
No more meta-parameter tuning in unsupervised sparse feature learning

Adriana Romero

ADRIANA.ROMERO@UB.EDU

Departament de Matemàtica Aplicada i Anàlisi, Universitat de Barcelona, Barcelona, Spain.

Petia Radeva

PETIA.IVANOV@UB.EDU

Departament de Matemàtica Aplicada i Anàlisi, Universitat de Barcelona, Barcelona, Spain.

Carlo Gatta

CGATTA@CVC.UAB.ES

Centre de Visió per Computador, Bellaterra, Spain.

Abstract

We propose a meta-parameter free, off-the-shelf, simple and fast unsupervised feature learning algorithm, which exploits a new way of optimizing for sparsity. Experiments on STL-10 show that the method presents state-of-the-art performance and provides discriminative features that generalize well.

1. Introduction

Significant effort has been devoted to handcraft appropriate feature representations of data in several fields. In tasks such as image classification and object recognition, unsupervised learned features have shown to compete well or even outperform manually designed ones (Ranzato et al., 2006; Yang et al., 2009; Coates et al., 2011). Unsupervised feature learning has also shown to be helpful in greedy layerwise pre-training of deep architectures (Hinton et al., 2006; Bengio et al., 2006; Larochelle et al., 2009; Erhan et al., 2010).

In (Bengio, 2009), the author claims that potentially interesting research involves pre-training algorithms, which “[...] would be proficient at extracting good features but involving an easier optimization problem.” In addition to that, one of the main criticisms to state-of-the-art methods is that they require a significant amount of meta-parameters (Bengio et al., 2013). As stated in (Snoek et al., 2012), the tuning of these meta-parameters is a laborious task that requires expert knowledge, rules of thumb or extensive search and, whose setting can vary for dif-

ferent tasks. Therefore, there is great interest for meta-parameter free methods (Ngiam et al., 2011) and automatic approaches to optimize the performance of learning algorithms (Snoek et al., 2012).

Nevertheless, little effort has been devoted to address this problem (see Table 1 for a comparison of meta-parameters required by unsupervised feature learning methods). To the best of our knowledge, work in this direction includes ICA (Hyvärinen & Oja, 2000; Hyvärinen et al., 2000) and sparse filtering (Ngiam et al., 2011). Although ICA provides good results at object recognition tasks (Le et al., 2011; Ngiam et al., 2011), the method scales poorly to large datasets and high input dimensionality.

Computational complexity is also a major drawback of many state-of-the-art methods. ICA requires an expensive orthogonalization to be computed at each iteration. Sparse coding has an expensive inference, which requires a prohibitive iterative optimization. Significant amount of work has been done in order to overcome this limitation (Lee et al., 2006; Kavukcuoglu et al., 2010). Predictive Sparse Decomposition (PSD) (Kavukcuoglu et al., 2010) is a successful variant of sparse coding, which uses a predictor to approximate the sparse representation and solves the sparse coding computationally expensive encoding step.

In this paper, we aim to solve some of the above-mentioned problems. We propose a **meta-parameter free, off-the-shelf, simple and fast** approach, which exploits a new way of optimizing for a sparsity, without explicitly modeling the data distribution. The method iteratively builds an *ideally* sparse target and optimizes the dictionary by minimizing the error between the system output and the *ideally* sparse target. *Defining sparsity concepts in terms of expected output allows to exploit a new strategy in unsupervised training.*

Table 1. Meta-parameters to tune of state-of-the-art unsupervised feature learning methods.

