

# Efficient online learning of a non-negative sparse autoencoder

Andre Lemme, R. Felix Reinhart and Jochen J. Steil

Research Institute for Cognition and Robotics (CoR-Lab), Bielefeld University  
{alemme, freinhar, jsteil}@CoR-Lab.Uni-Bielefeld.de

**Abstract.** We introduce an efficient online learning mechanism for non-negative sparse coding in autoencoder neural networks. In this paper we compare the novel method to the batch algorithm non-negative matrix factorization with and without sparseness constraint. We show that the efficient autoencoder yields to better sparseness and lower reconstruction errors than the batch algorithms on the MNIST benchmark dataset.

## 1 Introduction

Unsupervised learning techniques that can find filters for relevant parts in images are particularly important for visual object classification [1]. Lee proposed a non-negative matrix factorization (NMF) in [2] to produce a non-negative encoding of input images by combining a linear matrix factorization approach with non-negativity constraints. Hoyer introduced an additional sparseness constraint to the factorization process in order to increase the sparseness of the produced encoding [3]. Sparse codes enhance the probability of linear separability, a crucial feature for object classification [4]. However, the main drawbacks of matrix factorization approaches are high computational costs and the fact that the costly matrix factorization has to be redone each time new images are perceived. Therefore, NMF is not suited for problems that require processing in real-time and life-long learning. A non-negative and online version of the PCA was introduced recently [5]. This approach addresses the problem of non-negativity and computational efficiency, however, PCA is intrinsically a non-sparse method.

We propose a modified autoencoder model that encodes input images in a non-negative and sparse network state. We use logistic activation functions in the hidden layer, which map neural activities to strictly positive outputs. Sparseness of the representation is gained by a unsupervised self-adaptation rule, which adapts the non-linear activation functions based on the principle of intrinsic plasticity [6]. Non-negativity of connection weights is enforced by a novel, asymmetric regularization approach that punishes negative weights more than positive ones. The online learning rules use only local information and are very efficient. We compare the reconstruction performance as well as the generated codes of the introduced model to the non-negative offline algorithms NMF and non-negative matrix factorization with sparseness constraints (NMFSC) on the MNIST benchmark data set. The efficient online implementation outperforms the offline algorithms with respect to reconstruction errors and sparseness of filters as well as code vectors.

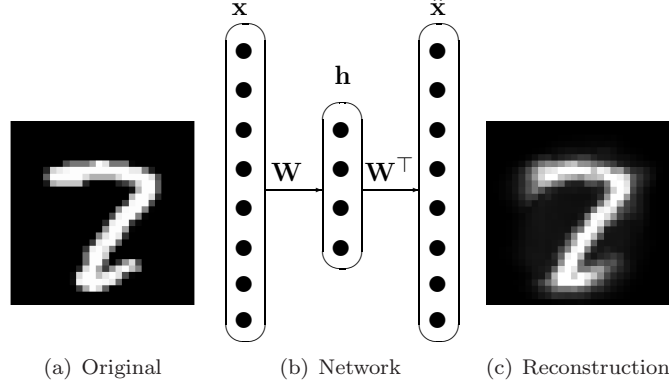


Fig. 1: Autoencoder network with input image  $\mathbf{x}$  (left) and reconstructed image  $\hat{\mathbf{x}}$  (right).  $\mathbf{W} \equiv \mathbf{W}^T$  is the tied weight matrix and  $\mathbf{h}$  the hidden network state.

## 2 Non-negative sparse autoencoder network (NNSAE)

We modify a autoencoder network in order to obtain non-negative, sparse encodings with only positive network weights. The network architecture is shown in Fig. 1. We denote the input image by  $\mathbf{x}$  and the network reconstruction by  $\hat{\mathbf{x}} = \mathbf{W}^T \mathbf{f}(\mathbf{W}\mathbf{x})$ , where  $\mathbf{W}$  is the weight matrix and  $\mathbf{f}(\cdot)$  are parameterized activation functions  $f_i(h_i) = (1 + \exp(-a_i h_i - b_i))^{-1} \in [0, 1]$  with slopes  $a_i$  and biases  $b_i$  that are applied component-wise to the neural activities  $\mathbf{h} = \mathbf{W}\mathbf{x}$ . Non-negative code vectors are already assured by the logistic activation functions.

