

# Biologically plausible single-layer networks for nonnegative independent component analysis

David Lipshutz<sup>1</sup>, Cengiz Pehlevan<sup>2</sup>, and Dmitri B. Chklovskii<sup>1,3</sup>

<sup>1</sup>Center for Computational Neuroscience, Flatiron Institute

<sup>2</sup>John A. Paulson School of Engineering and Applied Sciences, Harvard University

<sup>3</sup>Neuroscience Institute, NYU Medical Center

March 8, 2022

## Abstract

An important problem in neuroscience is to understand how brains extract relevant signals from mixtures of unknown sources, i.e., perform blind source separation. To model how the brain performs this task, we seek a biologically plausible single-layer neural network implementation of a blind source separation algorithm. For biological plausibility, we require the network to satisfy the following three basic properties of neuronal circuits: (i) the network operates in the online setting; (ii) synaptic learning rules are local; (iii) neuronal outputs are nonnegative. Closest is the work by Pehlevan et al. [Neural Computation, 29, 2925–2954 (2017)], which considers Nonnegative Independent Component Analysis (NICA), a special case of blind source separation that assumes the mixture is a linear combination of uncorrelated, nonnegative sources. They derive an algorithm with a biologically plausible 2-layer network implementation. In this work, we improve upon their result by deriving 2 algorithms for NICA, each with a biologically plausible *single-layer* network implementation. The first algorithm maps onto a network with indirect lateral connections mediated by interneurons. The second algorithm maps onto a network with direct lateral connections and multi-compartmental output neurons.

Keywords: Blind source separation, nonnegative independent component analysis, neural network, local learning rules

## 1 Introduction

Brains effortlessly extract relevant signals from mixtures of unknown sources [3, 4, 5, 21, 12, 1, 20, 11, 2], an unsupervised signal processing problem known as blind source separation. A classic example in audition is the cocktail party problem, in which a listener tries to follow a single conversation in the presence of multiple background conversations. We seek a model of how brains perform blind source separation.

A special case of blind source separation is Nonnegative Independent Component Analysis (NICA), which assumes a generative model in which the mixture of stimuli is a linear combination of uncorrelated, nonnegative sources; i.e.,  $\mathbf{x} = \mathbf{A}\mathbf{s}$ , where  $\mathbf{s}$  denotes the nonnegative vector of source intensities,  $\mathbf{A}$  is a mixing matrix and  $\mathbf{x}$  denotes the vector of mixed stimuli. The goal of NICA is

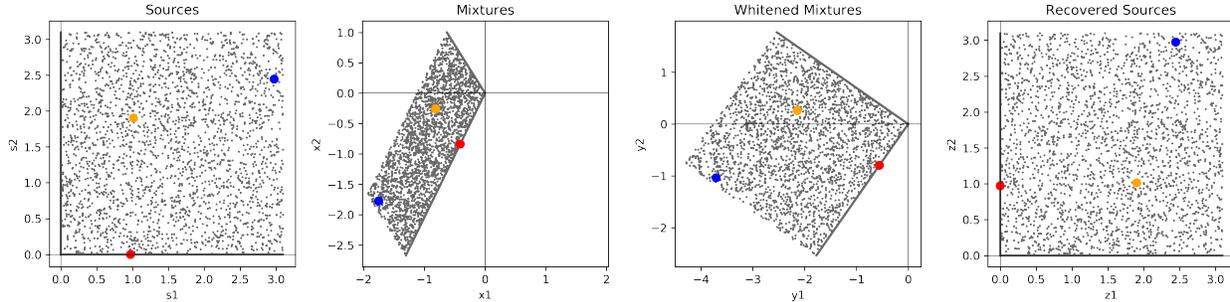


Figure 1: Illustration of Plumbley’s 2-step algorithm for NICA. The red, blue and oranges dots track three source vectors across the mixing, whitening and rotation steps. Our algorithms transform Mixtures into Recovered Sources in a single step implemented by single-layer neural networks.

to infer the source vectors  $\mathbf{s}$  from the mixture vectors  $\mathbf{x}$ . Both the linear additivity of stimuli and nonnegativity of the sources are reasonable assumptions in biological applications. For example, in olfaction, concentrations of odorants are both additive and nonnegative.

Plumbley [16] showed that when the sources are well-grounded (i.e., they have nonzero probability of taking infinitesimally small values), NICA can be solved in 2 steps; see Figure 1. In the first step, the mixture undergoes noncentered whitening; that is, the mixture is linearly transformed to have identity covariance matrix. The second step rotates the mixture until it lies in the nonnegative orthant. The result of these 2 steps must be a permutation of the original sources. This important observation led to a number of algorithms for implementing the rotation step [17, 18, 13, 22], many of which have neural network implementations.

Unfortunately, the above-mentioned networks do not offer a viable model of brain function because they do not satisfy one or more of the following three common requirements for biological plausibility [14]. First, the network operates in the online or streaming setting where it receives one input at a time and the output is computed before the next input arrives. Second, each synaptic update is local in the sense that it depends only on variables represented in the pre- and post-synaptic neurons. Third, the neuronal outputs are nonnegative.

Building on Plumbley’s method, Pehlevan et al. [15] proposed a 2-layer network for NICA, with each layer derived from a principled objective function. The first layer implements noncentered whitening and the second orthogonally rotates the whitened mixture. While their networks satisfies the requirements for biological plausibility, from a biological perspective, there are advantages to a single-layer network that economizes the number of neurons, which take up valuable resources such as space [19] and metabolic energy [8].

In this work, we derive 2 NICA algorithms (Algorithms 1 & 2) that can be implemented in biologically plausible single-layer networks. The first algorithm maps onto a network with point neurons and indirect lateral connections mediated by interneurons (Figure 2), and the second algorithm maps onto a network with 2-compartmental neurons and direct lateral connections (Figure 3). To derive our algorithms, we adopt a normative approach which relies on the fact that the original sources can be expressed (up to permutation) as optimal solutions of single objective functions that combine the 2 objectives from [15].

**Notation.** For integers  $p, q$ , let  $\mathbb{R}^p$  denote  $p$ -dimensional Euclidean space,  $\mathbb{R}_+^p$  denote the nonnegative orthant in  $\mathbb{R}^p$ ,  $\mathbb{R}^{p \times q}$  denote the set of  $p \times q$  real-valued matrices and  $\mathbb{R}_+^{p \times q}$  denote the subset of matrices with nonnegative entries. Let  $\mathcal{S}_{++}^p$  denote the set of  $p \times p$  positive definite matrices and let  $\mathbf{I}_p$  denote the  $p \times p$  identity matrix. Given  $T$  samples  $\mathbf{h}_1, \dots, \mathbf{h}_T$  of a time series, let

$$\langle \mathbf{h} \rangle := \frac{1}{T} \sum_{t=1}^T \mathbf{h}_t, \quad \mathbf{C}_{hh} := \frac{1}{T} \sum_{t=1}^T (\mathbf{h}_t - \langle \mathbf{h} \rangle)(\mathbf{h}_t - \langle \mathbf{h} \rangle)^\top$$

respectively denote the empirical mean and covariance of the time series. Let  $\bar{\mathbf{h}}_t := \frac{1}{t}(\mathbf{h}_1 + \dots + \mathbf{h}_t)$  denote the running sample mean. Given a data matrix  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_T]$ , let  $\delta\mathbf{H} := [\mathbf{h}_1 - \langle \mathbf{h} \rangle, \dots, \mathbf{h}_T - \langle \mathbf{h} \rangle]$  denote the centered data matrix.

## 2 Review of prior work

In this section, we review Plumbley's analysis [16] and the objective functions used by Pehlevan et al. [15] to derive a 2-layer network for NICA. Let  $d \geq 2$  and  $\mathbf{s}_1, \dots, \mathbf{s}_T \in \mathbb{R}_+^d$  be  $T$  samples of  $d$ -dimensional nonnegative source vectors whose components are uncorrelated. Since a constant factor multiplying a source can be absorbed into the associated column of the mixing matrix  $\mathbf{A}$ , we can assume, without loss of generality, that each component of the source vector has unit sample variance. In particular,  $\mathbf{C}_{ss} = \mathbf{I}_d$ . Let  $k \geq d$ ,  $\mathbf{A}$  be a full rank  $k \times d$  mixing matrix and define the  $k$ -dimensional mixture vectors by  $\mathbf{x}_t := \mathbf{A}\mathbf{s}_t$  for  $t = 1, \dots, T$ .

### 2.1 Plumbley's NICA method

Plumbley [16] proposed solving NICA in 2 steps: noncentered whitening followed by orthogonal transformation, which are depicted in Figure 1.

