

Article

The Fisher Information as a Neural Guiding Principle for Independent Component Analysis.

Rodrigo Echeveste^{1*}, Samuel Eckmann¹, and Claudius Gros¹

¹ Institute for Theoretical Physics, Goethe University Frankfurt, Germany

* Author to whom correspondence should be addressed; echeveste[[@](mailto:)]itp.uni-frankfurt.de

Version June 5, 2015 submitted to *Entropy*. Typeset by \LaTeX using class file *mdpi.cls*

Abstract: The Fisher information constitutes a natural measure for the sensitivity of a probability distribution with respect to a set of parameters. An implementation of the stationarity principle for synaptic learning in terms of the Fisher information results in a Hebbian self-limiting learning rule for synaptic plasticity. In the present work, we study the dependence of the solutions to this rule in terms of the moments of the input probability distribution and find a preference for non-Gaussian directions, making it a suitable candidate for independent component analysis (ICA). We confirm in a numerical experiment that a neuron trained under this rules is able to find the independent components in the non-linear bars problem.

The specific form of the plasticity rule depends on the transfer function used, becoming a simple cubic polynomial of the membrane potential when the rescaled error function is employed as transfer function. The cubic learning rule is also an excellent approximation for other transfer functions, as the standard sigmoidal, and can be used to show analytically that the proposed plasticity rules are selective for directions in the space of presynaptic neural activities characterized by a negative excess kurtosis.

Keywords: Fisher Information; Guiding Principle; Excess Kurtosis; Objective Functions; Synaptic Plasticity; Hebbian Learning; Independent Component Analysis

18 1. Introduction

19 Many living systems, such as neurons and neural networks as a whole, are guided by overarching
 20 constraints or principles. Energy [1] and metabolic costs [2] for the information processing of the brain
 21 act in this context as basic physiological constraints for the evolution of neural systems [3,4], making,
 22 e.g., efficient coding [5,6] a viable strategy.

23 Metabolic costs can be considered as special cases of objective functions [7], which are to be
 24 minimized. Objective functions can however be used in many distinct settings. For example to guide
 25 data selection [8] in engineering applications [9], or to guide learning within neural networks in terms
 26 of synaptic plasticity rules by minimizing an appropriate combination of moments of the neural activity
 27 [10].

28 Objective functions can also be formulated in terms of the probability distribution of the neural
 29 activity [11], allowing the formulation of information theoretical generative functionals for behavior
 30 in general [12–16], for neural activity [17,18], for the derivation of neural plasticity rules in terms of
 31 maximizing the relative information entropy [19], or the mutual information [20–22], and for variational
 32 Bayesian tasks of the brain through free energy minimization [23].

33 1.1. Combining objective functions

A fundamental question in the context of guiding principles for dynamical systems regards the
 combination of several distinct, and possibly competing objective functions. For a survey of optimization
 in the context of multiple objective function, see [9]. For discreteness, we start by considering a generic
 neural model in which the state of the system is determined by the neural activity y_i of neuron i , by the
 intrinsic parameters $a_i^k = (\hat{a})_i^k$ (with $k = 1, 2, \dots$ indexing the different internal degrees of freedom)
 of the neurons and by the inter-neural synaptic connectivity matrix $w_{ij} = (\hat{w})_{ij}$. Within the objective
 functional approach one considers evolution equations

$$\begin{cases} \dot{y}_i &= -\frac{\partial}{\partial y_i} \mathcal{F}^{act}(\mathbf{y}, \hat{a}, \hat{w}) \\ \dot{a}_i^k &= -\frac{\partial}{\partial a_i^k} \mathcal{F}^{int}(\mathbf{y}, \hat{a}, \hat{w}) \\ \dot{w}_{ij} &= -\frac{\partial}{\partial w_{ij}} \mathcal{F}^{syn}(\mathbf{y}, \hat{a}, \hat{w}) \end{cases} \quad (1)$$

for the full dynamical neural system, where there is a specific objective function $\mathcal{F}^\alpha(\mathbf{y}, \hat{a}, \hat{w})$ for every
 class $\alpha \in \{act, int, syn\}$ of dynamical variables. It is important to note, that an overarching objective
 function like

$$\mathcal{F}^{act}(\mathbf{y}, \hat{a}, \hat{w}) + \mathcal{F}^{int}(\mathbf{y}, \hat{a}, \hat{w}) + \mathcal{F}^{syn}(\mathbf{y}, \hat{a}, \hat{w}), \quad (2)$$

34 does generically not exist. In a biological system, each objective function \mathcal{F}^α may represent a different
 35 regulatory mechanism whose coupling occurs only through the biological agent itself [18]. Indeed,
 36 how exactly these mechanisms interact in neural systems, when formulated in terms of learning rules
 37 for intrinsic and synaptic plasticity, been subject of study in recent years [19,24]. Furthermore, for a
 38 stationary input distribution, such an overarching functional would result in a gradient system having
 39 only point attractors, since limit cycles are not possible in gradient systems [25]. Therefore, a formalism