Method	Meta-parameters to tune
Sparse RBM (Hinton et al., 2006; Lee et al., 2008)	weight decay, sparseness constant, sparsity penalty, momentum
Sparse auto-encoders (Ranzato et al., 2006)	weight decay, sparseness constant, sparsity penalty
Sparse Coding (Olshausen & Field, 1997)	sparsity penalty
RICA (Le et al., 2011)	reconstruction penalty
PSD (Kavukcuoglu et al., 2010)	sparsity penalty, prediction penalty
OMP-k (Pati et al., 1993; Blumensath & Davies, 2007; Coates & Ng, 2011)	k (non-zero elements)
ICA (Hyvärinen & Oja, 2000; Hyvärinen et al., 2000)	-
Sparse Filtering (Ngiam et al., 2011)	-

It is worth stressing that many optimization strategies can be used to minimize the above-mentioned error and that parameters of these optimization techniques must not be considered as belonging to our approach.

Experiments on STL-10 dataset show that the method outperforms state-of-the-art methods in single layer image classification, providing discriminative features that generalize well.

Linear feature extraction methods combined with sparse coding encodings are among best performers on object recognition datasets. The importance of properly combining training/encoding and encoding/pooling strategies has been argued in (Coates & Ng, 2011) and (Zeiler & Fergus, 2013) respectively. Since the goal of this paper is to propose a new method for unsupervised feature learning, dealing with all the possible combinations of encoding and pooling could mask the benefits of the method that we propose. However, for the sake of fair comparison with the state-of-the-art, we test the method with sparse coding and soft-threshold encodings combined with sum pooling, following the experimental pipeline of (Coates & Ng, 2011).

2. State-of-the-art

Commonly used algorithms for unsupervised feature learning include Restricted Boltzmann Machines (RBM) (Hinton et al., 2006), auto-encoders (Bengio et al., 2006), sparse coding (Raina et al., 2007) and hybrids such as PSD (Kavukcuoglu et al., 2010). Many other methods such as ICA (Hyvärinen & Oja, 2000; Hyvärinen et al., 2000), Reconstruction ICA (RICA) (Le et al., 2011), Sparse Filtering (Ngiam et al., 2011) and methods related to vector quantization such as Orthogonal Matching Pursuit (OMP-k) (Pati et al., 1993; Blumensath & Davies, 2007; Coates & Ng, 2011) have also been used in the literature to extract unsupervised feature representations. These algorithms could be divided into two categories: explicitly modeling or not the input distribution. Sparse auto-encoders (Ranzato et al., 2006), sparse RBM (Hinton et al., 2006; Lee et al., 2008; Hinton,

2010; Goh et al., 2012), sparse coding (Olshausen & Field, 1997), PSD (Kavukcuoglu et al., 2010), OMP-k (Pati et al., 1993; Blumensath & Davies, 2007; Coates & Ng, 2011) and Reconstruction ICA (RICA) (Le et al., 2011) explicitly model the data distribution by minimizing the reconstruction error. Although learning a good approximation of the data distribution may be desirable, approaches such as sparse filtering (Ngiam et al., 2011) show that this seems not so important if the goal is to have a discriminative sparse system. Sparse filtering does not attempt to explicitly model the input distribution but focuses on the properties of the output distribution instead.

Sparsity is among the desirable properties of a good output representation (Field, 1994; Olshausen & Field, 1997; Ranzato et al., 2006; Lee et al., 2008; Le et al., 2011; Ngiam et al., 2011; Bengio et al., 2013). Sparse features consist of a large amount of outputs, which respond rarely and provide high responses when they do respond. Sparsity can be described in terms of population sparsity and lifetime sparsity (Willmore & Tolhurst, 2001). Both lifetime and population sparsity are important properties of the output distribution. On one hand, lifetime sparsity plays an important role in preventing bad solutions such as numerous dead outputs. There seems to be a consensus to overcome such degenerate solutions, which is to ensure similar statistics among outputs (Field, 1994; Willmore & Tolhurst, 2001; Ranzato et al., 2006; Ngiam et al., 2011). On the other hand, population sparsity helps providing a simple interpretation of the input data such as the ones found in early visual areas. To the best of our knowledge, the definition of population sparsity remains ambiguous.