Noncentered whitening is a linear transformation  $\mathbf{y} := \mathbf{F}\mathbf{x}$  of the mixture, where  $\mathbf{y} \in \mathbb{R}^d$  and  $\mathbf{F}$  is a  $d \times k$  whitening matrix such that  $\mathbf{y}$  has identity covariance matrix, i.e.,  $\mathbf{C}_{yy} = \mathbf{I}_d$ . The combined effect of source mixing and prewhitening steps, which is encoded in the  $d \times d$  matrix  $\mathbf{FA}$  (since  $\mathbf{y} = \mathbf{F}\mathbf{x}$  and  $\mathbf{x} = \mathbf{A}\mathbf{s}$ ), is an orthogonal transformation. To see this, we use the facts that  $\mathbf{C}_{ss} = \mathbf{I}_d$ ,  $\mathbf{y} = \mathbf{F}\mathbf{A}\mathbf{s}$  and  $\mathbf{C}_{yy} = \mathbf{I}_d$  to write

$$\begin{aligned} (\mathbf{FA})(\mathbf{FA})^\top &= (\mathbf{FA})\mathbf{C}_{ss}(\mathbf{FA})^\top = \frac{1}{T} \sum_{t=1}^T \mathbf{FA}(\mathbf{s}_t - \langle \mathbf{s} \rangle)(\mathbf{s}_t - \langle \mathbf{s} \rangle)^\top (\mathbf{FA})^\top \\ &= \frac{1}{T} \sum_{t=1}^T (\mathbf{y}_t - \langle \mathbf{y} \rangle)(\mathbf{y}_t - \langle \mathbf{y} \rangle)^\top = \mathbf{C}_{yy} = \mathbf{I}_d. \end{aligned}$$

In the second step, one looks for an orthogonal matrix  $\mathbf{R}$  such that the transformation  $\mathbf{z} := \mathbf{R}\mathbf{y}$  is nonnegative. For the solution to be unique up to a permutation, each source  $s^i$  must be well grounded; that is,  $P(s^i < \delta) > 0$  for all  $\delta > 0$ . Then by [16, Theorem 1], the vector  $\mathbf{z}$  is equal to a permutation of the sources  $\mathbf{s}$ .

## 2.2 Similarity matching objectives for the 2-step algorithm

To obtain a biologically plausible network, Pehlevan et al. [15] proposed novel mathematical formulations of the noncentered whitening and rotation steps, which can be implemented by a biological plausible 2-layer network.

Here we recall the principled objective functions they use for each layer, which are closely related to the objective functions we use to derive our networks. To this end, define the  $k \times T$  concatenated data matrix  $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_T]$ . In the first step, Pehlevan et al. [15] optimize, with respect to the  $d \times T$  matrix  $\mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_T]$ , the following objective:

$$\arg \min_{\mathbf{Y} \in \mathbb{R}^{d \times T}} -\text{Tr}(\delta \mathbf{Y}^\top \delta \mathbf{Y} \delta \mathbf{X}^\top \delta \mathbf{X}) \quad \text{subject to} \quad \delta \mathbf{Y}^\top \delta \mathbf{Y} \preceq T \mathbf{I}_T \quad \text{and} \quad \mathbf{Y} = \mathbf{F} \mathbf{X}, \quad (1)$$

where  $\mathbf{F}$  is some  $d \times k$  matrix and the constraint enforces that the difference  $T \mathbf{I}_T - \delta \mathbf{Y}^\top \delta \mathbf{Y}$  is positive semidefinite. As shown in [15], objective (1) is optimized when  $\mathbf{Y}$  is a noncentered whitened transformation of  $\mathbf{X}$ .

For the second step, Pehlevan et al. [15] introduce the following Nonnegative Similarity Matching (NSM) objective:

$$\arg \min_{\mathbf{Z} \in \mathbb{R}_+^{d \times T}} \|\mathbf{Z}^\top \mathbf{Z} - \mathbf{Y}^\top \mathbf{Y}\|_{\text{Frob}}^2. \quad (2)$$

The objective minimizes the mismatch between similarities of the nonnegative outputs  $\mathbf{Z}$  and the noncentered whitened mixtures  $\mathbf{Y}$  (as measured by inner products). As shown in [15], any orthogonal transformation of  $\mathbf{Y}$  to the nonnegative orthant, which corresponds to a permutation of the original sources, is a solution of the NSM objective (2).

From objectives (1) and (2), Pehlevan et al. [15] derive a 2-step algorithm for NICA that can be implemented in a 2-layer neural network that operates in the online setting, uses local learning rules, and whose rotation layer has nonnegative neuronal outputs.

## 3 Combined objectives for NICA

We now modify objectives (1) and (2) to obtain 2 objectives for NICA, which will be the starting points for the derivations of our 2 online NICA algorithms with single-layer neural network implementations.

### 3.1 Adding a nonnegativity constraint to the noncentered whitening objective

We first modify the noncentered whitening objective (1). Note that the solution of objective (1) is not unique — left multiplying any solution  $\mathbf{Y}$  by an orthogonal matrix  $\mathbf{R}$  yields another noncentered whitened transformation of  $\mathbf{X}$ . In fact, the second step of Plumbley’s method [16] is to identify an orthogonal transformation  $\mathbf{R}$  that results in a *nonnegative* whitened transformation  $\mathbf{Z} = \mathbf{R} \mathbf{Y}$ . Here, we combine the 2 objectives by adding a nonnegativity constraint to the noncentered whitening objective (1), as follows:

$$\arg \min_{\mathbf{Y} \in \mathbb{R}_+^{d \times T}} -\text{Tr}(\delta \mathbf{Y}^\top \delta \mathbf{Y} \delta \mathbf{X}^\top \delta \mathbf{X}) \quad \text{subject to} \quad \delta \mathbf{Y}^\top \delta \mathbf{Y} \preceq T \mathbf{I}_T \quad \text{and} \quad \mathbf{Y} = \mathbf{F} \mathbf{X}, \quad (3)$$

where  $\mathbf{F}$  is some  $d \times k$  matrix.

### 3.2 Adding a whitening matrix to the NSM objective

Next, we alter the NSM objective (2) by replacing the Gram matrix  $\mathbf{Y}^\top \mathbf{Y}$  with terms that depend only on  $\mathbf{X}$ , which will avoid the need for the noncentered whitening step. Consider the eigendecomposition of the covariance matrix  $\mathbf{C}_{xx} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ , where  $\mathbf{U}$  is a  $k \times d$  matrix with orthonormal column vectors and  $\mathbf{\Lambda}$  is a  $d \times d$  diagonal matrix whose diagonal entries are the nonzero eigenvalues of  $\mathbf{C}_{xx}$ . Then the whitening matrix  $\mathbf{F}$  must be of the form  $\mathbf{Q}\mathbf{\Lambda}^{-\frac{1}{2}}\mathbf{U}^\top$ , where  $\mathbf{Q}$  is an arbitrary  $d \times d$  orthogonal matrix. Therefore,

$$\mathbf{Y}^\top \mathbf{Y} = \mathbf{X}^\top \mathbf{F}^\top \mathbf{F} \mathbf{X} = \mathbf{X}^\top \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^\top \mathbf{X} = \mathbf{X}^\top \mathbf{C}_{xx}^+ \mathbf{X},$$

where  $\mathbf{C}_{xx}^+ := \mathbf{U}\mathbf{\Lambda}^{-1}\mathbf{U}^\top$  is the Moore-Penrose inverse of  $\mathbf{C}_{xx}$ . Substituting in for  $\mathbf{Y}^\top \mathbf{Y}$  in the NSM objective (2) results in our second objective:

$$\arg \min_{\mathbf{Z} \in \mathbb{R}_+^{d \times T}} \|\mathbf{Z}^\top \mathbf{Z} - \mathbf{X}^\top \mathbf{C}_{xx}^+ \mathbf{X}\|_{\text{Frob}}^2. \quad (4)$$

## 4 Single-layer neural networks for NICA

Starting from objectives (3) and (4), we derive our 2 online NICA algorithms. The first algorithm maps onto a single-layer network with point neurons and *indirect* lateral connections. The second algorithm maps onto a single-layer network with 2-compartmental neurons and *direct* lateral connections.

### 4.1 Single-layer network with point neurons and indirect lateral connections

The derivation of our online algorithm starting from objective (3) closely follows the derivation of the whitening layer in the network derived in [15]. The main difference is that the neuronal outputs are constrained to be nonnegative. To begin, we introduce  $m$ -dimensional activity vectors  $\mathbf{n}_1, \dots, \mathbf{n}_T$ , with  $m \geq d$ , which we concatenate into the data matrix  $\mathbf{N} := [\mathbf{n}_1, \dots, \mathbf{n}_T]$ , and use the Gramian  $\delta \mathbf{N}^\top \delta \mathbf{N}$  as a Lagrange multiplier to enforce the constraint  $\delta \mathbf{Y}^\top \delta \mathbf{Y} \preceq T \mathbf{I}_T$ :

$$\min_{\mathbf{Y} \in \mathbb{R}_+^{d \times T}} \max_{\mathbf{N} \in \mathbb{R}^{m \times T}} \text{Tr} \left[ -\delta \mathbf{Y}^\top \delta \mathbf{Y} \delta \mathbf{X}^\top \delta \mathbf{X} + \delta \mathbf{N}^\top \delta \mathbf{N} (\delta \mathbf{Y}^\top \delta \mathbf{Y} - T \mathbf{I}_T) \right] \quad \text{subject to} \quad \mathbf{Y} = \mathbf{F} \mathbf{X}.$$