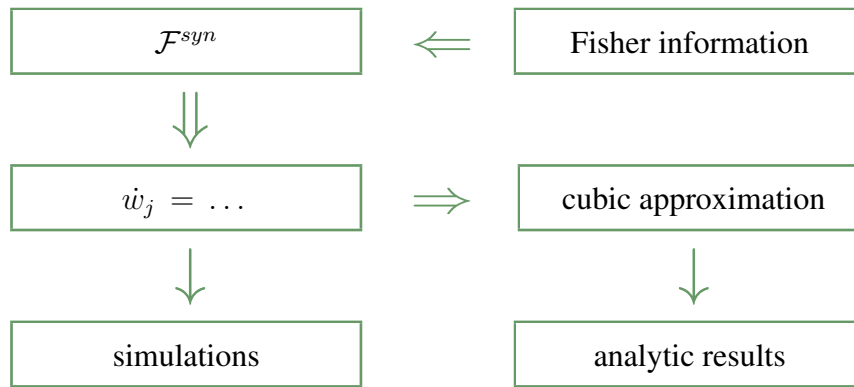


Figure 1. Organigram of the approach followed. The objective function \mathcal{F}^{syn} for synaptic plasticity studied here can be motivated by the Fisher information for the synaptic flux. The resulting plasticity rule \dot{w}_j for the synaptic weights will then be investigated both through simulations and using a cubic approximation in x (which becomes exact, when using the error functions as a transfer function $y(x) = \sigma(x - b)$, see Section 2.2), which allows to derive analytic results for the dependence of the synaptic adaption with respect to the kurtosis of the input statistics.

40 aiming to reproduce the wide variety of behaviors found in natural neural systems cannot be formulated
 41 as a gradient system of a single overarching objective function.

42 One needs to keep in mind, however, that it is often not possible to evaluate rigorously gradients of
 43 non-trivial objective functions, as in (1). Generating functionals are hence implemented, in many cases,
 44 only approximatively. For the case of information theoretical incentives, as considered in the present
 45 work, objective functions are, in addition, formulated in terms of time-averaged statistical properties and
 46 the gradient descent can hence be achieved only through a corresponding time average.

47 1.2. Hebbian learning in neural networks

48 Within the neurosciences, synaptic plasticity is generally studied under the paradigm of Hebbian
 49 learning [26], stating generically that neurons that fire together, wire together. Depending on whether
 50 one considers the frequency (also denoted firing rate), or the timing, of spikes, this principle can have
 51 different interpretations. In terms of the firing frequency of neurons, Hebbian learning is understood
 52 as a strengthening of the synaptic connection between two neurons (known as potentiation) when both
 53 neurons have simultaneously a high activity level or a weakening (depression) of the connectivity if the
 54 respective periods of high and low activity do not match [10,27]. In the case of Spike Timing Dependent
 55 Plasticity (STDP) [28,29], where synaptic modification is expressed as a function of the precise timing
 56 of spikes, on the other hand, the principle of Hebbian learning is understood in terms of causality, stating
 57 that a directional synaptic connection should be potentiated if two neurons fire in a causal order, and
 58 depressed otherwise [30,31].

59 In the present work we consider rate encoding neurons and therefore formulate plasticity in terms of
 60 an information theoretical measure of the activity distribution, or alternatively, in terms of the moments
 61 of this distribution. While the requirement of any such rule to respect the Hebbian principle of learning

will naturally constraint the manifold of learning rules, the particular details of each rule will determine its functionality. Oja's rule [27], for instance, is tailored to find the first principal component of a multi-dimensional input distribution. The rules we present in this paper, while able to find the first principal component of a distribution under certain conditions, as we show in [32], will generically perform an independent component analysis by selecting directions of maximal non-Gaussianness.

1.3. Instantaneous single neuron

In order to concentrate on the generating principle for synaptic plasticity, we consider here a single instantaneous point neuron, defined by an activity level y ,

$$y = \sigma(x - b), \quad \sigma(z) = \frac{1}{1 + e^{-z}}, \quad x = \sum_{j=1}^{N_w} w_j (y_j - \bar{y}_j), \quad (3)$$

representing the average firing rate of the neuron, where $\sigma(z)$ is a monotonically increasing sigmoidal transfer function, denoted in physics as the Fermi function, that converts the total weighed input x (also referred to as the membrane potential of the neuron) into an output activity. N_w represents the number of incoming inputs y_j , which represent in this case either an external input or the activities of other neurons in a network. b is a bias in the neuron's sensitivity and \bar{y}_j represents the (trailing) average of the input activity, such that only deviations from this average contribute to the integrated input. An objective function for the neural activity is, in this case, not present and the evolution equations (1) reduce to

$$\begin{cases} \dot{b} = -\epsilon_b \frac{\partial}{\partial b} \mathcal{F}^{int}(y, b, \mathbf{w}) \\ \dot{w}_j = -\epsilon_w \frac{\partial}{\partial w_j} \mathcal{F}^{syn}(y, b, \mathbf{w}) \end{cases} \quad \text{with} \quad y = \sigma \left(\sum w_j (y_j - \bar{y}_j) - b \right), \quad (4)$$

where we are left only with the objective function for the intrinsic \mathcal{F}^{int} and for the synaptic \mathcal{F}^{syn} plasticity. Here we have, with ϵ_b and ϵ_w , separated the adaption rates from the definition of the respective objective functions.