State-of-the-art methods optimize either for one or both sparsity forms in their objective function. The great majority seeks sparsity using the L_1 penalty and does not optimize for an explicit level of sparsity in their outputs. Sparse auto-encoders optimize for a target activation allowing to deal with lifetime sparsity; nevertheless, the target activation requires tuning and does not explicitly control the level of population sparsity. OMP-k defines the level of population sparsity by setting k to the maximum expected

number of non-zero elements per output code, whereas the methods in (Olshausen & Field, 1997; Ranzato et al., 2006; Lee et al., 2008; Le et al., 2011; Ngiam et al., 2011) do not explicitly define the proportion of outputs expected to be active at the same time.

3. Method

In this section, we describe how the proposed method learns a sparse feature representation of the data in terms of population and lifetime sparsity. The method iteratively builds an *ideally* sparse target and optimizes the dictionary by minimizing the error between the system output and the *ideally* sparse target. Subsection 3.1 highlights the algorithm to enforce lifetime and population sparsity in the ideally sparse target. Subsection 3.2 provides implementation details on the system and optimization strategies used to minimize the error between the system output and the ideally sparse target.

3.1. Enforcing Population and Lifetime Sparsity by defining an ideal target

We define population and lifetime sparsity as properties of an *ideal* sparse output. Given N training samples and an output of dimensionality N_h , we define the first property of the output as:

1. **Strong Lifetime Sparsity:** The output vectors must be composed solely of active and inactive units (no intermediate values between two fixed scalars are allowed) and all outputs must activate for an equal number of inputs. Activation is exactly distributed among the N_h outputs.

Our Strong Lifetime Sparsity definition is a more strict requirement than the high dispersal concept introduced in (Ngiam et al., 2011), since they only require that “the mean squared activations of each feature (output) [...] should be roughly the same for all features (outputs)”. While high dispersal attempts to diversify the learned bases, it does not guarantee the output distribution, in the lifetime sense, to be composed of only a few activations. Furthermore, our definition ensures the absence of dead outputs.

Given our definition of Strong Lifetime Sparsity, the population sparsity must require that, for each training sample, only one output element is active:

2. **Strong Population Sparsity:** For each training sample only one output must be active.

The rationale of our approach is to appropriately generate an ideal output target that fulfils properties (1) and (2), and then learn the parameters of the system by minimizing the

L_2 error between the output target and the output generated by the system during training. In this way, we seek a system optimized for both population and lifetime sparsity in an explicit way.

The key component of our approach is how to define the ideal output target based on the above-mentioned properties. However, to ensure that the optimization of the system parameters converges, we add a third property:

3. **Minimal Perturbation:** The ideal output target should be defined as the best approximation of the system output by means of L_2 error fulfilling properties (1) & (2).

Creating the output target that ensures the above-mentioned properties is analogous to solving an assignment problem. The Hungarian method (Kuhn, 1955) is a combinatorial optimization algorithm, which solves the assignment problem. However, its computational cost $\mathcal{O}((NN_h)^{3/2})$ is prohibitive. Therefore, in the next section we propose a simple and fast $\mathcal{O}(NN_h)$ algorithm to generate the ideal output target, which ensures sparsity properties (1) and (2) and provides an approximate solution for minimal perturbation property (3).

3.1.1. IDEAL TARGET GENERATION: THE ENFORCING POPULATION AND LIFETIME SPARSITY (EPLS) ALGORITHM

Let us assume that we have a system, which produces a row output vector \mathbf{h} . We use the notation \mathbf{h}_j to refer to one element of \mathbf{h} . We define an output matrix \mathbf{H} composed of N_b output vectors of dimensionality N_h , such that $N_b \leq N$. Likewise, we define an ideal target output matrix \mathbf{T} of the same size. Algorithm 1 details the EPLS algorithm to generate the *ideal* target \mathbf{T} from \mathbf{H} . For the sake of simplicity, every step of the algorithm where the subscript j appears must be applied $\forall j \in \{1, 2, \dots, N_h\}$.