Next, we normalize by  $T^2$  and substitute synaptic weight matrices  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$  in place of  $\frac{1}{T} \delta \mathbf{Y} \delta \mathbf{X}^\top$  and  $\frac{1}{T} \delta \mathbf{N} \delta \mathbf{Y}^\top$ , respectively:

$$\min_{\mathbf{Y} \in \mathbb{R}_+^{d \times T}} \max_{\mathbf{N} \in \mathbb{R}^{d \times T}} \min_{\mathbf{W}_{xy} \in \mathbb{R}^{d \times k}} \max_{\mathbf{W}_{yn} \in \mathbb{R}^{m \times d}} L_1(\mathbf{Y}, \mathbf{N}, \mathbf{W}_{xy}, \mathbf{W}_{yn}) \quad \text{subject to} \quad \mathbf{Y} = \mathbf{F} \mathbf{X},$$

where

$$\begin{aligned} L_1(\mathbf{Y}, \mathbf{N}, \mathbf{W}_{xy}, \mathbf{W}_{yn}) := & \frac{1}{T} \text{Tr} \left( 2\delta \mathbf{N}^\top \mathbf{W}_{yn} \delta \mathbf{Y} - 2\delta \mathbf{Y}^\top \mathbf{W}_{xy} \delta \mathbf{X} - \delta \mathbf{N}^\top \delta \mathbf{N} \right) \\ & - \text{Tr} \left( \mathbf{W}_{yn} \mathbf{W}_{yn}^\top \right) + \text{Tr} \left( \mathbf{W}_{xy} \mathbf{W}_{xy}^\top \right). \end{aligned}$$

The substitution can be readily justified by differentiating  $L_1$  with respect to  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$  and noting the minimum (resp. maximum) is achieved when  $\mathbf{W}_{xy} = \frac{1}{T} \delta \mathbf{Y} \delta \mathbf{X}^\top$  (resp.  $\mathbf{W}_{yn} = \frac{1}{T} \delta \mathbf{N} \delta \mathbf{Y}^\top$ ).

Since  $L_1$  satisfies the saddle point property with respect to  $\mathbf{N}$  and  $\mathbf{W}_{xy}$ , and with respect to  $\mathbf{Y}$  and  $\mathbf{W}_{yn}$ , we can interchange the order of optimization, as follows:

$$\min_{\mathbf{W}_{xy} \in \mathbb{R}^{d \times k}} \max_{\mathbf{W}_{yn} \in \mathbb{R}^{m \times d}} \min_{\mathbf{Y} \in \mathbb{R}_+^{d \times T}} \max_{\mathbf{N} \in \mathbb{R}^{d \times T}} L_1(\mathbf{Y}, \mathbf{N}, \mathbf{W}_{xy}, \mathbf{W}_{yn}) \quad \text{subject to} \quad \mathbf{Y} = \mathbf{F}\mathbf{X}. \quad (5)$$

We first solve objective (5) in the offline setting. In general, optimizing over  $(\mathbf{Y}, \mathbf{N})$  is challenging due to the constraint that  $\mathbf{Y}$  be a nonnegative linear transformation of  $\mathbf{X}$ . In appendix A, we show that when the synaptic weights  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$  are at their optimal values, we can optimize over  $(\mathbf{Y}, \mathbf{N})$  by repeating the following projected gradient descent steps until convergence:

$$\mathbf{Y} \leftarrow \left[ \mathbf{Y} + \gamma \left( \mathbf{W}_{xy} \mathbf{X} - \mathbf{W}_{yn}^\top \mathbf{N} \right) \right]_+, \quad \mathbf{N} \leftarrow \mathbf{N} + \gamma \left( \mathbf{W}_{yn} \mathbf{Y} - \mathbf{N} \right), \quad (6)$$

where  $\gamma > 0$  is a small step size and  $[\cdot]_+$  denotes taking the positive part elementwise, which ensures the nonnegativity of  $\mathbf{Y}$ . In the case the synaptic weights  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$  are not at their optimal values, we repeat the above projected gradient descent steps until convergence to obtain an approximation of the optimal  $(\mathbf{Y}, \mathbf{N})$ . We then perform a gradient descent-ascent step of the objective  $L_1$  with respect to  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$ :

$$\mathbf{W}_{xy} \leftarrow \mathbf{W}_{xy} + \eta \left( \frac{1}{T} \delta \mathbf{Y} \delta \mathbf{X}^\top - \mathbf{W}_{xy} \right) \quad (7)$$

$$\mathbf{W}_{yn} \leftarrow \mathbf{W}_{yn} + \eta \left( \frac{1}{T} \delta \mathbf{N} \delta \mathbf{Y}^\top - \mathbf{W}_{yn} \right). \quad (8)$$

Here  $\eta > 0$  is the learning rate for  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$ .

Next, we solve the objective (5) in the online setting. At each time step  $t$ , we approximate the optimization over  $(\mathbf{y}_t, \mathbf{n}_t)$  by taking the following projected gradient descent steps until convergence:

$$\mathbf{y}_t \leftarrow [\mathbf{y}_t + \gamma(\mathbf{W}_{xy} \mathbf{x}_t - \mathbf{W}_{ny} \mathbf{n}_t)]_+, \quad \mathbf{n}_t \leftarrow \mathbf{n}_t + \gamma(\mathbf{W}_{yn} \mathbf{y}_t - \mathbf{n}_t), \quad (9)$$

where we have defined  $\mathbf{W}_{ny} := \mathbf{W}_{yn}^\top$ . We then take stochastic gradient descent-ascent steps in  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$  by replacing the averages in equations (7) and (8) with their online approximations:

$$\begin{aligned} \mathbf{W}_{xy} &\leftarrow \mathbf{W}_{xy} + \eta \left( (\mathbf{y}_t - \bar{\mathbf{y}}_t)(\mathbf{x}_t - \bar{\mathbf{x}}_t)^\top - \mathbf{W}_{xy} \right) \\ \mathbf{W}_{yn} &\leftarrow \mathbf{W}_{yn} + \eta \left( (\mathbf{n}_t - \bar{\mathbf{n}}_t)(\mathbf{y}_t - \bar{\mathbf{y}}_t)^\top - \mathbf{W}_{yn} \right) \\ \mathbf{W}_{ny} &\leftarrow \mathbf{W}_{ny} + \eta \left( (\mathbf{y}_t - \bar{\mathbf{y}}_t)(\mathbf{n}_t - \bar{\mathbf{n}}_t)^\top - \mathbf{W}_{ny} \right). \end{aligned}$$

The symmetry of the updates for  $\mathbf{W}_{ny}$  and  $\mathbf{W}_{yn}$  ensures that  $\mathbf{W}_{ny} = \mathbf{W}_{yn}^\top$  after each iteration provided the constraint holds at initialization. As we show in appendix B, we can relax this initialization constraint, which yields our first online NICA algorithm, Algorithm 1.

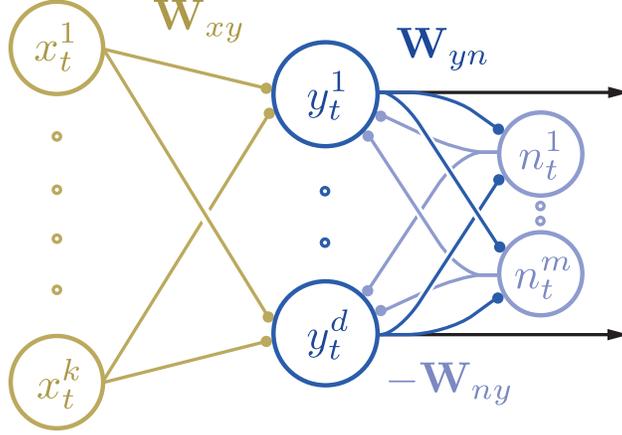


Figure 2: Single-layer network with interneurons for implementing Algorithm 1.

