

## Exploring the Role of Intrinsic Plasticity for the Learning of Sensory Representations

Nicholas J. Butko<sup>1</sup> and Jochen Triesch<sup>1,2</sup>

1- University of California, San Diego - Cognitive Science  
9500 Gilman Dr. 0515, La Jolla, CA, 92093-0515 - USA

2- Frankfurt Institute for Advanced Study  
Max-von-Laue-Str. 1, 60438 Frankfurt am Main - Germany

**Abstract.** Intrinsic plasticity (IP) refers to a neuron's ability to regulate its firing activity by adapting its intrinsic excitability. Previously, we showed that model neurons combining IP with Hebbian synaptic plasticity can adapt their weight vector to discover heavy-tailed directions in the input space. In this paper we consider networks of coupled model neurons and show how a population of such units can solve a standard non-linear ICA problem. We also present a simple model for the formation of maps of oriented receptive fields in primary visual cortex. Together, our results indicate that intrinsic plasticity may play an important role for learning efficient representations in populations of cortical neurons.

### 1 Introduction

Biological neurons regulate their firing activity by adapting their intrinsic excitability. Such *intrinsic plasticity* (IP) seems to be a ubiquitous phenomenon in the brain [1]. For example, Desai *et al.* showed that neurons that had been prevented from spiking for two days increased their response to current injection [2]. It is frequently assumed that IP contributes to the homeostasis of a neuron's firing activity. Baddeley *et al.* found that neurons in visual cortical areas show exponential distributions of their firing rate, which is thought to maximize a neuron's information transfer given a fixed energy budget [3]. This is because the exponential distribution has the maximum entropy among all distributions of a positive random variable with a fixed mean. It has been speculated that IP may be instrumental in achieving approximately exponential firing rate distributions in cortical neurons [4]. We have recently shown that IP that drives a neuron to exhibit an exponential firing rate distribution can synergistically interact with Hebbian learning at the synapses. These two processes lead to the discovery of heavy-tailed directions in the input space [5]. In this paper we extend these results to populations of neurons with IP. Our specific goal is to explore the potential role of IP for learning efficient map-like representations for sensory stimuli.

Computational models of the emergence of sensory representations in the brain abound. Frequently, they fall into one of two categories: *functional models* or *mechanistic models*. Mechanistic models start from neuroscientific data about the structure of cortical networks and cortical plasticity mechanisms (cell types,

connection patterns, plasticity rules, ...) which are distilled into simplified models. These models are trained on actual sensory data or noise patterns and the learned representations can be compared to neuroscientific data. If the resulting representations are similar to those found in the brain then this provides evidence that the processes in the brain have been accurately captured, but it does not clarify why the brain operates this way or in what sense the brain's solution may be optimal. An example of a model of this kind by Linsker is [6], where V1-style orientation columns are learned from random prenatal visual noise through Hebbian learning. Later Miller extended this work to learn many of the various map-structures in V1, and used model neurons that were somewhat more plausible [7].

Functional models focus on the abstract computational goal of the problem. For the case of learning sensory representations they start by asking: what is the *optimal* way to represent sensory stimuli, e.g. natural images, where optimality is usually defined with respect to certain statistical criteria (e.g., sparseness, independence, temporal coherence, ...) and additional constraints. Algorithms are derived to learn the optimal solution to the problem which can again be compared to neuroscientific data. If the found solution resembles the biological solution, then this provides evidence that the brain may in fact be trying to optimize a similar objective function. Through what mechanisms the brain may achieve this goal is typically not answered, however.

A central idea in many functional models is information maximization. According to this idea, individual neurons should spread out their responses in dense regions of the input space and compress responses in sparse regions. In effect, this maps the input to a uniform output distribution, maximizing entropy. Laughlin showed that blowfly Large Monopolar Cells have been adapted so that their input/output transfer functions nearly optimally represent the contrast statistics of their environment [8]. Bell & Sejnowski showed that the same information maximization principle can be applied to the independent component analysis (ICA) problem. They applied their technique to natural images and found oriented, bandpass sources [9] similar to those observed in V1. Olshausen & Field showed that oriented, bandpass receptive fields also arise when optimizing image reconstruction error subject to lifetime sparseness constraints [10]. They imposed a sparse prior on the contribution of each basis function in a generative model with the intuition that among the space of possible sources of an image, each one is present only rarely. This intuition is confirmed by Ruderger, who showed that bandpass filter responses to natural images follow an exponential distribution [11].

The model we present in the following attempts to bridge the gap between mechanistic and functional models. On the one hand, it has a clear connection to the idea of information maximization [5]. On the other hand, it has a mechanistic formulation that is biologically plausible because it makes use of information that is local in time and space, and uses patterns of lateral connections characteristic of neural populations. While similar bridges have been attempted before, e.g. [12], our model is the first to utilize IP as a fundamental

mechanism for the learning of sensory representations.