Algorithm 1 EPLS

Require: \mathbf{H} , \mathbf{a} , N

Ensure: \mathbf{T} , \mathbf{a}

```

1:  $\mathbf{T} = \mathbf{0}$ 
2: for  $n = 1 \rightarrow N_b$  do
3:    $\mathbf{h}_j = \mathbf{H}_{n,j}$ 
4:    $k = \arg \max_j (\mathbf{h}_j - \mathbf{a}_j)$ 
5:    $\mathbf{T}_{n,k} = 1$ 
6:    $\mathbf{a}_k = \mathbf{a}_k + \frac{N_h}{N} + \epsilon$ 
7: end for
8: Remap  $\mathbf{T}$  to active/inactive values of the corresponding function.
```

Starting with no activation in \mathbf{T} (line 1), the algorithm proceeds as follows. A row vector \mathbf{h} from \mathbf{H} is processed at each iteration (line 3). The crucial step is performed in line

4: the output k that has to be activated in the n^{th} row of \mathbf{T} is selected as the one that has the maximal activation value \mathbf{h}_j minus the inhibitor \mathbf{a}_j . The inhibitor \mathbf{a}_j can be seen as an accumulator that “counts” the number of times an output j has been selected, increasing its inhibition progressively by N_h/N until reaching maximal inhibition. This prevents the selection of an output that has already been activated N/N_h times. The rationale behind the equation in line 4 is that, while selecting the maximal responses in the matrix \mathbf{H} , we have to take care to distribute them evenly among all outputs (in order to ensure Strong Lifetime Sparsity). Using this strategy, it can be demonstrated that the resulting matrix \mathbf{T} perfectly fulfills properties (1) and (2). In line 5, the algorithm activates the k^{th} element of n^{th} row of the target matrix \mathbf{T} . By activating the “relative” maximum, we approximate property (3). Finally, the inhibitor \mathbf{a} is updated in line 6.

3.2. System and Optimization strategies

Let us assume that we have a system parameterized by $\Gamma = \{\mathbf{W}, \mathbf{b}\}$, with activation function f , which takes as input a data vector \mathbf{d} and produces an output vector $\mathbf{h} = f(\mathbf{d}, \Gamma)$. We use the same notation as in Section 3 and define a data matrix \mathbf{D} composed of N rows and N_d columns, where N_d is the input dimensionality.

To compare our training strategy to previous well known systems, we tested our algorithm using

$$\mathbf{H} = f(\mathbf{D}\mathbf{W} + \mathbf{b}), \quad (1)$$

where f is a logistic non-linearity.

3.2.1. OPTIMIZATION STRATEGY

The system might be trained by means of an off-the-shelf mini-batch Stochastic Gradient Descent (SGD) method with adaptive learning rates such as variance-based SGD (vSGD) (Schaul et al., 2013). Algorithm 2 details the latter training process. The mini-batch size N_b can be set to any value, in all the experiments we have set $N_b = N_h$. Starting with Γ set to small random numbers as in (LeCun et al., 1998) (line 1), at each epoch we shuffle the samples of the training set (line 3), reset the EPLS inhibitor \mathbf{a} to a flat activation (line 4) and process all mini-batches. For each mini-batch b , samples $\mathbf{D}^{(b)}$ are selected (line 6). Then, the output $\mathbf{H}^{(b)}$ is computed (line 7) and the EPLS is invoked to compute $\mathbf{T}^{(b)}$ and update \mathbf{a} (line 8). After that, the gradient of the error is computed (line 9) and the learning rate η is estimated as in (Schaul et al., 2013) (line 10). The system parameters are then updated to minimize the L_2 error $E^{(b)} = \|\mathbf{H}^{(b)} - \mathbf{T}^{(b)}\|_2^2$ (line 11). Finally, the bases \mathbf{W} in Γ are limited to have unit norm to avoid degenerate solutions (line 13). This procedure is repeated until a stop condition is met; in our experiments, the training stops when the rel-

ative decrement error between epochs is small ($< 10^{-6}$).