---

**Algorithm 1: Bio-NICA with interneurons**

---

**input** mixtures  $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ ; parameters  $\gamma, \eta$   
**initialize**  $\mathbf{W}_{xy}, \mathbf{W}_{yn}, \mathbf{W}_{ny}, \bar{\mathbf{x}}_0 = \mathbf{0}, \bar{\mathbf{y}}_0 = \mathbf{0}, \bar{\mathbf{n}}_0 = \mathbf{0}$   
**for**  $t = 1, 2, \dots, T$  **do**  
  **repeat**  
     $\mathbf{y}_t \leftarrow [\mathbf{y}_t + \gamma(\mathbf{W}_{xy}\mathbf{x}_t - \mathbf{W}_{ny}\mathbf{n}_t)]_+$   
     $\mathbf{n}_t \leftarrow \mathbf{n}_t + \gamma(\mathbf{W}_{yn}\mathbf{y}_t - \mathbf{n}_t)$   
  **until** convergence  
   $\bar{\mathbf{x}}_t \leftarrow \bar{\mathbf{x}}_{t-1} + \frac{1}{t}(\mathbf{x}_t - \bar{\mathbf{x}}_{t-1})$   
   $\bar{\mathbf{y}}_t \leftarrow \bar{\mathbf{y}}_{t-1} + \frac{1}{t}(\mathbf{y}_t - \bar{\mathbf{y}}_{t-1})$   
   $\bar{\mathbf{n}}_t \leftarrow \bar{\mathbf{n}}_{t-1} + \frac{1}{t}(\mathbf{n}_t - \bar{\mathbf{n}}_{t-1})$   
   $\mathbf{W}_{xy} \leftarrow \mathbf{W}_{xy} + \eta((\mathbf{y}_t - \bar{\mathbf{y}}_t)(\mathbf{x}_t - \bar{\mathbf{x}}_t)^\top - \mathbf{W}_{xy})$   
   $\mathbf{W}_{ny} \leftarrow \mathbf{W}_{ny} + \eta((\mathbf{y}_t - \bar{\mathbf{y}}_t)(\mathbf{n}_t - \bar{\mathbf{n}}_t)^\top - \mathbf{W}_{ny})$   
   $\mathbf{W}_{yn} \leftarrow \mathbf{W}_{yn} + \eta((\mathbf{n}_t - \bar{\mathbf{n}}_t)(\mathbf{y}_t - \bar{\mathbf{y}}_t)^\top - \mathbf{W}_{yn})$   
**end for**

---

Algorithm 1 can be implemented in a single-layer network with point neurons and indirect lateral connections mediated by interneurons, Figure 2, so we refer to the algorithm as ‘Bio-NICA with interneurons’. The network consists of  $k$  input neurons,  $d$  principal (output) output neurons and  $m$  interneurons. Feedforward synapses between the input and principal neurons encode the weight matrix  $\mathbf{W}_{xy}$  and lateral synapses between the principal neurons (resp. interneurons) and the interneurons (resp. principal neurons) encode the weight matrix  $\mathbf{W}_{yn}$  (resp.  $\mathbf{W}_{ny}$ ). At each time step  $t$ , the  $k$ -dimensional mixture  $\mathbf{x}_t$ , which is represented by the  $k$  input neurons, is multiplied by the weight matrix  $\mathbf{W}_{xy}$ , which yields the  $d$ -dimensional projection  $\mathbf{W}_{xy}\mathbf{x}_t$ . This is followed by the fast recurrent dynamics in equation (9). The equilibrium values of  $\mathbf{y}_t$  and  $\mathbf{n}_t$  respectively correspond to the nonnegative output of the principal neurons and the output of the interneurons.

We can write the elementwise synaptic updates as follows,

$$\begin{aligned} W_{xy}^{ij} &\leftarrow W_{xy}^{ij} + \eta \left( (y_t^i - \bar{y}_t^i)(x_t^j - \bar{x}_t^j) - W_{xy}^{ij} \right), & 1 \leq i \leq d, 1 \leq j \leq k, \\ W_{ny}^{ij} &\leftarrow W_{ny}^{ij} + \eta \left( (y_t^i - \bar{y}_t^i)(n_t^j - \bar{n}_t^j) - W_{ny}^{ij} \right), & 1 \leq i \leq d, 1 \leq j \leq m, \\ W_{yn}^{ij} &\leftarrow W_{yn}^{ij} + \eta \left( (n_t^i - \bar{n}_t^i)(y_t^j - \bar{y}_t^j) - W_{yn}^{ij} \right), & 1 \leq i \leq m, 1 \leq j \leq d, \end{aligned}$$

where we recall that  $\bar{\mathbf{x}}_t$ ,  $\bar{\mathbf{y}}_t$  and  $\bar{\mathbf{n}}_t$  are the running means of  $\mathbf{x}_t$ ,  $\mathbf{y}_t$  and  $\mathbf{n}_t$ , respectively. We assume that each neuron stores the running mean of its activity. Biologically, these means could be represented at the pre- and post-synaptic terminals by slowly changing calcium concentrations. From the elementwise updates, we see that the update for each synapse is local in the sense that it only depends on variables that are represented in the pre- and post-synaptic neurons.

## 4.2 Single-layer network with 2-compartmental neurons and direct lateral connections

The derivation of the our online algorithm starting form objective (4) is closely related to the derivation of the single-layer networks with multi-compartmental neurons for solving generalized eigenvalue problems [10, 9]. To begin, we expand the square, drop terms that do not depend on  $\mathbf{Z}$ , and normalize by  $T^2$ :

$$\min_{\mathbf{Z} \in \mathbb{R}_+^{d \times T}} \frac{1}{T^2} \text{Tr} \left( -2\mathbf{Z}^\top \mathbf{Z} \mathbf{X}^\top \mathbf{C}_{xx}^+ \mathbf{X} + \mathbf{Z}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{Z} \right). \quad (10)$$

Next, we introduce synaptic weight matrices  $\mathbf{W}_{xz}$  and  $\mathbf{W}_{zz}$  in place of  $\frac{1}{T} \mathbf{Z} \mathbf{X}^\top \mathbf{C}_{xx}^+$  and  $\frac{1}{T} \mathbf{Z} \mathbf{Z}^\top$ , respectively, which results in the minimax objective:

$$\min_{\mathbf{Z} \in \mathbb{R}_+^{d \times T}} \min_{\mathbf{W}_{xz} \in \mathbb{R}^{d \times k}} \max_{\mathbf{W}_{zz} \in \mathcal{S}_{++}^d} L_2(\mathbf{Z}, \mathbf{W}_{xz}, \mathbf{W}_{zz}), \quad (11)$$

where

$$L_2(\mathbf{Z}, \mathbf{W}_{xz}, \mathbf{W}_{zz}) := \frac{2}{T} \text{Tr} \left( \mathbf{Z}^\top \mathbf{W}_{zz} \mathbf{Z} - 2\mathbf{Z}^\top \mathbf{W}_{xz} \mathbf{X} \right) - \text{Tr} \left( \mathbf{W}_{zz}^2 - 2\mathbf{W}_{xz} \mathbf{C}_{xx} \mathbf{W}_{xz}^\top \right).$$

The equivalence between the minimization problem (10) and the minimax problem (11) can be seen by taking partial derivatives of  $L_2$  with respect to  $\mathbf{W}_{xz}$  (resp.  $\mathbf{W}_{zz}$ ) and noting the minimum (resp. maximum) is achieved when  $\mathbf{W}_{xz} = \frac{1}{T} \mathbf{Z} \mathbf{X}^\top \mathbf{C}_{xx}^+$  (resp.  $\mathbf{W}_{zz} = \frac{1}{T} \mathbf{Z} \mathbf{Z}^\top$ ). Since the objective  $L_2$  satisfies the strict saddle point property with respect to  $\mathbf{Z}$  and  $\mathbf{W}_{zz}$ , we can interchange the order of optimization, as follows:

$$\min_{\mathbf{W}_{xz} \in \mathbb{R}^{d \times k}} \max_{\mathbf{W}_{zz} \in \mathcal{S}_{++}^d} \min_{\mathbf{Z} \in \mathbb{R}_+^{d \times T}} L_2(\mathbf{Z}, \mathbf{W}_{xz}, \mathbf{W}_{zz}). \quad (12)$$

We first solve the minimax objective (12) in the offline setting by minimizing  $L_2$  over  $\mathbf{Z}$  and then taking gradient descent-ascent steps in  $\mathbf{W}_{xz}$  and  $\mathbf{W}_{zz}$ . The minimization over  $\mathbf{Z}$  can be approximated by repeating the following projected gradient descent steps until convergence:

$$\mathbf{Z} \leftarrow [\mathbf{Z} + \gamma(\mathbf{W}_{xz} \mathbf{X} - \mathbf{W}_{zz} \mathbf{Z})]_+,$$

where  $\gamma > 0$  is a small step size. Next, having minimized over  $\mathbf{Z}$ , we perform a gradient descent-ascend step of the objective function  $L_2$  with respect to  $\mathbf{W}_{xz}$  and  $\mathbf{W}_{zz}$ :

$$\mathbf{W}_{xz} \leftarrow \mathbf{W}_{xz} + 2\eta \left( \frac{1}{T} \mathbf{Z} \mathbf{X}^\top - \mathbf{W}_{xz} \mathbf{C}_{xx} \right), \quad (13)$$

$$\mathbf{W}_{zz} \leftarrow \mathbf{W}_{zz} + \frac{\eta}{\tau} \left( \frac{1}{T} \mathbf{Z} \mathbf{Z}^\top - \mathbf{W}_{zz} \right). \quad (14)$$

Here  $\tau > 0$  is the ratio between the learning rates for  $\mathbf{W}_{xz}$  and  $\mathbf{W}_{zz}$ , and  $\eta \in (0, \tau)$  is the learning rate for  $\mathbf{W}_{xz}$ . The upper bound  $\eta < \tau$  ensures that  $\mathbf{W}_{zz}$  remains positive definite given a positive definite initialization.

