

A Hebbian/Anti-Hebbian Network for Online Sparse Dictionary Learning Derived from Symmetric Matrix Factorization

Tao Hu¹, Cengiz Pehlevan^{2,3}, and Dmitri B. Chklovskii³

¹Texas A&M University
MS 3128 TAMUS
College Station, TX 77843
taohu@tees.tamus.edu

²Janelia Farm Research Campus
Howard Hughes Medical Institute
Ashburn, VA 20147
pehlevanc@janelia.hhmi.org

³Simons Center for Data Analysis
Simons Foundation
New York, NY 10010
mitya@simonsfoundation.org

Abstract—Olshausen and Field (OF) proposed that neural computations in the primary visual cortex (V1) can be partially modelled by sparse dictionary learning. By minimizing the regularized representation error they derived an online algorithm, which learns Gabor-filter receptive fields from a natural image ensemble in agreement with physiological experiments. Whereas the OF algorithm can be mapped onto the dynamics and synaptic plasticity in a single-layer neural network, the derived learning rule is nonlocal - the synaptic weight update depends on the activity of neurons other than just pre- and postsynaptic ones - and hence biologically implausible. Here, to overcome this problem, we derive sparse dictionary learning from a novel cost-function - a regularized error of the symmetric factorization of the input's similarity matrix. Our algorithm maps onto a neural network of the same architecture as OF but using only biologically plausible local learning rules. When trained on natural images our network learns Gabor-filter receptive fields and reproduces the correlation among synaptic weights hard-wired in the OF network. Therefore, online symmetric matrix factorization may serve as an algorithmic theory of neural computation.

Keywords—*sparse dictionary learning; neuron; online algorithm; matrix factorization; neuromorphic computing*

I. INTRODUCTION

In the quest to understand neural computation in mammals, the primary visual cortex (V1) has been an attractive and well-studied target system [1]. One of its major tasks is computing orientationally selective responses, or Gabor-filter receptive fields, out of orientationally nonselective inputs [2]. Such computation has been successfully modeled by Olshausen and Field (OF) who proposed a neural network that learns Gabor-filter receptive fields from an ensemble of natural images in an unsupervised fashion [3,4]. The OF network appeals as a model of V1 because it is both rigorously derived from a principled cost function and captures several salient anatomical and physiological features of V1 networks [5].

However, there remains an unanswered question regarding modeling neural computation in V1 by the OF algorithm. Whereas the original two-layer neural network

implementation of the OF algorithm [3,4] may model sensory periphery [11,12] the required symmetric feedback connections have not been observed in V1. At the same time, in the single-layer network implementation of the OF algorithm [6] appropriate for V1, the learning rule derived from the OF cost function is nonlocal - the synaptic weight update depends on the activity of neurons other than just pre- and postsynaptic ones - and therefore biologically implausible.

In this paper, we propose a novel cost function and demonstrate that from it one can derive neuronal dynamics and local learning rules, both Hebbian for feedforward and anti-Hebbian for lateral synaptic connections. We demonstrate that training the network on a natural image ensemble yields Gabor-filter receptive fields. We also demonstrate that the application of such rules yields lateral connection weights that obey the same relationship with feedforward weights as in the OF framework. In addition, our framework accounts for several salient properties of biological networks and predicts that the learning rate decays with time in an activity-dependent fashion agreeing with experiments. Therefore, we make a step towards understanding V1 and mammalian neural computation in general.

The paper is organized as follows. In the next Section we summarize the OF algorithm and its neural network implementation. In the Results: A) we present the new cost function, the regularized error squared between the input's and the output's similarity matrices, and a derivation of an online algorithm for sparse dictionary learning with local learning rules; B) we report the results of numerical simulations showing that our network performs similarly to OF; C) we derive analytically the observed relationship between feedforward and lateral synaptic connection weights in our network, which reproduces a hard-wired constraint in the OF network; D) we show that our online symmetric matrix factorization algorithm can discover independent components in the whitened input data. In the Discussion: A) we compare our model to biology; B) we suggest that matrix factorization may be a generic model of neural computation.

II. THE OLSHAUSEN-FIELD (OF) ALGORITHM

To motivate our work we briefly review the OF model [3,4] and point out the biologically implausible aspect of the single-layer implementation. The starting point of the OF model is the assumption that the vectorized image patches, $\mathbf{x}_t \in \mathbb{R}^n$, are represented by the neuronal feature vectors, i.e. columns of an overcomplete ($m > n$) dictionary, $\mathbf{W} \in \mathbb{R}^{n \times m}$, weighted by a sparse vector of neuronal activities, $\mathbf{y}_t \in \mathbb{R}^m$. To obtain such representation the OF model minimizes the squared representation error regularized by the l_1 -norm of activity:

$$\min_{\mathbf{W}} \sum_t \min_{\mathbf{y}_t} \left(\frac{1}{2} \|\mathbf{x}_t - \mathbf{W}\mathbf{y}_t\|_2^2 + \lambda \|\mathbf{y}_t\|_1 \right), \quad (1)$$

where λ reflects the relative importance of sparsity and representation accuracy.