1.4. Information theoretical incentives for synaptic plasticity

In the context of stochastic information processing systems, tools from information theory, such as the entropy of a given code or the mutual information between input and output [8], permit to formulate objective functions for learning and plasticity in terms the probability distributions of the stochastic elements that constitute the system [6,10–12,14,18,19]. Principles such as maximizing the output entropy of a system to improve the representational richness of the code [19], maximal information transmission for signal separation and deconvolution in networks [33], or maximal predictive information within the sensorimotor loop as a guiding principle to generate behavior [34], have proved successful in the past in both generating new approaches to learning and plasticity and in furthering the understanding of already available rules, integrating them into a broader context by formulating them in terms of a guiding principle [10].

In the present work we discuss a novel synaptic plasticity rule [32] resulting in self-limiting Hebbian learning and its interaction with known forms of intrinsic plasticity [19,35]. The novelty of this approach

84 relies on the objective function employed, which can be derived from the Fisher information [36] of the
 85 output probability distribution with respect to a synaptic flux operator, as shown in Section 2.3. In
 86 previous work [32,37], we had shown numerically, how a non-linear point neuron employing the Fermi
 87 function or the inverse tangent as its transfer function from membrane potential to output activity, was
 88 able to find the first principal component of an ellipsoidal input distribution but showed a preference for
 89 directions of large negative excess kurtosis otherwise. In the present work, however, we show how, by
 90 use of the rescaled error function as a transfer function $y(x) = \sigma(x - b)$, the resulting learning rule, while
 91 qualitatively equivalent to the ones previously studied, takes the form of a simple cubic polynomial in
 92 the membrane potential x . This fact, as we will show, represents an important step forward, since it allows
 93 us to study the attractors of the learning procedure and their stability analytically. In particular, the rule
 94 is shown to have interesting properties in terms of the moments of the input distributions, resulting in a
 95 useful tool for independent component analysis, which is thought to be of high relevance for biological
 96 cognitive systems [27,38,39].

97 It is worth mentioning at this point that, while the Fisher information is usually associated with the
 98 task of parameter estimation via the Cramér-Rao bound [40–42], it generically encodes the sensitivity
 99 of a probability distribution with respect to a given parameter, making it also a useful tool, both in
 100 the context of optimal population codes [43–45], or as here, for the formulation of objective functions.
 101 Indeed, this procedure has been successfully employed in the past in other fields, to derive, for instance,
 102 the Schrödinger Equation in Quantum Mechanics [46].

103 We will start in the following, as illustrated in Fig. 1, with the primary objective function \mathcal{F}^{syn}
 104 for synaptic plasticity, discussing its relation with the Fisher information for the synaptic flux later
 105 on in Sect. 2.3. Simulation results of the synaptic adaption rules will then be presented in Sect. 3,
 106 in comparison with the results obtained using an analytically treatable cubic approximation in the
 107 membrane potential, as presented in Sect. 2.1.

108 2. Objective functions for synaptic plasticity

Our primary objective function for synaptic plasticity is [32]:

$$\mathcal{F}^{syn} = E \left[(N + x(1 - 2y))^2 \right], \quad (5)$$

109 where $E[\cdot]$ denotes the expected value with respect to the input probability distribution, which can be
 110 equated to a time-average whenever the input probability distributions are stationary. The objective
 111 function \mathcal{F}^{syn} can be expressed entirely in terms of either x or y , which are related by (3). The current
 112 form (5) is chosen just for clarity. In Sect. 2.3, we show how \mathcal{F}^{syn} can be derived from the Fisher
 113 information with respect to an operator denoted the synaptic flux operator.

From (5) one can derive easily, via stochastic gradient descent, the update rule

$$\dot{w}_j = \epsilon_w G(x) H(x) (y_j - \bar{y}_j), \quad (6)$$

with $H(x) = -G'(x)$ and

$$G(x) = N + x(1 - 2y(x)), \quad H(x) = (2y(x) - 1) + 2x(1 - y(x))y(x). \quad (7)$$