When updating the system parameters, we assume that \mathbf{T} does not depend on Γ , thus $\frac{\partial \mathbf{T}}{\partial \Gamma} = 0$; we carried out experiments that show that this approximation does not significantly influence the gradient descent convergence nor the quality of the minimization. Moreover, this assumption makes the algorithm faster, since we remove the need of computing the numerical partial derivatives of \mathbf{T} .

The mini-batch vSGD allows to scale the algorithm easily, especially with respect to the number of samples N .

Algorithm 2 Standard EPLS training

Require: \mathbf{D}

Ensure: Γ

```

1:  $\Gamma =$  small random values
2: repeat
3:   Shuffle  $\mathbf{D}$  randomly
4:    $\mathbf{a} =$  flat activation
5:   for  $b = 1 \rightarrow \lfloor N/N_b \rfloor$  do
6:     Select mini-batch samples  $\mathbf{D}^{(b)}$ 
7:      $\mathbf{H}^{(b)} = f(\mathbf{D}^{(b)}, \Gamma)$ 
8:      $(\mathbf{T}^{(b)}, \mathbf{a}) = \text{EPLS}(\mathbf{H}^{(b)}, \mathbf{a}, N)$ 
9:      $\mathbf{G} = \nabla_{\Gamma} \|\mathbf{H}^{(b)} - \mathbf{T}^{(b)}\|_2^2$ 
10:    Estimate learning rate  $\eta$  as in (Schaul et al., 2013)
11:     $\Gamma = \Gamma - \eta \mathbf{G}$ 
12:  end for
13:  Limit the bases  $\mathbf{W}$  in  $\Gamma$  to have unit norm
14: until stop condition verified
```

4. Experiments

The performance of training and encoding strategies in single layer networks has been extensively analyzed in the literature (Coates et al., 2011; Coates & Ng, 2011; Ngiam et al., 2011) on STL-10¹ dataset. STL-10 dataset consists of 96x96 pixels color images belonging to 10 different classes. The dataset is divided into a large unlabeled training set containing 100K images and smaller labeled training and test sets, containing 5000 and 8000 images, respectively. It has to be considered that in STL-10, the primary challenge is to make use of the unlabeled data (100K images), which is 100 times bigger than the labeled data used to train the classifier (1000 images per fold). In this case, the supervised training must strongly rely on the ability of the unsupervised method to learn discriminative features. Moreover, since the unlabeled dataset contains other types of animals (bears, rabbits, etc.) and vehicles (trains, buses, etc.) in addition to the ones in the labeled set, the unsupervised method should be able to generalize well.

To validate our method, we follow the experimental pipeline of (Coates et al., 2011). We first extract random patches and normalize them for local brightness and con-

¹<http://www.stanford.edu/~acoates/stl10/>

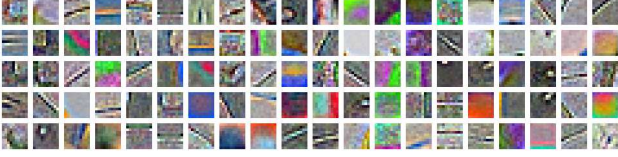


Figure 1. Random subset of bases learned by EPLS, a receptive field of 10 pixels and $N_h = 1600$ (better seen in color).

Table 2. Classification accuracy on STL-10.

Algorithm		Accuracy	
Single-Layer <i>with</i> meta-parameters			
RICA (Le et al., 2011) (1600/Natural)		52.9%	
OMP-1 (1600/Natural)		51.8% (0.47%)	
OMP-1 (whitening, 1600/Natural)		53.1% (0.52%)	
OMP-1 (whitening, 1600x2/Natural)		54.5% (0.66%)	
OMP-1 (whitening, 1600x2/SC)		59.0% (0.80%)	
Single-Layer <i>without</i> meta-parameters			
Raw pixels		31.8% (0.62%)	
ICA (whitening, Complete/Natural)		48.0% (1.47%)	
K-means-tri (whitening, 1600)		51.5% (1.73%)	
Sparse Filtering (1600/Natural)		53.5% (0.53%)	
	Natural (1600)	Natural (1600x2)	SC (1600x2)
EPLS	56.6% (0.66%)	56.9% (0.50%)	61.0% (0.58%)