To derive a neural network algorithm, OF minimized (1) in response to sequentially presented natural image patches, a so-called online setting. Specifically, for each presented image, \mathbf{x}_t , they i) find the optimal value of \mathbf{y}_t for fixed \mathbf{W} , ii) for fixed \mathbf{y}_t perform stochastic gradient descent with respect to feature vectors, \mathbf{W} . Next, we discuss these two steps in more detail.

i) To find \mathbf{y}_t for each image the algorithm minimizes (1) using stochastic (sub)gradient descent steps [7] with respect to \mathbf{y}_t :

$$\begin{cases} \mathbf{c}_t = \mathbf{W}'_t \mathbf{x}_t - \mathbf{W}'_t \mathbf{W}_t \mathbf{y}_{t-1} \\ \mathbf{y}_t = \text{ST}(\mathbf{c}_t, \lambda) \end{cases}, \quad (2)$$

where ST is a component-wise soft-threshold function [8], see Fig. 1A. Equation (2) can be viewed as dynamics of activity in a single-layer network with feedforward and lateral connections [6], Fig. 1B. Then, \mathbf{c}_t represents the total input currents into neurons and soft thresholding models a rectifying nonlinearity of a biological neuron. In such network, lateral connections implement “explaining away”, or competition between neurons in representing an input signal.

ii) After the network activity \mathbf{y}_t converges to a representation of an image, the algorithm updates feature vectors, \mathbf{W} :

$$\mathbf{W}_{t+1,i,j} = \mathbf{W}_{t,i,j} + \delta \left(x_{t,i} - \sum_k \mathbf{W}_{t,i,k} y_{t,k} \right) y_{t,j}, \quad (3)$$

where $\mathbf{W}_{t+1,i,j}$ can be viewed as the synaptic weight for the connection from neuron j to i , $\delta > 0$ is the learning rate.

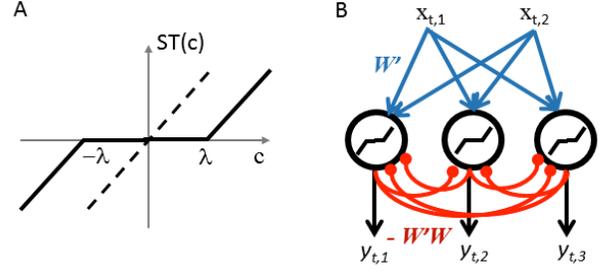


Fig. 1: A neural network implementation of the OF algorithm. A) Soft Thresholding (ST) function. B) A single-layer OF network. Each neuron applies ST on the inputs weighted by the feedforward connections, \mathbf{W}' , minus outputs weighted by the lateral connections, $\mathbf{W}'\mathbf{W}$. Connection weights are updated using nonlocal learning rules.

The OF model (2,3) successfully reproduces several salient features of the primary visual cortex (V1) anatomy and physiology such as the overcompleteness of cortical representation, sparsity of neural activity, nonlinearity of neural responses [4,5]. Perhaps, most impressively, the receptive fields (computed from the feature vectors and whitening matrix) learned by the network on the ensemble of whitened natural images are Gabor-filter patches resembling receptive fields of neurons in V1 [3,4].

However, a major problem with modeling V1 with the OF algorithm is that in the single-layer network implementation [6] the learning rules are nonlocal. Specifically, the proposed learning rule (3) requires that each synapse “knows” the weights of synapses belonging to neurons other than its pre- and postsynaptic neuron. Because no mechanism exists for such communication in the brain it is not clear how the OF model can describe learning in V1. In addition, lateral connection weights in the OF model (2) are not learned directly but computed from the feedforward connection weights, i.e. the lateral connection matrix \mathbf{M} satisfies

$$\mathbf{M} = -\mathbf{W}'\mathbf{W}. \quad (4)$$

Previously, this problem was addressed by a network of OF architecture but with local learning rules: Hebbian for feedforward and anti-Hebbian for lateral connections [9,10]. However, such local learning rules have been postulated rather than derived from any cost function.

III. RESULTS

Here, we derive a single-layer network for sparse overcomplete representation by minimizing the cost function comprising the squared difference between the similarity matrices of the input and the output data and a sparsity-inducing regularizer. Next, we demonstrate that this network learns Gabor patch receptive fields when trained on a natural image ensemble. Furthermore, we show that the relationship between lateral and feedforward connection weights agrees with that hard-wired into the OF network. Interestingly, our

framework also predicts the decay of learning rate with time as observed experimentally.

