbf{D}_I requires to solve the following optimization problem given examples $(\mathbf{x}, \mathbf{y}_I^*)$

$$\min_{\mathbf{d}_{I,k}} J(\mathbf{D}_I) \quad \text{s.t.} \quad \|\mathbf{d}_{I,k}\|_2^2 = 1 \quad \forall k \quad (7)$$

where $\mathbf{d}_{I,k}$ is the k th basis vector in \mathbf{D}_I . Taking the partial derivative with respect to $\mathbf{d}_{I,k}$ yields

$$\frac{\partial J(\mathbf{D}_I)}{\partial \mathbf{d}_{I,k}} = (\mathbf{D}_I \mathbf{y}_I^* - \mathbf{x}) y_{I,k}^*$$

where $\mathbf{y}_I^* = [y_{I,1}^* \dots y_{I,K_1}^*]^\top$ and $\mathbf{y}_I = [y_{I,1} \dots y_{I,K_1}]^\top$. We obtain the update rule to adjust \mathbf{D}_I by gradient descent

$$\mathbf{d}_{I,k} := \mathbf{d}_{I,k} - \beta (\mathbf{D}_I \mathbf{y}_I^* - \mathbf{x}) y_{I,k}^* \quad (8)$$

We denote the corrected dictionary \mathbf{D}_I^* . We redo sparse coding at hidden layer 1 with \mathbf{D}_I^* followed by max pooling. Similarly at hidden layer 2, we update \mathbf{D}_{II} to \mathbf{D}_{II}^* by

$$\mathbf{d}_{II,k} := \mathbf{d}_{II,k} - \gamma (\mathbf{D}_{II} \mathbf{y}_{II}^* - \mathbf{z}_I^\dagger) y_{II,k}^* \quad (9)$$

where \mathbf{z}_I^\dagger is the pooled sparse code over M_1 \mathbf{y}_I^\dagger 's from sparse coding redone with \mathbf{D}_I^* . Using corrected dictionary \mathbf{D}_{II}^* , we also redo sparse coding and max pooling at hidden layer 2. The resulting pooled sparse codes \mathbf{z}_{II}^\dagger are the output of the highest hidden layer, which will be used to retrain the classifier h_w . All of the steps just described are a single iteration of DSN backpropagation. We run multiple iterations until convergence.

The corrections made during down-propagation for DSN backpropagation are summarized by

$$\mathbf{z}_{II} \xrightarrow{\text{GD}} \mathbf{z}_{II}^* \xrightarrow{\text{putback}} \mathbf{y}_{II}^* \xrightarrow{\mathbf{D}_{II}} \mathbf{z}_I^* \xrightarrow{\text{putback}} \mathbf{y}_I^*.$$

The corrections by up-propagation follow

$$\mathbf{D}_I \xrightarrow{\text{GD}} \mathbf{D}_I^* \xrightarrow{\text{SC}} \mathbf{y}_I^\dagger \xrightarrow{\text{max pool}} \mathbf{z}_I^\dagger \xrightarrow{\text{GD}} \mathbf{D}_{II}^* \xrightarrow{\text{SC}} \mathbf{y}_{II}^\dagger \xrightarrow{\text{max pool}} \mathbf{z}_{II}^\dagger \xrightarrow{\text{GD}} h_w$$

where GD stands for gradient descent, and SC sparse coding. We present the backpropagation algorithm for general L -layer DSN in Algorithm 1.

Algorithm 1 DSN backpropagation

require Pretrained $\{\mathbf{D}_I, \mathbf{D}_{II}, \dots, \mathbf{D}_{L-2}\}$ and classifier h_w
input Labeled training examples $\{(\mathbf{X}_1, l_1), \dots, (\mathbf{X}_m, l_m)\}$
output Fine-tuned $\{\mathbf{D}_I^*, \mathbf{D}_{II}^*, \dots, \mathbf{D}_{L-2}^*\}$ and classifier h_w^*