To solve the minimax objective (12) in the online setting, we take stochastic gradient ascent-descent steps. At each time step  $t$ , analogous to the offline setting, we first minimize over the output  $\mathbf{z}_t$  by repeating the following projected gradient descent steps until convergence:

$$\mathbf{z}_t \leftarrow [\mathbf{z}_t + \gamma(\mathbf{c}_t - \mathbf{W}_{zz} \mathbf{z}_t)]_+, \quad (15)$$

where we have defined the projection  $\mathbf{c}_t := \mathbf{W}_{xz} \mathbf{x}_t$ . We then take stochastic gradient descent-ascend steps in  $\mathbf{W}_{xz}$  and  $\mathbf{W}_{zz}$ . To this end, we replace the averages  $\frac{1}{T} \mathbf{Z} \mathbf{X}^\top$  and  $\frac{1}{T} \mathbf{Z} \mathbf{Z}^\top$  in equations (13) and (14) with their respective online approximations  $(\mathbf{z}_t - \bar{\mathbf{z}}_t)(\mathbf{x}_t - \bar{\mathbf{x}}_t)^\top$  and  $(\mathbf{z}_t - \bar{\mathbf{z}}_t)(\mathbf{z}_t - \bar{\mathbf{z}}_t)^\top$ . While we could approximate the matrix  $\mathbf{W}_{xz} \mathbf{C}_{xx}$  in the online setting with  $\mathbf{W}_{xz}(\mathbf{x}_t - \bar{\mathbf{x}}_t)(\mathbf{x}_t - \bar{\mathbf{x}}_t)^\top$ , this does not lead to local learning rules. Instead, we observe that

$$\mathbf{W}_{xz} \mathbf{C}_{xx} = \frac{1}{T} \sum_{t=1}^T \mathbf{W}_{xz}(\mathbf{x}_t - \langle \mathbf{x} \rangle)(\mathbf{x}_t - \langle \mathbf{x} \rangle)^\top = \frac{1}{T} \sum_{t=1}^T (\mathbf{c}_t - \langle \mathbf{c} \rangle)(\mathbf{x}_t - \langle \mathbf{x} \rangle)^\top,$$

and replace  $\mathbf{W}_{xz} \mathbf{C}_{xx}$  with the online approximation  $(\mathbf{c}_t - \bar{\mathbf{c}}_t)(\mathbf{x}_t - \bar{\mathbf{x}}_t)^\top$ . This yields our second online algorithm for NICA, Algorithm 2.

---

**Algorithm 2:** Bio-NICA with 2-compartmental neurons

---

**input** mixtures  $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ ; parameters  $\gamma, \eta, \tau$   
**initialize**  $\mathbf{W}_{xz}, \mathbf{W}_{zz}, \bar{\mathbf{x}}_0 = \mathbf{0}, \bar{\mathbf{c}}_0 = \mathbf{0}$   
**for**  $t = 1, 2, \dots, T$  **do**  
     $\mathbf{c}_t \leftarrow \mathbf{W}_{xz} \mathbf{x}_t$   
    **repeat**  
         $\mathbf{z}_t \leftarrow [\mathbf{z}_t + \gamma(\mathbf{c}_t - \mathbf{W}_{zz} \mathbf{z}_t)]_+$   
    **until** convergence  
     $\bar{\mathbf{x}}_t \leftarrow \bar{\mathbf{x}}_{t-1} + \frac{1}{t}(\mathbf{x}_t - \bar{\mathbf{x}}_{t-1})$   
     $\bar{\mathbf{c}}_t \leftarrow \bar{\mathbf{c}}_{t-1} + \frac{1}{t}(\mathbf{c}_t - \bar{\mathbf{c}}_{t-1})$   
     $\mathbf{W}_{xz} \leftarrow \mathbf{W}_{xz} + 2\eta(\mathbf{z}_t \mathbf{x}_t^\top - (\mathbf{c}_t - \bar{\mathbf{c}}_t)(\mathbf{x}_t - \bar{\mathbf{x}}_t)^\top)$   
     $\mathbf{W}_{zz} \leftarrow \mathbf{W}_{zz} + \frac{\eta}{\tau}(\mathbf{z}_t \mathbf{z}_t^\top - \mathbf{W}_{zz})$   
**end for**

---

Algorithm 2 can be implemented in a single-layer network with 2-compartmental neurons and direct lateral connections, Figure 3, so we refer to the algorithm as ‘Bio-NICA with 2-compartmental neurons’. The network consists of  $k$  input neurons and  $d$  output neurons. Each output neuron has

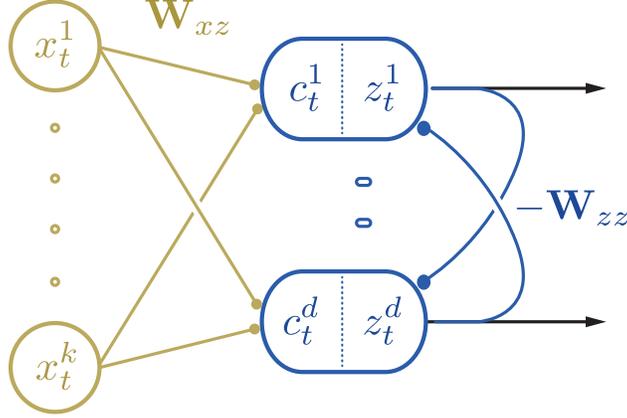


Figure 3: Single-layer network with 2-compartmental neurons for implementing Algorithm 2.

a dendritic compartment and a somatic compartment. Feedforward synapses between the input and output neurons encode the weight matrix  $\mathbf{W}_{xz}$  and recursive lateral synapses between the output neurons encode the weight matrix  $-\mathbf{W}_{zz}$ . At each time step  $t$ , the  $k$ -dimensional mixture  $\mathbf{x}_t$ , which is represented by the input neurons, is multiplied by the weight matrix  $\mathbf{W}_{xz}$ , which is encoded by the feedforward synapses connecting the input neurons to the output neurons. This yields the  $d$ -dimensional projection  $\mathbf{c}_t = \mathbf{W}_{xz}\mathbf{x}_t$ , which is computed in the dendritic compartments of the output neurons and then propagated to their somatic compartments. This is followed by the fast recurrent neural dynamics in equation (15). The equilibrium value of  $\mathbf{z}_t$  corresponds to the nonnegative somatic activity of the output neurons.

The elementwise synaptic updates are as follows,

$$\begin{aligned}
 W_{xz}^{ij} &\leftarrow W_{xz}^{ij} + 2\eta \left( z_t^i x_t^j - (\bar{c}_t^i - \bar{c}_t^i)(x_t^j - \bar{x}_t^j) \right), & 1 \leq i \leq d, 1 \leq j \leq k, \\
 W_{zz}^{ij} &\leftarrow W_{zz}^{ij} + \frac{\eta}{\tau} \left( z_t^i z_t^j - W_{zz}^{ij} \right), & 1 \leq i, j \leq d,
 \end{aligned}$$

where we recall that  $\bar{\mathbf{x}}_t$  and  $\bar{\mathbf{c}}_t$  are the running means of  $\mathbf{x}_t$  and  $\mathbf{c}_t$ , respectively. We assume that the input neurons and output neurons respectively store the running means  $\bar{\mathbf{x}}_t$  and  $\bar{\mathbf{c}}_t$ . Thus, we see that the update for each synapse is local; that is, the update depends only on variables that are represented in the pre- and post-synaptic neurons.

## 5 Numerical experiments

We evaluated Algorithms 1 and 2 on synthetic and real datasets and compare their performance to 2 state-of-the-art online NICA algorithms: Nonnegative PCA [18] and 2-layer NSM [15]. Nonnegative PCA requires (noncentered) pre-whitened inputs, which we implemented offline. To quantify the performance of the algorithms, we use the mean-squared error,

$$\text{error}(t) = \frac{1}{td} \sum_{t'=1}^t |\mathbf{s}_{t'} - \mathbf{P}\mathbf{y}_{t'}|^2,$$

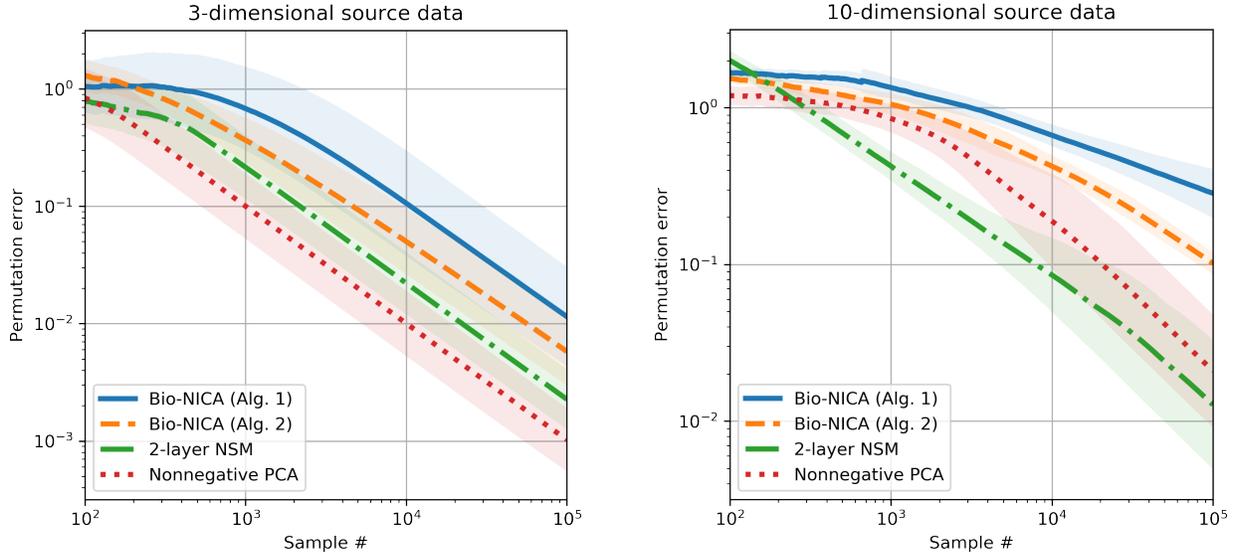


Figure 4: Performance of algorithms when presented with mixtures of sparse random uniform sources, in terms of permutation error. The lines and shaded regions denote the means and 90% confidence intervals over 10 runs.

where  $\mathbf{P}$  is the permutation matrix that minimizes the error at the final time point. For detailed descriptions of our implementations, see appendix C. The evaluation code is available at <https://github.com/flatironinstitute/bio-nica>.