A. Cost function and derivation of the algorithm

We start by introducing a data matrix notation for algorithm input:

$$\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_T) = \begin{pmatrix} x_{1,1} & \cdots & x_{T,1} \\ \vdots & \ddots & \vdots \\ x_{1,n} & \cdots & x_{T,n} \end{pmatrix}, \quad (5)$$

and for algorithm output:

$$\mathbf{Y} = (\mathbf{y}_1 \cdots \mathbf{y}_T) = \begin{pmatrix} y_{1,1} & \cdots & y_{T,1} \\ \vdots & \ddots & \vdots \\ y_{1,m} & \cdots & y_{T,m} \end{pmatrix} = \begin{pmatrix} \mathbf{y}_{\cdot,1} \\ \vdots \\ \mathbf{y}_{\cdot,m} \end{pmatrix}. \quad (6)$$

We denote a transpose of matrix \mathbf{A} as \mathbf{A}' and its Frobenius norm as $\|\mathbf{A}\|_F$.

We propose to model the online sparse dictionary learning by minimizing the following cost function:

$$\mathbf{y}_T = \arg \min_{\mathbf{y}_T} \|\mathbf{X}'\mathbf{X} - \mathbf{Y}'\mathbf{Y}\|_F^2 + \lambda \sum_i \|\mathbf{y}'_{\cdot,i} \mathbf{y}_{\cdot,i}\|_1 \quad (7)$$

where $\mathbf{y}_{\cdot,i}$ is an i -th row of matrix \mathbf{Y} (6), the activity of i -th output channel. The same loss term without the regularizer has been used previously, in the offline setting, in multi-dimensional scaling [13] and in symmetric nonnegative matrix factorization, where \mathbf{Y} is constrained to be element-wise nonnegative [14]. Whereas the regularizer may not look familiar, it induces sparsity on the outer product of rows of \mathbf{Y} and hence the activity of output channels. The motivation for choosing the particular form of the sparsity inducing regularizer will become clear below.

Let us derive an online algorithm temporarily ignoring the regularizer in (7). Such minimization problem can be solved by taking a derivative with respect to \mathbf{Y} and setting it to zero:

$$\begin{aligned} \left[\frac{\partial}{\partial \mathbf{Y}} \|\mathbf{X}'\mathbf{X} - \mathbf{Y}'\mathbf{Y}\|_F^2 \right]_{T,\bullet} &= \\ &= \left[\frac{\partial}{\partial \mathbf{Y}} \text{Tr}(\mathbf{Y}'\mathbf{Y}'\mathbf{Y} - 2\mathbf{X}'\mathbf{X}\mathbf{Y}'\mathbf{Y}) \right]_{T,\bullet} = \\ &= [4(\mathbf{Y}\mathbf{Y}'\mathbf{Y} - \mathbf{Y}\mathbf{X}'\mathbf{X})]_{T,\bullet} = 0, \end{aligned} \quad (8)$$

where the subscript T,\bullet denotes the T -th column. When $T > m$, the products $\mathbf{Y}\mathbf{Y}'$ and $\mathbf{Y}\mathbf{X}'$ change slowly with time and can be approximated by the matrices $\tilde{\mathbf{M}}_T$ and $\tilde{\mathbf{W}}_T'$ computed on the data available before the presentation of T -th sample. Then (8) can be linearized:

$$\tilde{\mathbf{M}}_T \mathbf{y}_T = \tilde{\mathbf{W}}_T' \mathbf{x}_T, \quad (9)$$

This linear system can be solved by coordinate descent (to avoid matrix division) leading to the following dynamics of neuronal activity:

$$y_{T,i} \leftarrow \mathbf{W}_{T,i}' \mathbf{x}_T - \mathbf{M}_{T,i} \mathbf{y}_T,$$

where

$$\mathbf{W}_{T,j,i} = \frac{\sum_{t=1}^{T-1} y_{t,i} x_{t,j}}{\sum_{t=1}^{T-1} y_{t,i}^2}; \quad \mathbf{M}_{T,i,j \neq i} = \frac{\sum_{t=1}^{T-1} y_{t,i} y_{t,j}}{\sum_{t=1}^{T-1} y_{t,i}^2}; \quad \mathbf{M}_{T,i,i} = 0. \quad (10)$$

These expressions lead to a natural single-layer network implementation of the algorithm, Fig 2A, where matrices \mathbf{W} and \mathbf{M} correspond to feedforward and lateral synaptic connection weights correspondingly. Interestingly, although the synaptic weights did not appear explicitly in the cost function (7), they arise naturally in the online minimization algorithm (10).