trast. Note that EPLS does not require any whitening of the input data, since it decorrelates the data during the training by means of the imposed strong sparsity properties of the output target. Then, we apply the system to retrieve sparse features of patches covering the input image, pool them into 4 quadrants and finally train a L_2 SVM for classification purposes. We tune the SVM parameter using 5-fold cross-validation. As in (Ngiam et al., 2011), we use a receptive field of 10x10 pixels and a stride of 1. The number of outputs is set to $N_h = 1600$ for fair comparison with the other state-of-the-art methods. We also provide the results of our method with sign split ($N_h = 1600 \times 2$, using \mathbf{W} and $-\mathbf{W}$ for encoding as in (Coates & Ng, 2011)) and using the sparse coding (SC) encoder, which (Coates & Ng, 2011) found to be the best when small number of labeled data is available. For this encoder, we searched over the same set of parameter values as (Coates & Ng, 2011), i.e., $\lambda = \{0.5, 0.75, 1.0, 1.25, 1.5\}$. The parameter λ is tuned to consider the use of sparse coding as encoder after the training and, thus, does not belong to the method that we propose.

Table 2 summarizes the results obtained on this dataset compared to other state-of-the-art methods. When pairing each training method with its associated natural encoding, EPLS outperforms all the other methods. When pairing the training methods with sparse coding, EPLS outperforms the state-of-the-art best performer in single layer networks as well, achieving 61.0% (0.58%) accuracy. More-

over, the standard deviation of the folds is lower than the one provided by OMP-1 with sparse coding encoding. Results are even more impressive if we compare them to meta-parameter free algorithms.

Figure 1 shows a subset of 100 randomly selected bases learned by our method, 10x10 pixel receptive field and a system of $N_h = 1600$ outputs. As shown in the figure, the method learns not only common bases such as oriented edges/ridges in many directions and colors but also corner detectors, tri-banded colored filters, center surrounds and Laplacian of Gaussians among others. This suggests that enforcing lifetime sparsity helps the system to learn a set of complex, rich and diversified bases.

5. Computational complexity

The EPLS algorithm requires the computation of \mathbf{T} , which has $\mathcal{O}(NN_h)$ cost, and therefore scales linearly on both N and N_h . Since we can use vSGD for optimization, the method scales linearly on N given a fixed number of epochs. Finally, applying the activation function, the cost of computing the derivative is linear with N_d , since we use a closed form for $\frac{\partial E}{\partial \mathbf{T}}$.

The memory complexity is related to the mini-batch size N_b . Consequently, the method can scale gracefully to very large datasets: theoretically, it requires to store in memory the mini-batch input data $\mathbf{D}^{(b)}$ ($N_b N_d$ elements), output $\mathbf{H}^{(b)}$ ($N_b N_h$ elements), target $\mathbf{T}^{(b)}$ ($N_b N_h$ elements) and the system parameters to optimize Γ ($N_h (N_d + 1)$ elements); a total amount of $N_h (N_d + 1) + N_b (N_d + 2N_h)$ elements.

6. Discussion

Our results show that simultaneously enforcing both population and lifetime sparsity helps in learning discriminative dictionaries, which reflect in better performance, especially when compared to meta-parameter free methods (Ngiam et al., 2011; Le et al., 2011). Experiments suggest that our algorithm is able to extract features that generalize well on unseen data. When comparing the performance STL-10 dataset, our algorithm outperforms state-of-the-art best performers. Results suggest that our algorithm helps the classifier in generalizing with a few training examples (1% of the dataset), gaining 2% accuracy w.r.t. the state-of-the-art best performer (OMP-1 paired with sparse coding) with a lower standard deviation across folds, suggesting more robustness to variations in the training folds.