## 2 Network Model with Intrinsic Plasticity

We consider a network of units learning to represent a sensory input signal  $\mathbf{x}$ . The activity of unit  $i$  in the network is given by:

$$y_i(h_i) = [1 + \exp(-a_i h_i - b_i)]^{-1}, \text{ where } h_i = \mathbf{x} \cdot \mathbf{w}_i, \quad (1)$$

where  $\mathbf{w}_i$  is the neuron's weight vector, and  $a_i$  and  $b_i$  are adjustable parameters of the neuron's transfer function that are controlled by IP. In particular  $a_i$  and  $b_i$  are subject to the following learning rule:

$$\begin{aligned} \Delta a_i &\stackrel{(a)}{=} \eta_{\text{IP}} [a_i^{-1} + h_i - (2 + \mu^{-1})h_i y_i + \mu^{-1}h_i y_i^2] \\ \Delta b_i &\stackrel{(b)}{=} \eta_{\text{IP}} [1 - (2 + \mu^{-1})y_i + \mu^{-1}y_i^2], \end{aligned} \quad (2)$$

where  $\eta_{\text{IP}}$  is a small learning rate and  $\mu$  is the desired mean activity of all units. As derived in [13], this learning rule has the effect of making the distribution of  $y_i$  a sparse, approximately exponential distribution, thereby maximizing the unit's entropy given a fixed average activity. Note that this rule is *local* in space and time, making it physiologically plausible.

Plasticity of the weight vectors  $\mathbf{w}_i$  is modeled with a Hebbian learning rule. In [5], we considered a single unit learning rule of the form  $\Delta \mathbf{w} \propto \mathbf{x}y$ . We showed that the coupling of IP with this form of Hebbian learning allowed the unit to discover heavy-tailed directions in the input. To extend this model to a population of model neurons, we introduce a *neighborhood function*  $\mathcal{N}$  as commonly used in self-organizing maps. The value of the neighborhood function for neuron  $i$  is determined by its activity  $y_i$  and the activities of all other neurons, i.e.  $\mathcal{N}(y_i; \mathbf{y})$ . Specific forms of  $\mathcal{N}$  are introduced below. After each stimulus presentation, the weights are updated according to:

$$\Delta \mathbf{w}_i \stackrel{(a)}{=} \mathbf{x}y\mathcal{N}(y_i; \mathbf{y}), \quad \mathbf{w}_i \stackrel{(b)}{\Leftarrow} \frac{\mathbf{w}_i + \eta_{\text{Hebb}}\Delta \mathbf{w}_i}{\|\mathbf{w}_i + \eta_{\text{Hebb}}\Delta \mathbf{w}_i\|}. \quad (3)$$

where  $\eta_{\text{Hebb}}$  is a learning rate,  $\Leftarrow$  denotes assignment, and the normalization in (3b) mimics competition between synapses on a neuron's dendritic tree [7].

## 3 Experiments

### 3.1 The "bars" problem

The bars problem is standard non-linear ICA problem introduced by Földiák [14]. Horizontal and vertical bars are presented on an R-by-R retina. The presence or absence of a bar is independent of that of any other bars. The unsupervised learning problem is to learn filters that correspond to the individual independent components, *i.e.* the bars. The problem is non-linear because the



Fig. 1: A population of 20 model neurons has learned all independent sources in the bars problem. The weight vector of each unit has discovered a single bar. Parameters were:  $\beta = 0.2$ ,  $\eta_{\text{Hebb}} = 0.01$ ,  $\eta_{\text{IP}} = 0.005$ , and  $\mu = 0.1$ .

pixel at the intersection of two bars is just as bright as any other pixel of the bars, not twice as bright. In our previous work [5], we showed that a single model neuron with IP and Hebbian learning discovers one of the bars when exposed to stimuli from the bars problem. Here we use a population of units to learn the complete problem. We use a retina of 10-by-10 pixels and the probability of any of the 20 bars occurring in a given stimulus is 10%. Since we want filters that respond highly when bars are present and not otherwise, the desired mean firing rate  $\mu$  is set at 10%.  $\mathcal{N}$  is chosen to enforce a winner-take-all competition between the units, so that the maximally activated neuron updates its weight vector in a standard Hebbian fashion, and all other units update their weight in an anti-Hebbian manner regulated by a decorrelation parameter  $\beta$ :

$$\mathcal{N}_{\text{bars}}(y_i; \mathbf{y}) = \begin{cases} 1 & : y_i = \max(\mathbf{y}) \\ -\beta & : \text{else} \end{cases} . \quad (4)$$

All units update their intrinsic parameters independently, as described in Equation 2(a&b). Figure 1 shows the parameters used and the results of learning. With the specified parameters, we found that learning was completely successful in 19 out of 20 trials. Varying the learning rates  $\eta_{\text{Hebb}}$  and  $\eta_{\text{IP}}$  affected learning little, provided both remained above 0. When  $\mu$  was 0.05, redundant filters were learned, and when it was 0.2, multiple bars were represented within single filters. This suggests that when the true mean of the components is unknown, it may be a better strategy to choose  $\mu$  too high rather than too low. This way, all true sources will likely be captured because individual filters each learn to represent several sources. Learning substantially worsened when  $\beta$  was less than 0.05 or greater than 0.3, with effects similar to  $\mu$  too low or too high respectively.

### 3.2 Modeling the Emergence of Orientation Maps

Receptive fields of simple cells in V1 are oriented, localized, and bandpass. Neurons are arranged in a columnar fashion that reflects their orientation preference. For modelling the emergence of orientation columns, we consider our populations of neurons to be located on a two-dimensional sheet, with neuron  $i$  at grid position  $(j, k)_i \in \mathbb{N} \times \mathbb{N}$  after the fashion of a Kohonen feature map. The most active unit exhibits a center-surround influence on learning in its neighbors according to a difference of Gaussians (DoG) neighborhood function centered at