Importantly, unlike in the single-layer neural network implementation of the OF model (3), here, the expressions for the synaptic weights are local, i.e. depend on the activities of

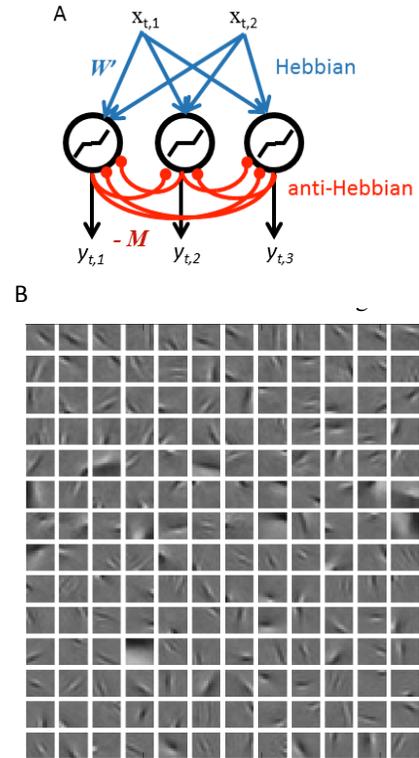


Fig. 2: A neural network implementation of the sparse matrix factorization algorithm. A) A single-layer network with local learning rules. Each neuron applies Soft Thresholding (ST) to the inputs weighted by the feedforward connections, \mathbf{W}' , minus outputs weighted by the lateral connections, \mathbf{M} . Connection weights are updated using Hebbian and anti-Hebbian learning rules correspondingly. B) Receptive fields learned on the whitened natural image ensemble.

only pre- and postsynaptic neurons, Fig 2A.

To avoid storing past input and output activity appearing in the sums (10) we rewrite learning rules in a recursive form that admits online implementation:

$$\begin{aligned}\hat{Y}_{T,i} &= \hat{Y}_{T-1,i} + y_{T-1,i}^2, \\ W_{T,j,i} &= W_{T-1,j,i} + y_{T-1,i} (x_{T-1,j} - W_{T-1,j,i} y_{T-1,i}) / \hat{Y}_{T,i}, \\ M_{T,i,j \neq i} &= M_{T-1,i,j} + y_{T-1,i} (y_{T-1,j} - M_{T-1,i,j} y_{T-1,i}) / \hat{Y}_{T,i}.\end{aligned}\quad (11)$$

Thus, the feedforward synaptic weights are updated according to the Oja's modification of the Hebb rule [15] with the activity dependent learning rate. To the best of our knowledge such single-neuron learning rule [16] has not been previously derived for the multi-neuron case. Moreover, for the first time, we were able to derive the Oja-like version of the anti-Hebbian rule, see also [17,18].

Including the regularizer in the cost function alters the derivation in that instead of the derivative one needs to take a sub-derivative [7]. This does not affect the learning rules but adds soft thresholding [8] of the inputs to the dynamics:

$$y_{T,i} \leftarrow \text{ST} \left(W_{T,i} x_T - M_{T,i} y_T, \eta_{T,i} \right), \quad (12)$$

where the threshold is:

$$\eta_{T,i} \equiv \frac{\lambda \sum_{t=1}^{T-1} |y_{t,i}|}{2 \sum_{t=1}^{T-1} y_{t,i}^2}. \quad (13)$$

Now, our motivation for the choice of the regularizer in (7) should become clear: the regularizer was chosen in order to preserve the magnitude of the threshold with time. When output activity is binary, 0 or 1, as in a spiking neuron, the threshold stays exactly the same. When output activity is real, corresponding to the firing rate model or graded potential neurons, the constancy is only approximate but has been confirmed by numerical simulations.

Thus, we derived an online algorithm that can be implemented by a single-layer network with OF architecture relying only on local learning rules. Next we simulate our algorithm numerically by training it on the ensemble of natural images.

B. Numerical simulations

We applied our algorithm (11,12) to a natural image ensemble. Specifically, 10^4 12×12 pixel patches randomly extracted from natural images [19] and whitened. The extracted principal components were presented sequentially to a network of 196 neuron with feedforward and lateral connections, Fig. 2A. While each patch was presented, the coordinate descent update (12) was repeated 50 times for each

neuron. Therefore, we simulate the neural dynamics with a total of 5×10^5 iterations. We initialize the network connection weights with Gaussian random variables and output activity with zeros. We set the initial synaptic learning rate to be $1 / \hat{Y}_{1,i} = 10^{-4}$ and the initial firing threshold to be $\eta_{1,i} = 1.0$.

As a result of training, the network learns the feedforward weight matrix, \mathbf{W}' . To plot neural filters (or receptive fields) acting on natural image patches, we right-multiply \mathbf{W}' by the whitening matrix \mathbf{Q} and plot the rows, Fig. 2B. One can see that the receptive fields have the appearance of Gabors filters of varying orientation and spatial frequency. We fit the receptive fields with 2D Gabor functions:

$$G(\tilde{x}, \tilde{y}) = g \exp\left(-\tilde{x}^2 / 2\sigma_x^2 - \tilde{y}^2 / 2\sigma_y^2\right) \cos(2\pi f\tilde{x} + \varphi),$$

where $\tilde{x} = (x - x_0) \cos \theta + (y - y_0) \sin \theta$ and $\tilde{y} = -(x - x_0) \sin \theta + (y - y_0) \cos \theta$ are obtained by a translation of the original coordinate system (x_0, y_0) followed by a rotation by angle θ . In this equation, g is the amplitude, σ_x and σ_y represent the widths of the Gaussian envelope, f is the spatial frequency of the sinusoidal grating, and φ is phase offset. We present the measured distribution of spatial frequencies and orientation in Fig. 3A and B respectively. Both statistics were similar to that in the OF network [19] and in physiological measurements [28]. Furthermore, the distribution of output activity, \mathbf{y} , has a strong peak at zero (sparsity) and a heavy tail, Fig. 3C.