### 5.1 Mixture of sparse random uniform sources

We first compare the algorithms on a synthetic dataset generated by independent and identically distributed samples. Following [15], each source sample was set to zero with probability  $1/2$  or sampled uniformly from the interval  $(0, \sqrt{48/5})$  with probability  $1/2$ . We used random square mixing matrices whose elements were independent standard normal random variables. In Figure 4, we plot the performance of each algorithm on mixtures of 3- and 10-dimensional sources.

### 5.2 Mixture of natural images

We apply the NICA algorithms to the problem of recovering images from their mixtures, see Figure 5 (left). Three image patches of size  $252 \times 252$  pixels were chosen from a set of images of natural scenes [6] previously used in [7, 18, 15]. Each image is treated as one source, with the pixel intensities (shifted and scaled to be well-ground and have unit variance) representing the  $252^2 = 63504$  samples. The source vectors were multiplied by a random  $3 \times 3$  mixing matrix to generate 3-dimensional mixtures, which were presented to the algorithms 5 times with a randomly permuted order in each presentation. In Figure 5 (right), we show the performance of each algorithm.

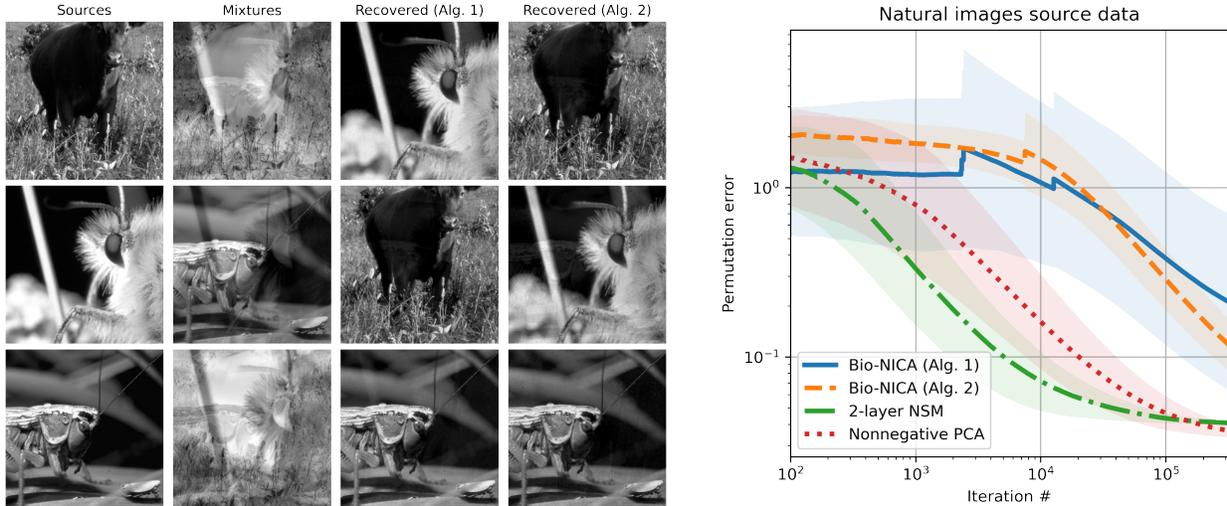


Figure 5: Performance of algorithms when presented with mixtures of natural images. The left image shows the sources, mixtures, and recovered sources (from Algorithms 1 and 2). The right plot shows the performance of the algorithms in terms of permutation error. The lines and shaded regions denote the means and 90% confidence intervals over 10 runs.

## 6 Summary

In this work, we derived 2 algorithms for NICA, each of which can be implemented by biologically plausible single-layer networks. Our networks respectively use two-thirds and one-third as many neurons as the 2-layer biologically plausible network derived in [15].

Our numerical experiments suggest that Algorithms 1 and 2 are outperformed by Nonnegative PCA and the 2-layer NSM. However, a direct comparison is not entirely fair because Nonnegative PCA requires prewhitened inputs and its neural network implementation does not use local learning rules, and the 2-layer NSM network requires 2 layers of neurons. Our algorithms perform both the whitening and the rotation steps in a single layer, which leads to a trade-off in performance. We view this as consistent with the fact that biological systems must make trade-offs between performance and resource efficiency.

Finally, we do not prove convergence guarantees for Algorithms 1 and 2. In general, establishing theoretical guarantees for gradient descent-ascent problems is challenging and is further complicated by the non-smoothness of the projected gradient descent steps in Algorithms 1 and 2.

**Acknowledgements.** We thank Siavash Golkar, Johannes Friedrich, Tiberiu Tesileanu, Alex Genkin, Jason Moore and Yanis Bahroun for helpful comments and feedback on an earlier draft of this work. We especially thank Siavash Golkar for pointing out that, in Sec. 4.2,  $\mathbf{W}(\mathbf{x}_t - \bar{\mathbf{x}}_t)(\mathbf{x}_t - \bar{\mathbf{x}}_t)$  is not equal to  $(\mathbf{c}_t - \bar{\mathbf{c}}_t)(\mathbf{x}_t - \bar{\mathbf{x}}_t)$  due to the (suppressed) time-dependency of the weights  $\mathbf{W}$ .

## References

- [1] Mark A Bee and Christophe Michey. The cocktail party problem: what is it? How can it be solved? And why should animal behaviorists study it? *Journal of Comparative Psychology*,

- 122(3):235, 2008.
- [2] Adelbert W Bronkhorst. The cocktail-party problem revisited: early processing and selection of multi-talker speech. *Attention, Perception, & Psychophysics*, 77(5):1465–1487, 2015.
  - [3] E Colin Cherry. Some experiments on the recognition of speech, with one and with two ears. *The Journal of the Acoustical Society of America*, 25(5):975–979, 1953.
  - [4] Robert Desimone and John Duncan. Neural mechanisms of selective visual attention. *Annual Reviews Neuroscience*, 18(1):193–222, 1995.
  - [5] Stewart H Hulse, Scott A MacDougall-Shackleton, and Amy B Wisniewski. Auditory scene analysis by songbirds: Stream segregation of birdsong by European starlings (*Sturnus vulgaris*). *Journal of Comparative Psychology*, 111(1):3, 1997.
  - [6] Aapo Hyvärinen and Patrik Hoyer. Emergence of phase-and shift-invariant features by decomposition of natural images into independent feature subspaces. *Neural computation*, 12(7):1705–1720, 2000.
  - [7] Aapo Hyvärinen and Erkki Oja. Independent component analysis: Algorithms and applications. *Neural Networks*, 13(4-5):411–430, 2000.
  - [8] Simon B Laughlin and Terrence J Sejnowski. Communication in neuronal networks. *Science*, 301(5641):1870–1874, 2003.
  - [9] David Lipshutz, Yanis Bahroun, Siavash Golkar, Anirvan M Sengupta, and Dmitri B Chklovskii. A biologically plausible neural network for multichannel canonical correlation analysis. *Neural Computation*, 33(9):2309–2352, 2021.
  - [10] David Lipshutz, Charles Windolf, Siavash Golkar, and Dmitri B Chklovskii. A biologically plausible neural network for slow feature analysis. *Advances in Neural Information Processing Systems*, 33:14986–14996, 2020.
  - [11] Josh H McDermott. The cocktail party problem. *Current Biology*, 19(22):R1024–R1027, 2009.
  - [12] Rajiv Narayan, Virginia Best, Erol Ozmeral, Elizabeth McClaine, Micheal Dent, Barbara Shinn-Cunningham, and Kamal Sen. Cortical interference effects in the cocktail party problem. *Nature Neuroscience*, 10(12):1601–1607, 2007.
  - [13] Erkki Oja and Mark Plumbley. Blind separation of positive sources by globally convergent gradient search. *Neural Computation*, 16(9):1811–1825, 2004.
  - [14] Cengiz Pehlevan and Dmitri B Chklovskii. Neuroscience-inspired online unsupervised learning algorithms: Artificial neural networks. *IEEE Signal Processing Magazine*, 36(6):88–96, 2019.
  - [15] Cengiz Pehlevan, Sreyas Mohan, and Dmitri B Chklovskii. Blind nonnegative source separation using biological neural networks. *Neural Computation*, 29(11):2925–2954, 2017.
  - [16] Mark Plumbley. Conditions for nonnegative independent component analysis. *IEEE Signal Processing Letters*, 9(6):177–180, 2002.