- 1: **repeat**
- 2: **subalgorithm** Down-propagation
- 3: **for** $J := L - 2$ **to** 1
- 4: **if** $J == L - 2$
- 5: Compute classifier error $\epsilon^{(i)} = l^{(i)} - h_w(\mathbf{z}_J^{(i)}) \quad \forall i$
- 6: Compute $\mathbf{z}_J^{*(i)}$ by $z_{J,k}^{*(i)} = z_{J,k}^{(i)} + \alpha \cdot \epsilon^{(i)} \cdot w_k \quad \forall i, k$
- 7: **else**
- 8: Compute $\mathbf{z}_J^{*(i)} = \mathbf{D}_{J+1} \mathbf{y}_{J+1}^{*(i)} \quad \forall i$
- 9: **end**
- 10: Estimate $\mathbf{y}_J^{*(i)}$ from $\mathbf{z}_J^{*(i)}$ via putback $\quad \forall i$
- 11: **end**
- 12: **end**
- 13: **subalgorithm** Up-propagation
- 14: **for** $J := 1$ **to** $L - 2$
- 15: **if** $J == 1$
- 16: Compute \mathbf{D}_I^* by Eq. (8)
- 17: Compute $\mathbf{y}_I^{\dagger(i)}$ by sparse coding with $\mathbf{D}_I^* \quad \forall i$
- 18: Compute $\mathbf{z}_I^{\dagger(i)}$ by max pooling $\quad \forall i$
- 19: **else**
- 20: Compute \mathbf{D}_J^* by Eq. (9)
- 21: Compute $\mathbf{y}_J^{\dagger(i)}$ by sparse coding with $\mathbf{D}_J^* \quad \forall i$
- 22: Compute $\mathbf{z}_J^{\dagger(i)}$ by max pooling $\quad \forall i$
- 23: **end**
- 24: **end**
- 25: Retrain classifier h_w with $\{\mathbf{z}_{L-2}^{\dagger(i)}, l^{(i)}\} \quad \forall i$
- 26: **end**
- 27: **until** converged

IV. EXPERIMENTS

We present a comparative performance analysis of single-layer sparse coding, deep stacked autoencoder (SAE) [13], and DSN on the multi-class image classification task using CIFAR-10. Both SAE and DSN have 4 layers.

A. Data processing and training for sparse coding

Instead of using the full CIFAR-10 dataset, we uniformly sample 20,000 images and cut to five folds for cross validation. We use four folds for training and the remaining fold for testing. We enforce exactly 2,000 images per class. For output layer, we have trained a 1-vs-all linear classifiers for each of ten classes in CIFAR-10. Each datum in CIFAR-10 is a $32 \times 32 \times 3$ color image. We consider a *per-image* feature vector from densely overlapping patches drawn from a receptive field with width $w = 6$ pixels and stride $s = 2$. Thus, each patch (vectorized) has size $N = 3 \times 6 \times 6 = 108$. We preprocess patches by ZCA-whitening before sparse coding. We use a couple of different sparsity configurations for each LARS and OMP. We configure hidden layer 1 sparse coding more densely with $\lambda = 0.1$ (regularization penalty) for LARS and $S = 0.2N$ (sparsity bound) for OMP. For hidden layer 2, we use $\lambda = 0.2$ and $S = 0.1N$.

Figure 4 illustrates sparse coding and max pooling at hidden layer 1. Each image is divided into four quadrants. For each quadrant, there are four (pooling) groups of 9 patches. Hidden layer 1 uses a dictionary size $K_1 = 4N = 432$ and max pooling factor $M_1 = 9$. Hidden layer 1 produces $\{\mathbf{z}_I^{(1)}, \dots, \mathbf{z}_I^{(4)}\}$, $\{\mathbf{z}_I^{(5)}, \dots, \mathbf{z}_I^{(8)}\}$, $\{\mathbf{z}_I^{(9)}, \dots, \mathbf{z}_I^{(12)}\}$, and $\{\mathbf{z}_I^{(13)}, \dots, \mathbf{z}_I^{(16)}\}$ (4 pooled sparse codes per quadrant), which will be passed to hidden layer 2.

Figure 5 illustrates sparse coding and max pooling at hidden layer 2. We use $K_2 = 2K_1 = 864$ and $M_2 = 4$. To prevent further expansion of features, we have computed the two averages and encoded the final per-image feature vector $\phi_{DSN} = [\text{mean}(\mathbf{z}_{II}^{(1)}, \mathbf{z}_{II}^{(2)}); \text{mean}(\mathbf{z}_{II}^{(3)}, \mathbf{z}_{II}^{(4)})]$. Thus, the final feature vector for DSN has a dimensionality $2K_2 = 1,728$. For single-layer sparse coding, we have averaged $\mathbf{z}_I^{(i)}$'s per quadrant and stacked them to form the final per-image feature vector ϕ_{SLSC} , which also has 1,728 features.