Finally, we found that the feedforward and lateral connection weights are strongly correlated, Fig. 4. Whereas in the OF network such correlation, (4), is predetermined by the algorithm (2,3), in our network it appeared as a result of independently acting learning rules.

C. Derivation of the relationship between feedforward and lateral connections

In this Section we present an analytical derivation of the relationship between connection matrices \mathbf{W} and \mathbf{M} in the steady state solution of the sparse symmetric matrix factorization cost function. Because the dictionary is overcomplete, when the regularization constant, λ , is not too large, the steady state solution satisfies approximately:

$$\mathbf{X}'\mathbf{X} = \mathbf{Y}\mathbf{Y}' \quad (14)$$

The SVD of the data matrix \mathbf{X} can be written in a standard form:

$$\mathbf{X} = \mathbf{U}_x \Sigma_x \mathbf{V}_x', \quad (15)$$

where as usual singular vectors are orthonormal, $\mathbf{U}_x' \mathbf{U}_x = \mathbf{I}$, $\mathbf{V}_x' \mathbf{V}_x = \mathbf{I}$ and Σ_x is a diagonal matrix. Similarly,

$$Y = U_Y \Sigma_Y V_Y', \quad (16)$$

Next, we substitute (15,16) into (14):

$$V_X \Sigma_X U_X' U_X \Sigma_X V_X' = V_Y \Sigma_Y U_Y' U_Y \Sigma_Y V_Y'.$$

By taking into account the orthonormality of the singular vectors:

$$V_X \Sigma_X \Sigma_X V_X' = V_Y \Sigma_Y \Sigma_Y V_Y'$$

From this we conclude that the right singular vectors of X and Y are equal, and Σ_X and Σ_Y share the same nonzero diagonal values:

$$V_X = V_Y = V, \text{ and } \Sigma_Y \Sigma_Y = \Sigma_X \Sigma_X. \quad (17)$$

Then, the unnormalized connection weight matrices for feedforward and lateral connections:

$$\tilde{W}' = YX' = U_Y \Sigma_Y V' V \Sigma_X U_X' = U_Y \Sigma_Y \Sigma_X U_X', \quad (18)$$

$$\tilde{M} = YY' = U_Y \Sigma_Y V' V \Sigma_Y U_Y' = U_Y \Sigma_Y \Sigma_Y U_Y'. \quad (19)$$

Note that

$$\begin{aligned} \tilde{W}' \tilde{W} &= U_Y \Sigma_Y \Sigma_X \Sigma_X \Sigma_Y U_Y' \\ &= U_Y \Sigma_Y \Sigma_Y \Sigma_Y \Sigma_Y U_Y' = U_Y (\Sigma_Y \Sigma_Y)^2 U_Y'. \end{aligned} \quad (20)$$

If the input matrix is properly whitened Σ_X contains only 1's or 0's on the diagonal. Since Σ_X and Σ_Y share the same nonzero diagonal values, Σ_Y also contains only 1's or 0's on the diagonal. Therefore, (19) and (20) are identical establishing a relationship between feedforward and lateral connection weights.

To obtain synaptic connection weights, W' and M , one has to normalize \tilde{W}' and \tilde{M} by the cumulative postsynaptic