- [17] Mark D Plumbley. Algorithms for nonnegative independent component analysis. *IEEE Transactions on Neural Networks*, 14(3):534–543, 2003.
- [18] Mark D Plumbley and Erkki Oja. A “nonnegative PCA” algorithm for independent component analysis. *IEEE Transactions on Neural Networks*, 15(1):66–76, 2004.
- [19] Marta Rivera-Alba, Hanchuan Peng, Gonzalo G de Polavieja, and Dmitri B Chklovskii. Wiring economy can account for cell body placement across species and brain areas. *Current Biology*, 24(3):R109–R110, 2014.
- [20] Barbara G Shinn-Cunningham. Object-based auditory and visual attention. *Trends in Cognitive Sciences*, 12(5):182, 2008.
- [21] Rachel I Wilson and Zachary F Mainen. Early events in olfactory processing. *Annual Reviews Neuroscience*, 29:163–201, 2006.
- [22] Zhijian Yuan and Erkki Oja. A fastICA algorithm for non-negative independent component analysis. In *International Conference on Independent Component Analysis and Signal Separation*, pages 1–8. Springer, 2004.

## A Optimization over neural activity matrices ( $\mathbf{Y}, \mathbf{N}$ ) in the derivation of Algorithm 1

In this section, we show that when  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$  are at their optimal values, the optimal neural activity matrices ( $\mathbf{Y}, \mathbf{N}$ ) can be approximated via projected gradient descent. We first compute that optimal values of  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$ .

**Lemma 1.** *Suppose  $(\mathbf{W}_{xy}^*, \mathbf{W}_{yn}^*, \mathbf{Y}^*, \mathbf{N}^*)$  is an optimal solution of objective (5). Then*

$$\mathbf{W}_{xy}^* = \mathbf{P}\mathbf{A}^\top, \quad \mathbf{W}_{yn}^{*,\top} \mathbf{W}_{yn}^* = \mathbf{P}\mathbf{A}^\top \mathbf{A}\mathbf{P}^\top,$$

for some permutation matrix  $\mathbf{P}$ .

*Proof.* From [15, Theorem 3], we know that every solution of the objective

$$\arg \min_{\mathbf{Y} \in \mathbb{R}^{d \times T}} -\text{Tr}(\delta \mathbf{Y}^\top \delta \mathbf{Y} \delta \mathbf{X}^\top \delta \mathbf{X}) \quad \text{subject to} \quad \delta \mathbf{Y}^\top \delta \mathbf{Y} \preceq T\mathbf{I}_T \quad \text{and} \quad \mathbf{Y} = \mathbf{F}\mathbf{X}, \quad (16)$$

is of the form  $\mathbf{Y} = \mathbf{F}\mathbf{X}$ , where  $\mathbf{F}$  is a whitening matrix. In particular, since  $\mathbf{Y} = \mathbf{F}\mathbf{A}\mathbf{S}$  and  $\mathbf{S}$  also has identity covariance matrix,  $\mathbf{Y}$  is an orthogonal transformation of  $\mathbf{S}$ . Furthermore, since  $\mathbf{S}$  is well grounded, by [16, Theorem 1],  $\mathbf{Y}$  is nonnegative if and only if  $\mathbf{F}\mathbf{A}$  is a permutation matrix. Therefore, every solution  $\mathbf{Y}^*$  of the objective

$$\arg \min_{\mathbf{Y} \in \mathbb{R}_+^{d \times T}} -\text{Tr}(\delta \mathbf{Y}^\top \delta \mathbf{Y} \delta \mathbf{X}^\top \delta \mathbf{X}) \quad \text{subject to} \quad \delta \mathbf{Y}^\top \delta \mathbf{Y} \preceq T\mathbf{I}_T \quad \text{and} \quad \mathbf{Y} = \mathbf{F}\mathbf{X}, \quad (17)$$

is of the form  $\mathbf{Y}^* = \mathbf{P}\mathbf{X}$  for some permutation matrix  $\mathbf{P}$ . In addition, differentiating the expression

$$-\text{Tr}(\delta \mathbf{Y}^\top \delta \mathbf{Y} \delta \mathbf{X}^\top \delta \mathbf{X} + \delta \mathbf{N}^\top \delta \mathbf{N} (\delta \mathbf{Y}^\top \delta \mathbf{Y} - T\mathbf{I}_T)), \quad (18)$$

with respect to  $\delta\mathbf{Y}$  and setting the derivative equal to zero, we see that at the optimal value,  $\delta\mathbf{N}^{*\top}\delta\mathbf{N}^* = \delta\mathbf{X}^\top\delta\mathbf{X} = \delta\mathbf{S}^\top\mathbf{A}^\top\mathbf{A}\delta\mathbf{S}$ .

Differentiating  $L_1$  with respect to  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$ , we see that the optimal values for the synaptic weight matrices are achieved at  $\mathbf{W}_{xy} = \frac{1}{T}\delta\mathbf{Y}\delta\mathbf{X}^\top$  and  $\mathbf{W}_{yn} = \frac{1}{T}\delta\mathbf{N}\delta\mathbf{Y}^\top$ . Thus,

$$\mathbf{W}_{xy}^* = \frac{1}{T}\delta\mathbf{Y}^*\delta\mathbf{X}^\top = \frac{1}{T}\mathbf{P}\delta\mathbf{S}\delta\mathbf{S}^\top\mathbf{A}^\top = \mathbf{P}\mathbf{A}^\top,$$

and

$$\mathbf{W}_{yn}^{*\top}\mathbf{W}_{yn}^* = \frac{1}{T^2}\delta\mathbf{Y}\delta\mathbf{N}^{*\top}\delta\mathbf{N}^*\delta\mathbf{Y}^\top = \frac{1}{T^2}\mathbf{P}\delta\mathbf{S}\delta\mathbf{S}^\top\mathbf{A}^\top\mathbf{A}\delta\mathbf{S}\delta\mathbf{S}^\top\mathbf{P}^\top = \mathbf{P}\mathbf{A}^\top\mathbf{A}\mathbf{P}^\top.$$

□

Next, we show that when  $\mathbf{W}_{xy}$  and  $\mathbf{W}_{yn}$  are at their optimal values, the optimal  $(\mathbf{Y}^*, \mathbf{N}^*)$  can be approximated by running the projected gradient dynamics in Eq. (6).

**Lemma 2.** *Suppose  $\mathbf{W}_{xy} = \mathbf{P}\mathbf{A}^\top$  and  $\mathbf{W}_{yn}^\top\mathbf{W}_{yn} = \mathbf{P}\mathbf{A}^\top\mathbf{A}\mathbf{P}^\top$  for some permutation matrix  $\mathbf{P}$ . Then*

$$\mathbf{Y}^* = (\mathbf{W}_{yn}^\top\mathbf{W}_{yn})^{-1}\mathbf{W}_{xy}\mathbf{X} = \mathbf{P}\mathbf{S}, \quad \mathbf{N}^* = \mathbf{W}_{yn}\mathbf{Y}^*. \quad (19)$$

is a solution of the min-max problem

$$\min_{\mathbf{Y} \in \mathbb{R}_+^{d \times T}} \max_{\mathbf{N} \in \mathbb{R}^{m \times T}} \frac{2}{T} \text{Tr} \left( \delta\mathbf{N}^\top \mathbf{W}_{yn} \delta\mathbf{Y} - \delta\mathbf{Y}^\top \mathbf{W}_{xy} \delta\mathbf{X} - \delta\mathbf{N}^\top \delta\mathbf{N} \right) \quad \text{s.t.} \quad \mathbf{Y} = \mathbf{F}\mathbf{X}. \quad (20)$$

In particular,  $(\mathbf{Y}^*, \mathbf{N}^*)$  is the unique solution of the min-max problem

$$\min_{\mathbf{Y} \in \mathbb{R}_+^{d \times T}} \max_{\mathbf{N} \in \mathbb{R}^{m \times T}} \frac{2}{T} \text{Tr} \left( \mathbf{N}^\top \mathbf{W}_{yn} \mathbf{Y} - \mathbf{Y}^\top \mathbf{W}_{xy} \mathbf{X} - \mathbf{N}^\top \mathbf{N} \right), \quad (21)$$

which can be approximated by running the projected gradient dynamics in Eq. (6).