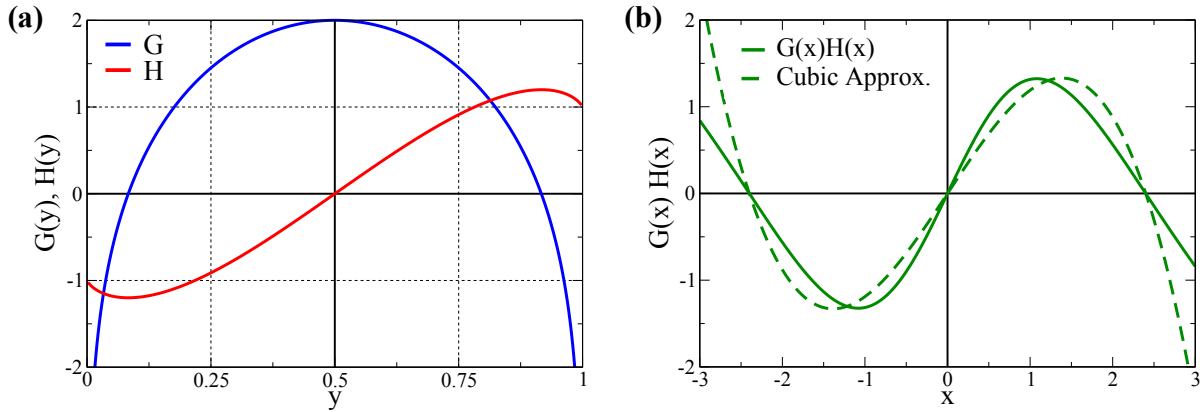


Figure 2. (a) The plasticity functions G and H , as defined by (7), here expressed entirely in terms of the output activity $y \in [0, 1]$, for clarity. H represents the Hebbian contribution of the rule, with G acting as a limiting factor, reverting the sign of (6) for activity values close to 0/1. (b) Plot of the learning rule (6) together with the cubic approximation (8), expressed this time as a function of the membrane potential x . Parameters: $b = 0$ and $N = 2$.

N is a parameter that allows to shift the positions of the roots of G . The synaptic functions G and H can also be entirely expressed either in terms of x or y , as shown in Fig. 2. For $N = 2$ the two synaptic functions are proportional to each other's derivatives, $G(x) = 2y(1 - y)H'(x)$, viz they are conjugate to each other [32].

$H(y)$ is an essentially linear function of positive slope throughout most of the activity range of the neuron, see Fig. 2 (a), saturating only for $y \rightarrow 1/0$. The product $H(y)(y_j - \bar{y}_j)$ constitutes hence the Hebbian part of the plasticity rule (6), resulting in an increase of the synaptic weight whenever the input y_j and the output y are correlated.

The plasticity function $G(y)$, however, reverts the sign of the learning rule if the activity level approaches the extremes, $y \rightarrow 1/0$, serving hence as a limiting factor. The process hence adapts the synaptic weights over time such that the the membrane potential x remains close to the roots of $G(x)$. The synaptic weight will consequently also remain finite, making the adaption rules (6) self-limiting.

2.1. Cubic approximation

In order to describe the stationary solutions of (6) in terms of the moments of the input probability distributions, we consider a polynomial expansion in x . The two roots $\pm x_0$ of the limiting function $G(x)$, compare Fig. 2 (a), are symmetric for the case $b = 0$, considered in the following, and scale $\sim N$ for large N [32]. The Hebbian function $H(x)$ has, on the other hand, only a single root at $x = 0$ (viz at $y = 0.5$), for $b = 0$. We are then led to the cubic approximation

$$\begin{aligned} \dot{w}_j &= \epsilon_w G(x)H(x)(y_j - \bar{y}_j) \approx -\epsilon_w x(x - x_0)(x + x_0)(y_j - \bar{y}_j)/N^2 \\ &= \epsilon_w x(x_0^2 - x^2)(y_j - \bar{y}_j)/N^2 \end{aligned} \quad (8)$$

of (6). Note, that the scaling factor $1/N^2 > 0$ could also be absorbed into the adaption rate ϵ_w . In Fig. 2 (b) the learning rule (6) is compared to the cubic approximation (8).

For convenience we denote $\gamma_j = (y_j - \bar{y}_j)$, and compute with

$$\langle \dot{w}_j \rangle = \epsilon_w \frac{1}{N^2} E \left[\gamma_j \left[\left(\sum_{i=1}^{N_w} w_i \gamma_i \right) x_0^2 - \left(\sum_{i=1}^{N_w} w_i \gamma_i \right)^3 \right] \right], \quad (9)$$

the time-averaged expectation value of the synaptic weight changes, equating the time average with the statistical average $E[\cdot]$ over the distributions $p(y_j)$ of the input activities y_j . We now assume uncorrelated and symmetric input distributions,

$$E[\gamma_i \gamma_j] = 0 = E[\gamma_i^k], \quad k = 1, 3, 5, \dots$$

129 The odd moments hence vanish. Here it is important to note that any learning rule defined purely in terms
130 of the overall input $x = \mathbf{w} \cdot \boldsymbol{\gamma}$, will be fully rotational invariant. Therefore, the result does not depend on
131 the direction one chooses for the PCs. In particular, if one chooses the principal components to lie along
132 the axes of reference, one can eliminate the linear correlation terms, without loss of generality.