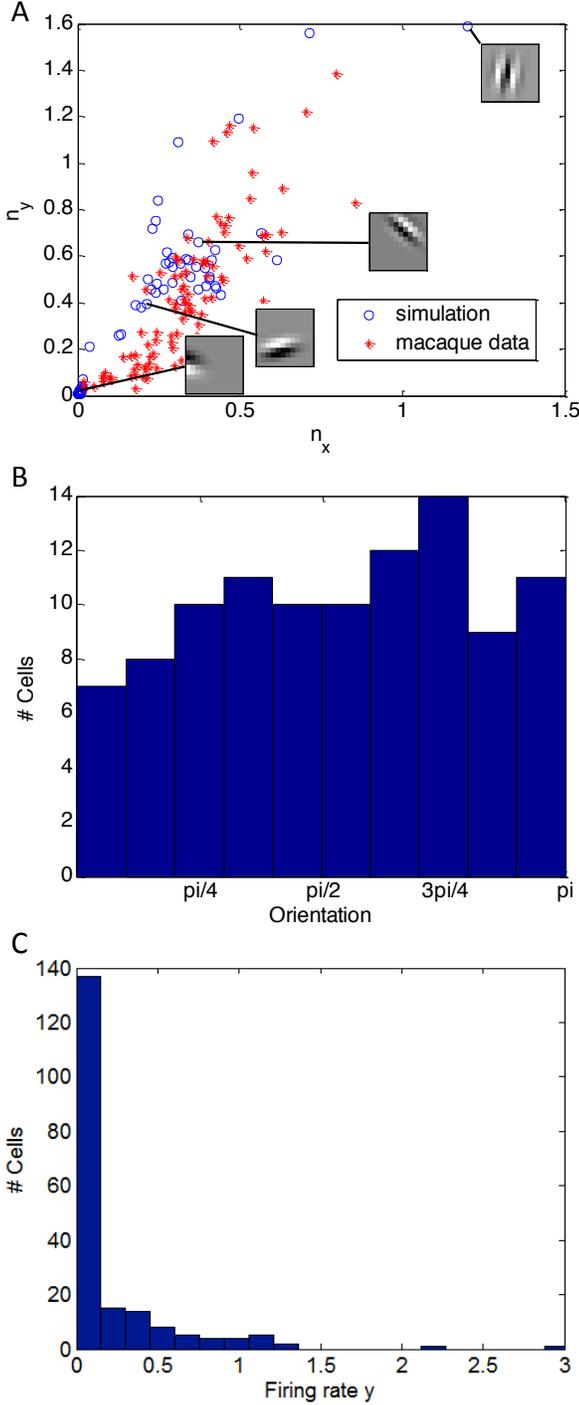


Fig. 3: Statistics of receptive fields and neuronal activity computed in our network matches that of OF model [19] and mammalian physiology [20, 28]. A. Spatial frequencies of Gabor fits, where $n_x = f\sigma_x$ and $n_y = f\sigma_y$. B. The distribution of orientation preference, θ . C. The distribution of activity, γ , among output units is sparse and heavy-tailed.

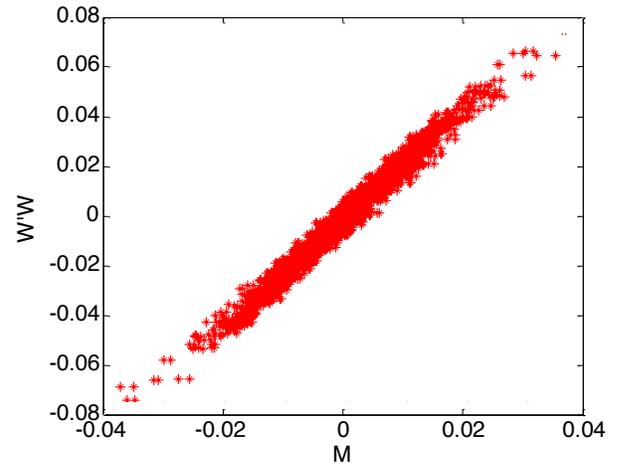


Fig. 4: Correlation between the lateral connection weights and the Gram matrix of the feedforward connections (off-diagonal elements only) in the sparse matrix factorization network.

activity (10). The normalization accounts for the variation in slope in Fig.4.

D. Sparse symmetric matrix factorization can discover independent components

Here we argue that symmetric matrix factorization can be used to discover independent components in their whitened mixture, i.e. perform independent component analysis (ICA). ICA has been successful in recovering Gabor filters from natural images [22]. The argument given below can be seen as an alternative explanation of our numerical simulation results given in Section III B.

The goal of ICA is to recover the sources, given that the input data are generated by the following linear model [21]:

$$\mathbf{x}_T = \mathbf{A}\mathbf{s}_T, \quad (21)$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is the mixing matrix, assumed to be invertible, and the random source vector, \mathbf{s}_T has statistically independent elements. Each source is assumed to have zero mean and sparse, e.g. Laplace distributed.

To establish a connection between sparse symmetric matrix factorization and ICA, we first show that the whitened input (21) is an orthogonal rotation of the original sources [19]. To see this, we rewrite the whitened input, $\tilde{\mathbf{x}}_T$ in terms of the assumed to be known whitening matrix, \mathbf{Q} , and by substituting (21) find:

$$\tilde{\mathbf{x}}_T = \mathbf{Q}\mathbf{x}_T = \mathbf{Q}\mathbf{A}\mathbf{s}_T. \quad (22)$$

By denoting $\mathbf{QA} \equiv \mathbf{G}$ we obtain from (22):

$$\tilde{\mathbf{x}}_T = \mathbf{G}\mathbf{s}_T.$$

The orthonormality of the square matrix \mathbf{G} follows from the orthonormality of whitened data [19]:

$$\mathbf{I}_n = \langle \tilde{\mathbf{x}}_T, \tilde{\mathbf{x}}_T' \rangle_T = \mathbf{G} \langle \mathbf{s}_T, \mathbf{s}_T' \rangle_T \mathbf{G}' = \mathbf{G}\mathbf{G}'. \quad (23)$$