It is important to highlight that OMP-1 can be seen as a special case of our algorithm, where the activation function is $|\mathbf{DW}|$ and lifetime sparsity is not taken into account in the optimization process (potentially leading to dead out-

puts). Our algorithm has several advantages over OMP-1: (1) It can use any activation function; (2) by enforcing lifetime sparsity it does not suffer of the dead output problem, thus not requiring ad-hoc tricks to avoid it; (3) it does not require whitening, which can be a problem if the input dimensionality is large (Le et al., 2011).

With our proposal, we advance in the meta-parameter free line of ICA (Hyvärinen & Oja, 2000) and sparse filtering (Ngiam et al., 2011). It is clear that the advantage of sparse filtering over ICA comes from removing the orthogonality constraint, and imposing some sort of “competition” between outputs, which also permits overcomplete representations. Following this spirit, our algorithm imposes an even more strict form of competition to prevent dead outputs by means of Strong Lifetime Sparsity and confirms the trend of (Ngiam et al., 2011; Hyvärinen & Oja, 2000) that data reconstruction seems not so important if the goal is to have a discriminative sparse system.

Last and most importantly, it is worth highlighting five interesting properties of the EPLS algorithm. First, the method is meta-parameter free, which highly simplifies the training process for practitioners, especially when used as a greedy pre-training method in deep architectures. Second, the method is fast and scales linearly with the number of training samples and the input/output dimensionalities. Third, EPLS is easy to implement. We implemented the EPLS in Algorithm 1 in less than 50 lines of C code. The mini-batch vSGD is a general purpose optimizer; our Matlab implementation of vSGD plus the EPLS mex source will be publicly available after publication. Fourth, the proposed learning strategy is not limited to perceptrons. Fifth, there is an interest in the literature in avoiding redundancy in the image representation by using the algorithms in a convolutional fashion (Kavukcuoglu et al., 2010). For this purpose, the EPLS can be slightly modified to apply the procedure to a whole image at once and consider the mini-batch size to be the image divided into patches. This aspect is not considered in the paper and is left for future investigation.

7. Conclusion

In this paper, we introduced the Enforcing Population and Lifetime Sparsity method. The algorithm provides a **meta-parameter free, off-the-shelf, simple and computationally efficient** approach for unsupervised sparse feature learning. It seeks both lifetime and population sparsity in an explicit way in order to learn discriminative features, thus preventing dead outputs.

Results show that the method significantly outperforms all state-of-the-art methods on STL-10 dataset with lower standard deviation across folds, suggesting more robust-

ness across training sets.

References

- Bengio, Yoshua. Learning deep architectures for AI. *Foundations and Trends in Machine Learning*, 2(1):1–127, 2009.
- Bengio, Yoshua, Lamblin, Pascal, Popovici, Dan, and Larochelle, Hugo. Greedy layer-wise training of deep networks. In *NIPS*, pp. 153–160, 2006.
- Bengio, Yoshua, Courville, Aaron C., and Vincent, Pascal. Representation learning: A review and new perspectives. *IEEE TPAMI*, 35(8):1798–1828, 2013.
- Blumensath, T. and Davies, M. E. On the difference between orthogonal matching pursuit and orthogonal least squares. *Unpublished manuscript*, 2007. URL <http://www.see.ed.ac.uk/~tblumens/papers/BDOMPvsOLS07.pdf>.
- Coates, A., Lee, H., and Ng, A. Y. An analysis of single-layer networks in unsupervised feature learning. In *AISTATS*, pp. 214–223, 2011.
- Coates, Adam and Ng, Andrew. The importance of encoding versus training with sparse coding and vector quantization. In *ICML*, pp. 921–928, 2011.
- Erhan, Dumitru, Courville, Aaron, Bengio, Yoshua, and Vincent, Pascal. Why does unsupervised pre-training help deep learning? In *AISTATS*, volume 9, pp. 201–208, May 2010.
- Field, D. J. What is the goal of sensory coding? *Neural Computation*, 6(4):559–601, July 1994.
- Goh, Hanlin, Thome, Nicolas, Cord, Matthieu, and Lim, Joo-Hwee. Unsupervised and supervised visual codes with restricted boltzmann machines. In *ECCV*, pp. 298–311, 2012.
- Hinton, G. E. A Practical Guide to Training Restricted Boltzmann Machines. Technical report, University of Toronto, 2010.
- Hinton, Geoffrey E., Osindero, Simon, and Teh, Yee-Whye. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, July 2006.
- Hyvärinen, A. and Oja, E. Independent component analysis: algorithms and applications. *Neural Networks*, 13(4–5):411–430, 2000.
- Hyvärinen, A., Karhunen, J., and Oja, E. *Independent component analysis*. Wiley Interscience, 2000.