The synaptic weights are quasi-stationary for small adaption rates $\epsilon_w \rightarrow 0$ and we obtain

$$\langle \dot{w}_j \rangle = \epsilon_w \frac{1}{N^2} w_j \sigma_j^2 (x_0^2 - w_j^2 \sigma_j^2 K_j - 3\Phi) \quad (10)$$

from (9), where we have defined with

$$\sigma_j^2 = E[\gamma_j^2], \quad K_j = \frac{E[\gamma_j^4]}{\sigma_j^4} - 3, \quad \Phi = \sum_j w_i^2 \sigma_j^2 \quad (11)$$

133 the standard deviation (SD) σ_j of the j -the input, the excess kurtosis K_j , and the weighed average Φ of
134 the afferent standard deviations.

135 2.1.1. Scaling of dominant components

The stationary solutions w_j^* of (10) satisfy

$$w_j^* = 0 \quad \vee \quad w_j^{*2} \sigma_j^2 K_j = x_0^2 - 3\Phi, \quad (12)$$

136 which implies, that there is a competition between small components $w_j^* \approx 0$ of the synaptic weight
137 vector and large components.

In [32], the authors trained a neuron with ellipsoidal distributions, consisting of normal distributions truncated to $[0, 1]$, with one direction having a large SD σ_1 (the first principal component, or FPC) and the rest of the directions having a small SD. In this context, the weight vector aligns to the FPC, resulting in one large weight (w_1). All other synaptic weight adapt to small values. Solving (12) for the large component yields

$$|w_1^{cub}| = \frac{x_0}{\sigma_1 \sqrt{K_1 + 3}}. \quad (13)$$

138 We note that the excess kurtosis is bounded from below [47], $K \geq -2$ (the probability distribution
139 having the lowest possible excess kurtosis of -2 is the bimodal distribution made of two δ -peaks) and
140 that consequently $K + 3 > 0$.

141 In Sect. 3.1, a quantitative comparison between (13) and the numerical result of the learning rule is
142 presented.

143 2.1.2. Sensitivity to the excess kurtosis

We are interested now in examining the stability of the solutions obtained via the cubic approximation, in order to explain why a particular solution could be selected in a given setting and not others. To simplify the computations, we study the case of two competing inputs with standard deviations σ_i and excess kurtosis K_i , for $i = 1, 2$. Three types of solutions can then, in principle, exist:

$$(0, 0), \quad (w_1^* \neq 0, 0), \quad (w_1^* \neq 0, w_2^* \neq 0),$$

144 with the $(0, w_2^* \neq 0)$ being the analog of $(w_1^* \neq 0, 0)$. One can compute the eigenvalues $\lambda_{1,2}$ in each
145 case and evaluate the stability of the fixpoints. A sketch of the fixpoints and their stability is presented
146 in Fig. 3.

- The trivial fixpoint $(0, 0)$ is always unstable, with positive eigenvalues

$$\lambda_{1,2}(0, 0) = \epsilon_w \frac{x_0^2}{N^2} (\sigma_1^2, \sigma_2^2). \quad (14)$$

- For $(w_1^* \neq 0, 0)$ one finds the eigenvalues

$$\lambda_{1,2}(w_1^* \neq 0, 0) = \epsilon_w \frac{x_0^2}{N^2} \left(-2\sigma_1^2, \frac{\sigma_2^2 K_1}{K_1 + 3} \right). \quad (15)$$

147 The first eigenvalue λ_1 is hence always negative with the sign of the second eigenvalue λ_2
148 depending exclusively on K_1 . The fixpoint $(w_1^* \neq 0, 0)$ is hence stable / unstable for negative
149 / positive K_1 .

- The last term $3\Phi - x_0^2$ in (12) is identical for all synapses. Two non-zero synaptic weights $(w_1^* \neq 0, w_2^* \neq 0)$ can hence only exist for identical signs of the respective excess kurtosis, $K_1 K_2 \geq 0$. It is easy to show that $(w_1^* \neq 0, w_2^* \neq 0)$ is unstable/stable whenever both $K_{1,2}$ are negative/positive, in accordance with (15).

150 The solutions of the type $(w_1^* \neq 0, 0)$ are hence the only stable fixpoints when the excess kurtosis of the
151 corresponding direction, in this case K_1 , is negative.

156 2.1.3. Principal component analysis

157 The observation [32], that the update rules (6) perform a principal component analysis (PCA) can be
158 understood on two levels.

159 It is, firstly, evident from (10) that $\langle \dot{w}_j \rangle \sim \sigma_j^2$, and that synaptic weight tends hence to grow fast
160 whenever the corresponding presynaptic input has a large variance σ_j^2 .

161 Alternatively one can consider the respective phase-space contractions $\lambda_1^{(\alpha)} + \lambda_2^{(\alpha)}$, see Eq. (15), around
162 the two competing fixpoints $\mathbf{w}^{(\alpha=1)} = (w_1^* \neq 0, 0)$ and $\mathbf{w}^{(\alpha=2)} = (0, w_2^* \neq 0)$. Using the expression (15)
163 for the case $K_1 = K_2 < 0$, one finds that the phase space contracts faster around $\mathbf{w}^{(1)}$ when $\sigma_1^2 > \sigma_2^2$,
164 and vice versa.