**Learning to reconstruct:** Learning of the autoencoder is based on minimizing the reconstruction error  $\|\mathbf{x} - \hat{\mathbf{x}}\|^2$  on the training set. Note that  $\mathbf{W} \equiv \mathbf{W}^T$  is the same matrix for encoding and decoding, i.e. the autoencoder has *tied* weights [7, 8], which reduces the parameters to adapt. Learning of the autoencoder boils down to a simple error correction rule for all training examples  $\mathbf{x}$

$$\Delta w_{ij} = \eta (x_i - \hat{x}_i) h_j + d(w_{ij}), \quad (1)$$

where  $w_{ij}$  is the connection from hidden neuron  $j$  to the output neuron  $i$ . The decay function  $d(w_{ij})$  will be discussed in the next section. We use an adaptive learning rate  $\eta = \tilde{\eta} (\|\mathbf{h}\|^2 + \epsilon)^{-1}$  that scales the gradient step width  $\tilde{\eta}$  by the inverse network activity and is similar to the backpropagation-decorrelation learning rule introduced for reservoir networks in [9].

**Asymmetric decay enforces non-negative network parameter:** In order to enforce non-negative weights, we introduce an asymmetric regularization approach. Typically, a Gaussian prior for the network weights is assumed, which is expressed by a quadratic cost term added to the error function. Using gradient descent, the Gaussian prior results in a so called decay term  $d(w_{ij}) = -\lambda w_{ij}$ . We assume a virtually deformed Gaussian prior that is skewed with respect to the sign of the weight  $w_{ij}$ . Gradient descent then yields the following asymmetric,

piecewise linear decay function:

$$d(w_{ij}) = \begin{cases} -\alpha w_{ij} & \text{if } w_{ij} < 0 \\ -\beta w_{ij} & \text{else} \end{cases} \quad (2)$$

The main advantage of (2) is parameterized decay behavior depending on the sign of the weight: it is possible to allow a certain degree of negative weights ( $0 < \alpha \ll 1$ ) or to prohibit negative weights completely ( $\alpha = 1$ ). The decay function falls back to the case of a symmetric prior for  $\alpha = \beta$ .

**Adapting non-linearities to gain sparse codes:** In the context of the non-negative autoencoder, it is crucial to adapt the neuron’s activation functions to the strictly non-negative input distribution caused by the non-negative weight and data matrix. We therefore apply an efficient, unsupervised learning scheme to the parameterized activation functions  $f_i(a_i)$  that is known as intrinsic plasticity (IP) [6]. In addition, IP enhances the sparseness of the input representation in the hidden layer [10]. For more details about IP see [6, 10].

We use a small learning rate  $\eta_{IP} = 0.0001$  throughout the experiments and initialize the parameters of the activation functions  $\mathbf{a} = \mathbf{1}$  and  $\mathbf{b} = -\mathbf{3}$  in order to accelerate convergence. IP provides a parameter to adjust the desired mean network activation, which we set to  $\mu = 0.2$  corresponding to a sparse network state. In addition, we use a constant bias decay  $\Delta \mathbf{b} = \Delta \mathbf{b}_{IP} - 0.00001$  to further increase sparseness of the encoding.

### 3 Encoding and decoding of handwritten digits

We analyze the behavior of the non-negative sparse autoencoder (NNSAE) and compare the new NNSAE approach to the batch algorithms NMF and NMFSC on a benchmark data set.

**Dataset and network training:** The MNIST data set is commonly used to benchmark image reconstruction and classification methods [11]. For the following experiments the “MNIST-basic” set [12] is used, which comprises 12000 trainings and 50000 test images of handwritten digits from 0 to 9 in a centered and normalized  $28 \times 28$  pixel format. We subdivide the training set into 10000 samples for learning and use the remaining 2000 examples as validation set and we use all test images in the test set. Some example digits from the test set are shown in the top row of Fig. 2.

We use an auto-encoder with 784 input neurons and 100 hidden neurons to encode and decode the input images. The weight matrix is initialized randomly according to a uniform distribution in  $[0.95, 1.05]$ . We set  $\tilde{\eta} = 0.01$  and  $\epsilon = 0.002$  in (1). Each auto-encoder is trained for 100 epochs, where the whole training set is presented to the model per epoch. To account for the random network initialization, we present all results averaged over 10 trials.