To demonstrate that our algorithm can be used as online ICA, we rewrite the cost function (7) for the whitened input by using the orthonormality \mathbf{G} (23):

$$\begin{aligned} \mathbf{y}_T &= \arg \min_{\mathbf{y}_T} \left\| \tilde{\mathbf{X}}\tilde{\mathbf{X}} - \mathbf{Y}\mathbf{Y}' \right\|_F^2 + \lambda \sum_i \left\| \mathbf{y}'_{\cdot,i} \mathbf{y}_{\cdot,i} \right\|_1 \\ &= \arg \min_{\mathbf{y}_T} \left\| \mathbf{S}'\mathbf{S} - \mathbf{Y}\mathbf{Y}' \right\|_F^2 + \lambda \sum_i \left\| \mathbf{y}'_{\cdot,i} \mathbf{y}_{\cdot,i} \right\|_1. \end{aligned} \quad (24)$$

Because \mathbf{Y} has the same dimensionality as \mathbf{S} , $\mathbf{Y} = \mathbf{S}$ is a minimum of the unregularized cost in (24). However, this minimum is not unique: \mathbf{Y} can be left-rotated by an orthonormal matrix without affecting the cost. Then, the role

of the sparsity-inducing regularizer is to favor a sparse \mathbf{Y} , allowing the recovery of the original sparse \mathbf{S} .

The analysis of this section can be extended straightforwardly to the offline ICA problem. Symmetric matrix factorization or whitened input, with a suitably chosen sparsity inducing regularizer can be used as an ICA cost function.

IV. DISCUSSION

In this paper, by introducing a novel cost-function we derived an online algorithm that reproduces many features of the OF model but can be implemented by a single-layer neural network relying only on local learning rules. Therefore, we proposed a more biologically plausible implementation of the sparse coding hypothesis.

A. Biological relevance

1. Weighted summation of inputs and soft thresholding. Our online algorithm maps onto a neural network where each unit performs soft thresholding of the weighted sum of its inputs (both feedforward and lateral). Such computation corresponds to a commonly used basic model of biological neurons. Although, the two-sided thresholding our algorithm requires is not encountered in biological neurons, it may be implemented by a pair of neurons each responsible for positive or negative inputs. Such ON and OFF neurons exist in the peripheral visual system of both vertebrates and invertebrates [27].

2. Local Hebbian and anti-Hebbian synaptic learning rules. The learning rules we derived are consistent with those previously abstracted from biological observations of synaptic plasticity. Crucially, these learning rules do not require any synapse to keep track of the activity of neurons other than the pre- and postsynaptic pair it connects. Anti-Hebbian learning could be implemented indirectly via a Hebbian update of the synaptic weights of inhibitory interneurons.

3. Dependence of learning rate on cumulative activity. The learning rate in the synaptic weight update is inversely proportional to the cumulative activity of the postsynaptic neuron (11). Such variation of plasticity with time corresponds to the reports of LTP decaying with age in an activity dependent manner [23-25].

4. Sparsity of neuronal activity. The distribution of neuronal firing has a peak at zero and a heavy tail, Fig. 3C in agreement with physiological measurements [16,26].

B. Symmetric matrix factorization as a generic model of neural computation

We believe that the significance of online symmetric matrix factorization goes beyond deriving sparse dictionary learning with local Hebbian and anti-Hebbian learning rules. We speculate that it serves as a powerful and versatile

elementary building block of neural computation. Indeed, symmetric matrix factorization with (and without) various constraints can solve multiple computational objectives. We argued above that ICA can be formulated as a symmetric matrix factorization problem. Furthermore, unconstrained symmetric matrix factorization can compute the principal subspace of the streamed data [17]. Nonnegative symmetric matrix factorization can be viewed as a clustering algorithm capable of nonlinear feature discovery [18]. Jointly, these tools represent a formidable arsenal for modeling neural computation.

ACKNOWLEDGMENTS

The authors would like to thank Sanjeev Arora, Alex Genkin, Bruno Olshausen, Eftychios Pnevmatikakis, Christopher Rozell, and Zaid Towfic for helpful discussions.