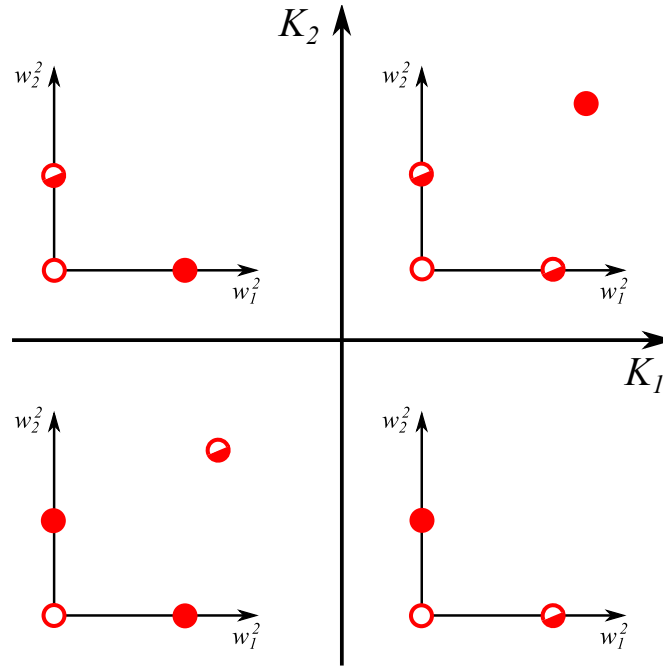


Figure 3. Sketch of the fixpoints of (10), which approximates (9), for two competing weights w_1 and w_2 as a function of the kurtosis K_1 and K_2 of the respective input directions. Open, full and half-full circles represent unstable fixpoints, stable fixpoints and saddles respectively. The axes are expressed in terms of w_i^2 , since the solutions are determined only up to a sign change.

165 **2.2. Alternative transfer functions**

The objective function (5) can be expressed generically [37] as

$$\mathcal{F}^{syn} = E \left[(N + A(x))^2 \right], \quad A(x) = \frac{xy''}{y'}, \quad (16)$$

166 where y' and y'' represent the first and second derivative of the transfer function $y(x) = \sigma(x - b)$ with
 167 respect to x . This expression, which appears as an intermediate step in the derivation of the objective
 168 function from the Fisher information, compare Section 2.3, reduces to (5) for the sigmoidal transfer
 169 function defined in (3).

The qualitative behavior of the learning rule remains unchanged when considering alternative functional forms for the transfer function $y(x)$, whenever they fulfill the basic requirement of being smooth monotonic functions with $\lim_{x \rightarrow \mp\infty} y(x) = 0/1$. For example, in [37], the authors showed that this is indeed the case for an arc-tangential transfer function. An interesting transfer function to consider in this context is the rescaled error function $\text{erf}(x - b)$,

$$y = \frac{1}{2} + \frac{1}{2} \text{erf} \left(\frac{x - b}{s\sqrt{2}} \right) = \frac{1}{2} + \frac{1}{\sqrt{\pi}} \int_{-\infty}^{(x-b)/(s\sqrt{2})} e^{-z^2} dz = \frac{1}{\sqrt{2\pi}s} \int_{-\infty}^{x-b} e^{-\frac{z^2}{2s^2}} dz, \quad (17)$$

as defined by the integral of the normal distribution of variance s . The constant s sets the slope of the transfer function and if one wants to have the same slope as for the original transfer function (3), one simply sets $s = 4/\sqrt{2\pi}$. The derivatives of (17) are:

$$y' = \frac{1}{\sqrt{2\pi}s} e^{-\frac{(x-b)^2}{2s^2}}, \quad y'' = -\frac{(x-b)}{s^2} y'. \quad (18)$$

Using (16) one obtains

$$\mathcal{F}^{syn} = E \left[\left(N - \frac{x(x-b)}{s^2} \right)^2 \right] \quad (19)$$

170 for the objective function and consequently

$$\begin{aligned} \dot{w}_j &= \epsilon_w (x - b/2) (Ns^2 - x(x-b)) (y_j - \bar{y}_j) \\ &= \epsilon_w (x - b/2) (x_0^2 - x(x-b)) (y_j - \bar{y}_j) \end{aligned} \quad (20)$$

for the synaptic plasticity rule, where we have replaced Ns^2 by x_0^2 , the squared roots for $b = 0$. Equation (20) reduces, interestingly and apart from an overall scaling factor, to the cubic approximation (20) for $b = 0$:

$$\dot{w}_j = -\epsilon_w x (x - x_0) (x + x_0) (y_j - \bar{y}_j) = \epsilon_w x (x_0^2 - x^2) (y_j - \bar{y}_j). \quad (21)$$

For non-vanishing b , we can rewrite (20), as:

$$\dot{w}_j = -\epsilon_w (x - b/2) (x - x^-) (x - x^+) (y_j - \bar{y}_j). \quad (22)$$

with

$$x^\pm = -\frac{b}{2} \pm \sqrt{b^2/4 + x_0^2} \approx -\frac{b}{2} \pm x_0, \quad (23)$$

171 where the last expression holds for small b . The whole learning rule (21) is therefore, for small bias b ,
172 simply shifted by a factor $b/2$.