B. Data processing and training for autoencoder

For fairness, we have matched the total number of trained weights between the DSN and SAE schemes. Given an input $\mathbf{x} \in \mathbb{R}^N$, an autoencoder layer [1] trains a set of encoding weights that transforms \mathbf{x} into the activations $\mathbf{a} \in \mathbb{R}^H$ and a set of decoding weights that recovers $\hat{\mathbf{x}}$, an estimate of \mathbf{x} , from \mathbf{a} , all through backpropagation. Similar to sparse coding, we use the encoding \mathbf{a} as the feature vector for \mathbf{x} . If $H < N$, the autoencoder layer is forced to learn a *compressed* representation of the input. Also, even when $H > N$ (i.e., more hidden units than input dimension), we can still learn meaningful features by imposing a *sparsity* constraint on \mathbf{a} such that only $S \ll H$ neurons are activated. We have experimented with both approaches and will report the best result of the two for deep SAE.

TABLE I
AVERAGE 1-VS-ALL CLASSIFICATION ACCURACY FOR SINGLE-LAYER SPARSE CODING AND AUTOENCODER

	Classification accuracy
Autoencoder	69.8%
OMP ($S = 0.1N$)	75.3%
OMP ($S = 0.2N$)	76.9%
LARS ($\lambda = 0.2$)	78.4%
LARS ($\lambda = 0.1$)	80.1%

TABLE II
AVERAGE 1-VS-ALL CLASSIFICATION ACCURACY COMPARISON BETWEEN DSN AND SAE

	Classification accuracy
Deep SAE (pretraining only)	71.8%
Deep SAE (pretraining+backprop)	78.9%
DSN-OMP (pretraining only)	79.6%
DSN-OMP (pretraining+backprop)	84.3%
DSN-LARS (pretraining only)	83.1%
DSN-LARS (pretraining+backprop)	87.5%

C. Results

We report cross-validated 1-vs-all classification accuracy of single-layer sparse coding, DSN, and deep SAE. In Table I, we present the single-layer autoencoder and sparse coding performances. For classification, we have used two standard, off-the-shelf algorithms, SVM and logistic regression. The table summarizes the better of the two. LARS with $\lambda = 0.1$ has achieved the best single-layer accuracy at 80.1%.

In Table II, we compare the performance of DSN against deep SAE in various configurations. All DSN schemes show improvements from their respective single-layer configurations. Optimization by backpropagation over the whole network is critical for deep SAE as evidenced in the accuracy gain of 7.7%, which is significantly higher than the 2% gain from multilayering. For DSN, multilayering improves about 3%, and the proposed backpropagation additional 4%. Importantly, DSN-OMP with only pretraining is already 0.7% better than the backpropagation-finetuned deep SAE. DSN-OMP improves by 4.7% on backpropagation whereas the improvement is slightly less for DSN-LARS with a 4.4% gain. Overall, we find DSN-LARS the best performer.

V. CONCLUSION

Motivated by superior feature learning performance of single-layer sparse coding, we have presented Deep Sparse-coded Network (DSN), a deep architecture for sparse coding. We have discussed the benefit of DSN and described training methods including a novel backpropagation algorithm that effectively traverses and optimizes multiple layers of sparse coding and max pooling.

We stress that this paper limits to report an evaluation of DSN that confirms superior classification accuracy in a