**Impact of asymmetric decay:** We first investigate the impact of the asymmetric decay function (2). The following variates are calculated to quantify the

model properties depending on  $\alpha$ : (a) pixelwise mean square error

$$MSE = \frac{1}{MN} \sum_{i=0}^N \sum_{j=0}^M (x_j^i - \hat{x}_j^i)^2 \quad (3)$$

of the reconstructions, where  $M$  is the dimension of the data and  $N$  the number of samples. (b) average sparseness of the code matrix  $\mathbf{H}$  and filters  $\mathbf{w}_i$  which we estimate by

$$s(\mathbf{x}) = (\sqrt{n} - (\sum |x_i|) / \sqrt{\sum x_i^2}) (\sqrt{n} - 1)^{-1}, \quad (4)$$

where  $n$  is the number of dimensions of  $\mathbf{x}$ . This function was introduced by Hoyer [13] and rates the energy in a vector  $\mathbf{x}$  with values between zero and one, where sparse vectors  $\mathbf{x}$  are mapped to  $sp(\mathbf{x}) \gg 0$ .

	$\alpha=0$	$\alpha=0.1$	$\alpha=0.2$	$\alpha=0.3$	$\alpha=1.0$
$MSE$	0.0107	0.0150	0.0150	0.0150	0.0150
$STD$ 1.0e-04 *	0.5966	0.4079	0.7153	0.7143	0.6893
$s(\mathbf{H})$	0.39	0.35	0.35	0.35	0.35
$s(\mathbf{W})$	0.53	0.93	0.93	0.93	0.93

Tab. 1:  $MSE$ , standard deviation of  $MSE$  and sparseness of the encodings  $\mathbf{H}$  as well as network weights  $\mathbf{W}$  depending on the asymmetric decay rate  $\alpha$ .

Tab. 1 shows the  $MSE$  and the sparseness of hidden states as well as connection weights depending on  $\alpha$ . The results in Tab. 1 show that a certain amount of negative entries seem to be useful in terms of the reconstruction error. However, if we commit to non-negativity it does not matter whether there are just a small set of negative or only positive entries in the weight matrix. This indicates that a rather moderate non-negativity constraint causes the weight dynamics and the final solution to change completely. We therefore set  $\alpha=1$  in the following experiments.

#### 4 Online NNSAE versus offline NMFSC

The batch algorithms NMF and NMFSC solve the equation  $\mathbf{X} = \mathbf{WH}$  starting from the data matrix  $\mathbf{X}$  and initial filter and code matrices  $\mathbf{W}$  and  $\mathbf{H}$ . Then, iterative updates the code and weight matrix are conducted in order to reduce the reconstruction error. We use 100 iteration steps to make the comparison with the NNSAE fair. We further set the sparseness parameters provided by NMFSC for the weight and code matrix to  $s(\mathbf{W}) = 0.9$  and  $s(\mathbf{H}) = 0.35$  respectively, according to the values obtained for the NNSAE. For more details about NMF and NMFSC see [2, 13]. To test the NMF(SC)  $\mathbf{X}$  denotes the test set,  $\mathbf{W}$  contains the trained filters and is set fix. The code matrix  $\mathbf{H}$  is iteratively updated the same way as in the training.

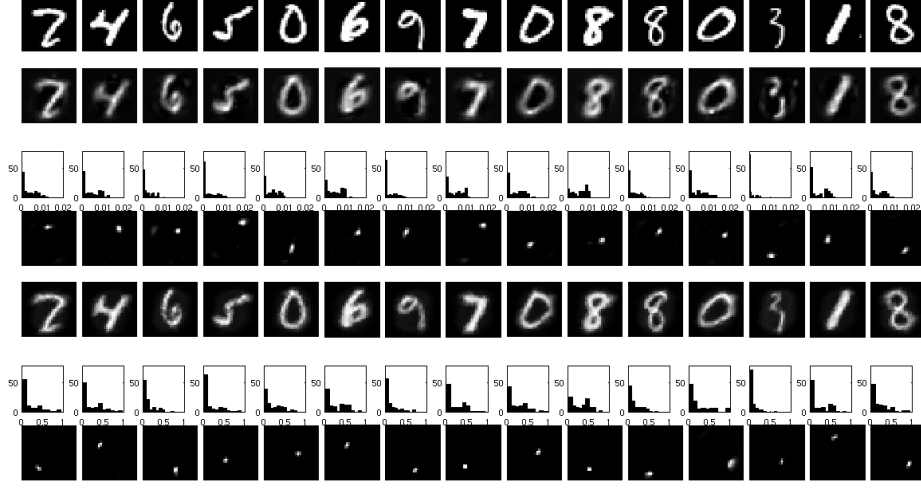


Fig. 2: 15 images from the test set (top row) with reconstructions (2nd row), code histogram (3rd row) and some filters (4rd row) performed by the NMFSC. Reconstruction (5rd row), code histogram (6rd row) and some filters (7rd row) performed by the NNSAE.