Finally, in analogy to (7), we can again write (20) as a product of a Hebbian term (H) and a self-limiting term (G):

$$\dot{w}_j = \epsilon_w H(x) G(x) (y_j - \bar{y}_j), \quad H(x) = (x - b/2), \quad G(x) = (x_0^2 - x(x-b)). \quad (24)$$

173 In order to compare to Figure 2 (a), functions G and H , now as defined in (24), are plotted as a function
174 of the activity level y in Figure 4 (a).

We can now easily compute the average weight change for (20) in the same way we did for (8), obtaining

$$\langle \dot{w}_j \rangle = \epsilon_w w_j \sigma_j^2 \left[\left(x_0^2 - \frac{b^2}{2} \right) + \frac{3b}{2} w_j \sigma_j S_j - w_j^2 \sigma_j^2 K_j - 3\Phi \right], \quad (25)$$

where S_j is the skewness of input distribution y_j , as defined by

$$S_j = \frac{E[\gamma_j^3]}{\sigma_j^3}. \quad (26)$$

175 In expression (25) the interaction between intrinsic and synaptic plasticity becomes evident through b .
176 We note that for symmetric input distributions ($S_j = 0$) as the ones we have been treating, small values
177 of b produce only a shift in the effective x_0 (provided that b^2 is smaller than $x_0^2/2$).

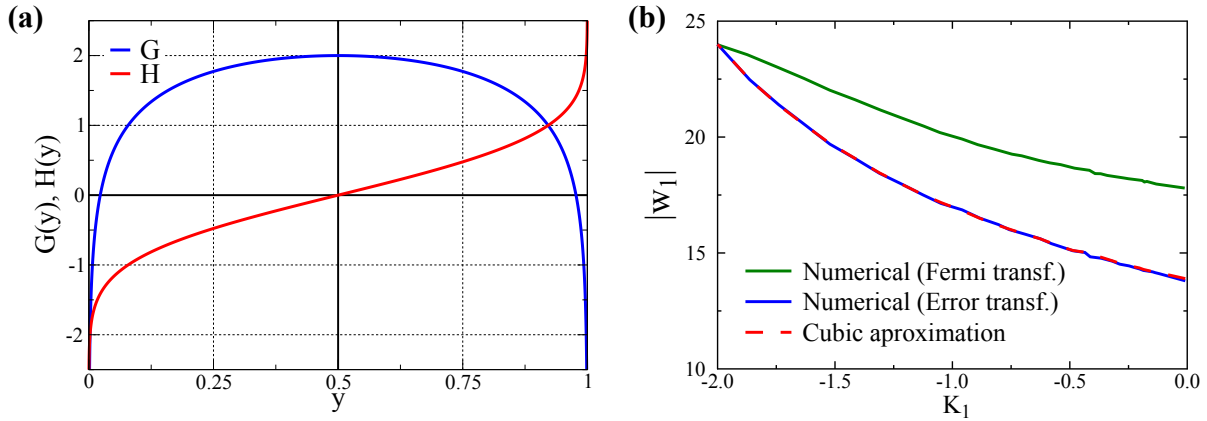


Figure 4. (a) Functions G and H , as in Figure 2 (a), now for (7), here expressed entirely in terms of the output activity $y \in [0, 1]$, for clarity. (b) Final absolute value of the weight w_1 after training, with both learning rules (6) and (20), together with the prediction (13) from the cubic approximation, as a function of the kurtosis K_1 for the direction of the principal component. For $b = 0$, $\sigma_1 = 0.1$, $\sigma_{i \neq 1} = \sigma_1/2$ and $N_w = 100$. One observes that the prediction is practically exact in the case of the error transfer function, remaining qualitatively similar for the case of the Fermi transfer function (6).

178 We note that the trivial solution $w_j = 0$ would become stable for negative $x_0^2 - b^2/2$. This has however
 179 not happened for the numerical simulations we performed, which resulted in values of $b \approx 1$, for target
 180 activity levels $\langle y \rangle$ as low as 0.1, while $x_0 = 2.4$ for $N = 2$ (which we used). Even sparser activity levels
 181 $\langle y \rangle \ll 1$ would require larger firing thresholds $b \gg 1$ and stable synaptic plasticity would be achieved by
 182 selecting then appropriately large N , corresponding to values of x_0 such that $x_0^2 - b^2/2$ remains positive.

183 For non-symmetric distributions the skewness of the input distribution, together with the sign of b will
 184 determine the sign w , which was before undetermined, since the learning rules are rotationally invariant.