- Kavukcuoglu, Koray, Sermanet, Pierre, Boureau, Y-Lan, Gregor, Karol, Mathieu, Michaël, and LeCun, Yann. Learning convolutional feature hierarchies for visual recognition. In *NIPS*, 2010.
- Kuhn, Harold W. The hungarian method for the assignment problem. *Naval Research Logistics Quarterly*, 2:83–97, 1955.
- Larochelle, Hugo, Bengio, Yoshua, Louradour, Jérôme, and Lamblin, Pascal. Exploring strategies for training deep neural networks. *J. Mach. Learn. Res.*, 10:1–40, June 2009.
- Le, Q. V., Karpenko, A., Ngiam, J., and Ng, A. Y. ICA with reconstruction cost for efficient overcomplete feature learning. In *NIPS*, pp. 1017–1025, 2011.
- LeCun, Yann, Bottou, Leon, Orr, Genevieve, and Müller, Klaus. Efficient backprop. In *Neural Networks: Tricks of the Trade*, pp. 9–50. Springer Berlin, 1998.
- Lee, H., Battle, A., Raina, R., and Ng, A. Y. Efficient sparse coding algorithms. In *NIPS*, pp. 801–808, 2006.
- Lee, H., Ekanadham, C., and Ng, A. Y. Sparse deep belief net model for visual area v2. In *NIPS*, pp. 873–880, 2008.
- Ngiam, J., Koh, P. W., Chen, Z., Bhaskar, S., and Ng, A. Y. Sparse filtering. In *NIPS*, pp. 1125–1133, 2011.
- Olshausen, B. and Field, D. J. Sparse coding with an overcomplete basis set: a strategy employed by v1? *Vision Research*, 37(23):3311–3325, 1997.
- Pati, Y. C., Rezaifar, R., and Krishnaprasad, P. S. Orthogonal matching pursuit: recursive function approximation with applications to wavelet decomposition. In *ACSSC*, pp. 40–44, November 1993.
- Raina, Rajat, Battle, Alexis, Lee, Honglak, Packer, Benjamin, and Ng, Andrew Y. Self-taught learning: Transfer learning from unlabeled data. In *ICML*, pp. 759–766, 2007.
- Ranzato, M. A., Poultney, C., Chopra, S., and Lecun, Y. Efficient learning of sparse representations with an energy-based model. In *NIPS*, pp. 1137–1144, 2006.
- Schaul, Tom, Zhang, Sixin, and LeCun, Yann. No More Pesky Learning Rates. In *ICML*, 2013.
- Snoek, J., Larochelle, H., and Adams, R. P. Practical bayesian optimization of machine learning algorithms. In *NIPS*, pp. 2960–2968, 2012.
- Willmore, B. and Tolhurst, D. J. Characterizing the sparseness of neural codes. *Network*, 12(12):255–270, January 2001.
- Yang, J., Yu, K., Gong, Y., and Huang, T. Linear spatial pyramid matching using sparse coding for image classification. In *IEEE CVPR*, pp. 1794–1801, 2009.
- Zeiler, Matthew D. and Fergus, Rob. Stochastic pooling for regularization of deep convolutional neural networks. *CoRR*, abs/1301.3557, 2013.