**Comparison of reconstruction errors and sparseness:** Tab. 2 gives a review of the performance of the NNSAE and compares it to NMF/NMFSC. The NNSAE outperforms NMF/NMFSC with respect to the  $MSE$  on the test set, sparseness of the code matrix  $\mathbf{H}$  and sparseness of the weight matrix  $\mathbf{W}$ . It is surprising that even NMFSC does not reach the sparse codes that we achieve with the NNSAE, although the constraint is set to be equal ( $s(\mathbf{H})=0.35$ ). We show for each example the corresponding code histogram of NMFSC and NNSAE in Fig. 2. The sparseness of the basis vectors  $\mathbf{w}_i$  does not differ significantly for the different methods, which seems to be the effect of the non-negative constraint itself. In Fig. 2 we show in the 4th and the last row some filters of NMFSC and NNSAE. The filters of both methods cover only blob like areas, i.e. are local feature detectors, which corresponds to a sparse representation.

	$MSE$			Sparseness	
	Test	Validation	Training	code $\mathbf{H}$	filter $\mathbf{W}$
NMF	0.0473	0.2220	0.0113	0.2617	0.8867
NMFSC	0.1089	0.1079	0.0143	0.2852	0.9000
NNSAE	0.0150	0.0156	0.0124	0.35	0.93

Tab. 2: Mean square reconstruction errors for training, validation and test set as well as sparseness of code and weight matrices for NMFSC, NMF and NNSAE.

## 5 Discussion and conclusion

We present an autoencoder network for efficient online learning of sparse and non-negative encodings. We therefore combine a mechanism to enhance the sparseness of encodings and an asymmetric decay function to create non-negative weights. The reconstruction performance and code sparseness are compared to the offline techniques NMF and NMFSC. The online trained non-negative autoencoder outperforms the offline methods when reconstructing images from the test set, underlining the generalization capabilities of the autoencoder. Additionally, the autoencoder produces more sparse codes compared to the offline methods while having similar filters at the same time. Besides the life-long learning capability of the autoencoder, the main advantage of the proposed approach is the efficient encoding of novel inputs: the state of the network has simply to be updated with the new input, where the matrix factorization approaches in contrast require a full optimization step on the new data.

It is of particular interest to use the NNSAE in a stacked auto-encoder network for pattern recognition in the future. Does the non-negative representation improve classification performance, and how is fine-tuning of the hierarchy, e.g. by backpropagation learning, affected by the asymmetric decay function?

## References

- [1] M.W.Spratling. Learning Image Components for Object Recognition. *Machine Learning*, 2006.
- [2] Lee, D. D. and Seung, H. S. Learning the parts of objects by nonnegative matrix factorization. *Nature*, pp. 788–791, 1999.
- [3] Hoyer, Patrik O and Centre, Neural Networks Research. Non-negative sparse coding. *Neural Networks for Signal Processing XII*, pp. 557–565, 2002.
- [4] Y-lan Boureau and Yann Lecun. Sparse Feature Learning for Deep Belief Networks. *NIPS*, pp. 1–8, 2007.
- [5] Mark D. Plumbley and Erkki Oja. A "nonnegative PCA" algorithm for independent component analysis. *Neural Networks*, pp. 66–76, 2004.
- [6] Jochen Triesch. A Gradient Rule for the Plasticity of a Neuron's Intrinsic Excitability. *Neural Computation*, pp. 65–70, 2005.
- [7] Hinton, Geoffrey E and Osindero, Simon and Teh, Yee-Whye A fast learning algorithm for deep belief nets. *Neural Computation*, pp. 1527–54, 2006.
- [8] Vincent, Pascal and Larochelle, Hugo and Bengio, Yoshua and Manzagol, Pierre-Antoine Extracting and composing robust features with denoising autoencoders *ICML '08*, pp. 1096–1103, 2008.
- [9] Steil, J. J. Backpropagation-Decorrelation: online recurrent learning with O(N) complexity *Proc. IJCNN*, pp. 843–848, 2004.
- [10] Jochen Triesch. Synergies Between Intrinsic and Synaptic Plasticity Mechanisms. *Neural Computation*, pp. 885–909, 2007.
- [11] Yann Lecun and Corinna Cortes. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist>.
- [12] Variations on MNIST The MNIST database of handwritten digits. <http://www.iro.umontreal.ca/lisa/ptwiki/>.
- [13] Hoyer, Patrik O Non-negative Matrix Factorization with Sparseness Constraints. *Journal of Machine Learning Research*, pp. 1457–1469, 2004.