*Proof.* We first relax the condition that  $\mathbf{Y}$  be a nonnegative linear transformation of  $\mathbf{X}$  and consider the min-max problem

$$\min_{\mathbf{Y} \in \mathbb{R}^{d \times T}} \max_{\mathbf{N} \in \mathbb{R}^{m \times T}} \frac{2}{T} \text{Tr} \left( \delta\mathbf{N}^\top \mathbf{W}_{yn} \delta\mathbf{Y} - \delta\mathbf{Y}^\top \mathbf{W}_{xy} \delta\mathbf{X} - \delta\mathbf{N}^\top \delta\mathbf{N} \right).$$

After differentiating with respect to  $\delta\mathbf{Y}$  and  $\delta\mathbf{N}$ , we see that this objective is optimized when the centered matrices  $\delta\mathbf{Y}$  and  $\delta\mathbf{N}$  are given by

$$\delta\mathbf{Y} = (\mathbf{W}_{yn}^\top\mathbf{W}_{yn})^{-1}\mathbf{W}_{xy}\delta\mathbf{X}, \quad \delta\mathbf{N} = \mathbf{W}_{yn}\delta\mathbf{Y}.$$

Next, we see that the above relations for the centered matrices hold when  $\mathbf{Y}$  and  $\mathbf{N}$  are given by Eq. (19), where we have used the fact that  $\mathbf{W}_{xy} = \mathbf{P}\mathbf{A}^\top$  and  $\mathbf{W}_{yn}^\top\mathbf{W}_{yn} = \mathbf{P}\mathbf{A}^\top\mathbf{A}\mathbf{P}^\top$ . Note that  $\mathbf{Y}$  is a linear transformation of  $\mathbf{X}$  and  $\mathbf{Y}$  is nonnegative since it is a permutation of the nonnegative sources. It follows that  $(\mathbf{Y}, \mathbf{N})$  is also a solution to the *constrained* min-max problem (20). Finally, differentiating the objective in Eq. (21) with respect to  $\mathbf{Y}$  and  $\mathbf{N}$ , we see that the optimal  $\mathbf{Y}$  and  $\mathbf{N}$  are again given by Eq. (19). □

## B Decoupling the interneuron synapses

The NICA algorithm derived in section 4.1 requires the interneuron-to-output neuron synaptic weight matrix  $\mathbf{W}_{ny}$  to be the the transpose of the output neuron-to-interneuron synaptic weight matrix  $\mathbf{W}_{yn}$ . Enforcing this symmetry via a centralized mechanism is not biologically plausible, and is commonly referred to as the weight transport problem.

Here, we show that the symmetry of the 2 weights asymptotically follows from the learning rules in Algorithm 1, even when the symmetry does not hold at initialization. Let  $\mathbf{W}_{ny,0}$  and  $\mathbf{W}_{yn,0}$  denote the initial values of  $\mathbf{W}_{ny}$  and  $\mathbf{W}_{yn}$ . Then, in view of the updates rules Algorithm 2, the difference  $\mathbf{W}_{ny} - \mathbf{W}_{yn}^\top$  after  $t$  updates is given by

$$\mathbf{W}_{ny} - \mathbf{W}_{yn}^\top = (1 - \eta)^t (\mathbf{W}_{ny,0} - \mathbf{W}_{yn,0}^\top).$$

In particular, the difference decays exponentially.

## C Details of numerical experiments

The simulations were performed on an Apple machine with a 2.8 GHz Quad-Core Intel Core i7 processor.

### C.1 Mixing matrices

We used the  $3 \times 3$  mixing matrix for the 3-dimensional random uniform sources that was used in [15]:

$$\mathbf{A} = \begin{bmatrix} 0.031518 & 0.38793 & 0.061132 \\ -0.78502 & 0.16561 & 0.12458 \\ 0.34782 & 0.27295 & 0.67793 \end{bmatrix}.$$

The  $10 \times 10$  mixing matrix for the 10-dimensional random uniform sources is as follows (entries are rounded to 2 decimal places for space considerations):

$$\mathbf{A} = \begin{bmatrix} -1.61 & 0.11 & 0.11 & 1.26 & -0.01 & -1.66 & 0.45 & 0.48 & 0.93 & -0.57 \\ -0.95 & -0.05 & 0.35 & -0.68 & 1.14 & 0.71 & -0.38 & -0.20 & -0.20 & 2.02 \\ 0.54 & 2.16 & 0.06 & -0.08 & 0.36 & -0.16 & -0.22 & -1.82 & -0.22 & 0.40 \\ -0.98 & -0.12 & -1.45 & -0.58 & -0.56 & 0.34 & -0.51 & 0.19 & -0.44 & -0.15 \\ -0.87 & 0.54 & 0.68 & 1.28 & 0.63 & 1.04 & -0.81 & 1.08 & -0.65 & -0.30 \\ 0.91 & 0.84 & 0.45 & -0.31 & -0.14 & -1.46 & -0.18 & 0.48 & -0.41 & 0.75 \\ -1.20 & 1.29 & 0.39 & -1.40 & 0.84 & -2.32 & -1.54 & -0.26 & -1.99 & -0.34 \\ 1.34 & 0.75 & -1.29 & -0.63 & -1.63 & -1.05 & 0.07 & 0.09 & -0.67 & 0.28 \\ -0.32 & -0.38 & -0.11 & 1.18 & -0.41 & 0.58 & -0.92 & 1.09 & 0.41 & 1.29 \\ 2.04 & 2.00 & -0.50 & 0.78 & -0.65 & -0.93 & 0.42 & -1.69 & -1.16 & -0.68 \end{bmatrix}.$$

The  $3 \times 3$  mixing matrix for the 3-dimensional natural image sources is given by:

$$\mathbf{A} = \begin{bmatrix} 0.71964649 & -1.55757433 & -1.94561985 \\ -1.77115767 & -0.99092683 & 0.35559978 \\ -0.78408667 & 1.09213136 & -1.36539258 \end{bmatrix}.$$

## C.2 Implementation of algorithms

For each of the algorithms that we implement, we use a time-dependent learning rate of the form:

$$\eta_t = \frac{\eta_0}{1 + \gamma t}. \quad (22)$$

To choose the parameters, we perform a grid search over  $\eta_0 \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$  and over  $\gamma \in \{10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}\}$ . In Table 1 we report the best performing hyperparameters we found for each algorithm. We now detail our implementation of each algorithm.

1. **Bio-NICA with interneurons (Algorithm 1):** The neural outputs were computed using the quadratic convex optimization function `solve_qp` from the Python package `quadprog`. After each iteration, we checked if any output neuron had not been active up until that iteration. If so, we flipped the sign of its feedforward inputs. In addition, if the norm of one of the row vectors of  $\mathbf{W}_{xy}$  fell below 0.1, we would replace the row vector with a random vector to avoid the row vector becoming degenerate; and if a singular value of  $\mathbf{W}_{xy}$ ,  $\mathbf{W}_{yn}$  or  $\mathbf{W}_{ny}$  fell below 0.01, we replaced the singular value with 1 (we checked every 100 iterations).
2. **Bio-NICA with 2-compartmental neurons (Algorithm 2):** The neural outputs were computed using the quadratic convex optimization function `solve_qp` from the Python package `quadprog`. We used the time-dependent learning rate of Eq. (22) and included  $\tau \in \{0.01, 0.03, 0.05, 0.08, 0.1, 0.3, 0.5, 0.8, 1, 3\}$  in the grid search to find the best performance. After each iteration, we checked if any output neuron had not been active up until that iteration. If so, we flipped the sign of its feedforward inputs. In addition, if a eigenvalue of  $\mathbf{W}_{zz}$  fell below 0.01, we replaced the eigenvalue with 1 to prevent  $\mathbf{W}_{zz}$  from becoming degenerate (we checked every 100 iterations).
3. **2-layer NSM:** We implemented the algorithm in [15] with time-dependent learning rates. For the whitening layer, we used the optimal time-dependent learning rate reported in [15]:  $\zeta_t = 0.01/(1 + 0.01t)$ . For the NSM layer, we used the time-dependent learning rate of Eq. (22). To compute the neuronal outputs, we used the quadratic convex optimization function `solve_qp` from the Python package `quadprog`. After each iteration, we checked if any output neuron had not been active up until that iteration. If so, we flipped the sign of its feedforward inputs.
4. **Nonnegative PCA (NPCA):** We use the online version given in [18]. The algorithm assumes the inputs are noncentered and whitened. We performed the noncentered whitening offline. After each iteration, we checked if any output neuron had not been active up until that iteration. If so, we flipped the sign of its feedforward inputs.

	Alg. 1 ( $\eta_0, \gamma$ )	Alg. 2 ( $\eta_0, \gamma, \tau$ )	2-layer NSM ( $\eta_0, \gamma$ )	NPCA ( $\eta_0, \gamma$ )
$d = 3$	$(10^{-2}, 10^{-3})$	$(10^{-1}, 10^{-2}, 0.8)$	$(10^{-1}, 10^{-7})$	$(10^{-2}, 10^{-5})$
$d = 10$	$(10^{-2}, 10^{-3})$	$(10^{-3}, 10^{-4}, 0.03)$	$(10^{-1}, 10^{-6})$	$(10^{-2}, 10^{-5})$
Images	$(10^{-3}, 10^{-6})$	$(10^{-2}, 10^{-4}, 0.5)$	$(10^{-1}, 10^{-6})$	$(10^{-3}, 10^{-5})$

Table 1: Optimal hyperparameters used for Alg. 1, Alg. 2, 2-layer NSM, and NPCA.