REFERENCES

- [1] D. H. Hubel, *Eye, Brain and Vision*, W. H. Freeman; 2nd edition, 1995.
- [2] D. H. Hubel, & T. N. Wiesel. Receptive fields and functional architecture of monkey striate cortex. *Journal of Physiology*, London, 195, 215-243.
- [3] B. Olshausen and D. Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images," *Nature*, vol. 381, pp. 607-609, 1996.
- [4] B. Olshausen and D. Field, "Sparse coding with an overcomplete basis set: a strategy employed by V1?," *Vision Res.*, vol. 37, no. 23, pp. 3311-3325, 1997.
- [5] B. A. Olshausen and D. J. Field, "Sparse coding of sensory inputs," *Curr Opin Neurobiol*, vol. 14, no. 4, pp. 481-7, Aug. 2004.
- [6] C. J. Rozell, D. H. Johnson, R. G. Baraniuk, and B.A. Olshausen, "Sparse coding via thresholding and local competition in neural circuits," *Neural Computation*, vol. 20, pp. 2526-2563, 2008.
- [7] S. Boyd, & L. Vandenberghe. *Convex optimization*. Cambridge university press (2009).
- [8] T. Hastie, R. Tibshirani, J. Friedman., *The elements of statistical learning* (Vol. 2, No. 1). New York: Springer (2009).
- [9] P. Foldiak. Forming Sparse representations by local anti-Hebbian learning, *Biol. Cybern.* 64, 165-170 (1990)
- [10] J. Zylberberg, J. T. Murphy, & M. R. DeWeese, A sparse coding model with synaptically local plasticity and spiking neurons can account for the diverse shapes of V1 simple cell receptive fields. *PLoS computational biology*, 7(10), e1002250 (2011).
- [11] A. A. Koulakov and D. Rinberg. "Sparse incomplete representations: a potential role of olfactory granule cells." *Neuron* 72.1 (2011): 124-136.
- [12] S. Druckmann, T. Hu, & D. B. Chklovskii. A mechanistic model of early sensory processing based on subtracting sparse representations. In *Advances in Neural Information Processing Systems* (pp. 1979-1987). 2012.
- [13] J. Carroll and J. Chang. Idioscal (individual differences in orientation scaling) A generalization of indscal allowing idiosyncratic reference systems as well as an analytic approximation to indscal. In *Psychometric meeting, Princeton, NJ* (1972).
- [14] C. Ding, X. He, and H. Simon, "On the Equivalence of Nonnegative Matrix Factorization and Spectral Clustering," *SDM*, no. 4, 2005.
- [15] E. Oja, "A simplified neuron model as a principal component analyzer," *Journal of Mathematical Biology*, vol. 15, pp. 267-273, 1982.
- [16] T. Hu; Z.J. Towfic, C. Pehlevan, A. Genkin, D.B. Chklovskii. "A neuron as a signal processing device," *Asilomar Conference on Signals, Systems and Computers, 2013*, vol., no., pp.362,366, 3-6 Nov. 2013
- [17] C. Pehlevan, T. Hu, & D.B. Chklovskii. "A Hebbian/Anti-Hebbian Network for Linear Subspace Tracking: A Derivation from Multi-dimensional Scaling of Streaming Data," *Neural Computation*, submitted.
- [18] C. Pehlevan D.B. Chklovskii. "A Hebbian/Anti-Hebbian Network for Dervied from Online Nonnegative Matrix Factorization Can Cluster and Discover Sparse Features," *Asilomar Conference on Signals, Systems and Computers*, 2-5 Nov. 2014
- [19] A. Hyvärinen, J. Hurri and P.P. Hoyer. *Natural Image Statistics-A probabilistic approach to early computational vision* (Springer-Verlag, London), 2009.
- [20] D.L. Ringach. Spatial structure and symmetry of simple-cell receptive fields in macaque primary visual cortex. *J Neurophysiol.* 88(1):455-63 (2002).
- [21] A. Hyvärinen and E. Oja, "Independent component analysis: algorithms and applications," *Neural networks*, vol. 13, no. 4-5, pp. 411-430, Jun. 2000.
- [22] A. Bell and T. J. Sejnowski, "The 'independent components' of natural scenes are edge filters.," *Vision Res.*, vol. 37, no. 23, pp. 3327-38, Dec. 1997.
- [23] M. Crair and R. Malenka, "A critical period for long-term potentiation at thalamocortical synapses," *Nature*, vol. 375, pp. 325-327, 1995.
- [24] A. Kirkwood, H. K. Lee, and M. F. Bear, "Co-regulation of long-term potentiation and experience-dependent synaptic plasticity in visual cortex by age and experience," *Nature*, vol. 375, pp. 328-331, 1995.
- [25] C. Poo and J. S. Isaacson, "An early critical period for long-term plasticity and structural modification of sensory synapses in olfactory cortex," *J. Neurosci*, vol. 27, pp. 7553-7558, 2007.
- [26] T. Hromadka, MR Deweese, AM Zador. Sparse representation of sounds in the unanesthetized auditory cortex. *PLoS Biol* 6:e16 (2008).
- [27] R.H. Masland, "The fundamental plan of the retina," *Nature Neuroscience*, 4(9), 877-886, 2001
- [28] B Li, MR Peterson, RD Freeman, *Journal of Neurophysiology*, 90(1) 204-217, Jul 2003,.