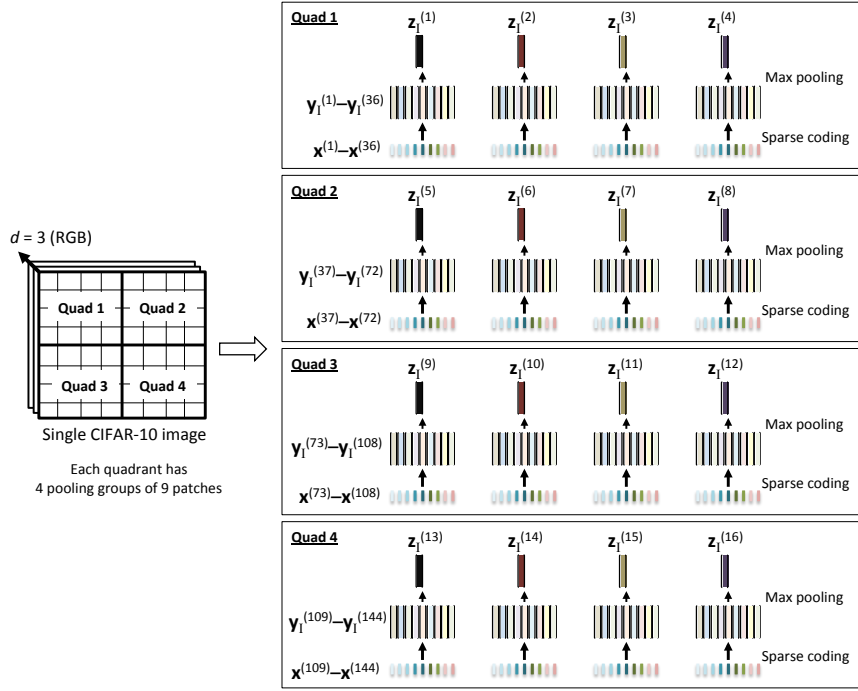


Fig. 4. Sparse coding and max pooling at hidden layer 1 for single CIFAR-10 image

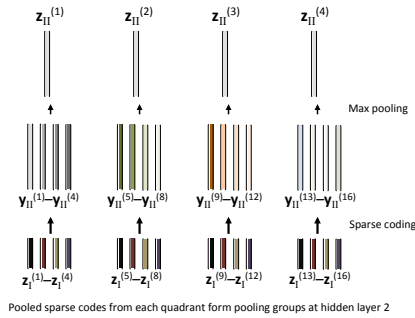


Fig. 5. Sparse coding and max pooling at hidden layer 2

medium-sized setup with CIFAR-10 images. We are currently experimenting with DSN more rigorously using larger datasets CIFAR-100, Caltech-101, and Caltech-256. In future work, we will test DSN in broader scope for text, sound, and wireless signal classification tasks.

ACKNOWLEDGMENTS

This work is supported in part by gifts from the Intel Corporation and in part by the Naval Supply Systems Command award under the Naval Postgraduate School Agreements No. N00244-15-0050 and No. N00244-16-1-0018.

REFERENCES

[1] G. E. Hinton and R. R. Salakhutdinov, "Reducing the Dimensionality of Data with Neural Networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.

[2] R. Salakhutdinov, A. Mnih, and G. Hinton, "Restricted Boltzmann Machines for Collaborative Filtering," in *ICML*, 2007.

[3] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-Based Learning Applied to Document Recognition," in *Proc. of the IEEE*, vol. 86, no. 11, 1998, pp. 2278–2324.

[4] A. Coates, A. Y. Ng, and H. Lee, "An Analysis of Single-layer Networks in Unsupervised Feature Learning," in *AISTATS*, 2011.

[5] Y. Karklin and M. S. Lewicki, "Learning Higher-order Structures in Natural Images," *Network*, vol. 14, no. 3, pp. 483–99, 2003.

[6] C. Cadieu and B. A. Olshausen, "Learning Transformational Invariants from Natural Movies," in *NIPS*, 2008.

[7] R. Tibshirani, "Regression Shrinkage and Selection via the Lasso," *Journal of Royal Statistical Society, Series B*, vol. 58, pp. 267–288, 1994.

[8] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani, "Least Angle Regression," *Annals of Statistics*, vol. 32, pp. 407–499, 2004.

[9] Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad, "Orthogonal Matching Pursuit: Recursive Function Approximation with Applications to Wavelet Decomposition," in *Asilomar Conference on Signals, Systems and Computers*, 1993.

[10] B. Olshausen and D. Field, "Sparse Coding with an Overcomplete Basis Set: Strategy Employed by V1?" *Vision research*, 1997.

[11] G. Hinton, S. Osindero, and Y. Teh, "A Fast Learning Algorithm for Deep Belief Nets," *Neural Computation*, vol. 18, no. 7, pp. 1527–1554, 2006.

[12] M. Aharon, M. Elad, and A. Bruckstein, "K-SVD: An Algorithm for Designing Overcomplete Dictionaries for Sparse Representation," *IEEE Trans. on Sig. Proc.*, vol. 54, no. 11, 2006.

[13] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, "Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion," *Journal of Machine Learning Research*, vol. 11, pp. 3371–3408, 2010.