185 2.3. The stationarity principle of statistical learning

186 In statistical learning one considers an agent trying to extract information from data input streams
 187 having stationary statistical properties. In a neural setting, this corresponds to a neuron adapting the
 188 afferent synaptic weights and learning is complete when the synaptic weights do not change any more.
 189 At this point, the probability distribution function $p(y)$ of the neural activity y also becomes stationary
 190 and its sensitivity with respect to changes in the afferent synaptic weight vector vanishes. This is the
 191 stationarity principle of statistical learning.

192 2.3.1. The Fisher information with respect to the synaptic flux

The Fisher information [36]

$$\mathcal{F}_\theta = \int p_\theta(y) \left(\frac{\partial}{\partial \theta} \ln(p_\theta(y)) \right)^2 dy \quad (27)$$

193 encodes the average sensitivity of a given probability distribution $p_\theta(y)$ with respect to a certain
 194 parameter θ , becoming minimal whenever θ does not influence the statistics of y . The Fisher information

195 is hence a suitable information theoretical functional for the implementation of the stationary principle
196 of statistical learning.

We drop the index θ in the following and consider in a first step a neuron with $N_w = 1$ afferent neurons. We define with

$$\mathcal{F}_{N_w=1}^{syn} = \int \left(w_1 \frac{\partial}{\partial w_1} \ln(p(y(y_1))) \right)^2 p(y_1) dy_1 \quad (28)$$

197 an objective function for synaptic plasticity, measuring the sensibility of the neural activity with respect
198 to the afferent synaptic weight w_1 . Here $y(y_1)$ is given by $\sigma(w_1(y_1 - \bar{y}_1) - b)$, as defined in Eq. (3).
199 There are two changes with respect to the bare Fisher information (27).

- 200 • The operator $w_1 \partial / \partial w_1$ corresponds to a dimensionless differential operator and hence to the
201 log-derivative. The whole objective function $\mathcal{F}_{N_w=1}^{syn}$ is hence dimensionless.
- 202 • The average sensitivity is computed as an average over the probability distribution $p(y_1)$ of the
203 presynaptic activity y_1 , since we are interested in minimizing the time average of the sensitivity
204 of the postsynaptic activity with respect to synaptic weight changes in the context of a stationary
205 presynaptic activity distribution $p(y_1)$.

For a distribution $p(y(y_1))$, for which y is a monotonic function of y_1 , we have

$$p(y(y_1)) dy = p(y_1) dy_1, \quad p(y(y_1)) = \frac{p(y_1)}{\partial y / \partial y_1}, \quad (29)$$

which allows us to rewrite (28) as:

$$\mathcal{F}_{N_w=1}^{syn} = \int \left(w_1 \frac{\partial}{\partial w_1} \ln \left(\frac{p(y_1)}{\partial y / \partial y_1} \right) \right)^2 p(y_1) dy_1. \quad (30)$$

Defining with $\mathbf{y} = (y_1, \dots, y_{N_w})$ the vector of afferent synaptic weights and with $p(\mathbf{y})$ the corresponding probability distribution function we may generalize (30) as

$$\mathcal{F}_{N_w}^{syn} = \int \left(\sum_{j=1}^{N_w} w_j \frac{\partial}{\partial w_j} \ln \left(\frac{p(y_j)}{\partial y / \partial y_j} \right) \right)^2 p(\mathbf{y}) d\mathbf{y}, \quad (31)$$

where we have replaced $p(y(\mathbf{y}))$ from (28) by $\frac{p(y_j)}{\partial y / \partial y_j}$, in what constitutes the independent synapse extension, and which represents the Fisher information with respect to the flux operator:

$$\frac{\partial}{\partial \theta} \rightarrow \sum_j w_j \frac{\partial}{\partial w_j} = \mathbf{w} \cdot \nabla_{\mathbf{w}}, \quad (32)$$

206 which is a dimensionless scalar. Some comments:

- 207 • Minimizing $\mathcal{F}_{N_w}^{syn}$, in accordance with the stationarity principle for statistical learning, leads
208 to a synaptic weight vector \mathbf{w} which is perpendicular to the gradient $\nabla_{\mathbf{w}}(\log(p))$, restricting
209 consequently the overall growth of the modulus of \mathbf{w} .

- In $\mathcal{F}_{N_w}^{syn}$ there is no direct cross talk between different synapses. Expression (32) is hence adequate for deriving Hebbian-type learning rules in which every synapse has access only to locally available information, together with the overall state of the postsynaptic neuron in terms of its firing activity y , or its membrane potential x . We call (32) the *local synapse extension* with respect to other formulations allowing for inter-synaptic cross talk.
- It is straightforward to show [37], that (31) reduces to (5), when using the relations (3), viz $\mathcal{F}_{N_w}^{syn} = \mathcal{F}^{syn}$ when we identify $N \rightarrow N_w$. We have, however, opted to retain N generically as a free parameter in (5), allowing to shift appropriately the roots of $G(x)$.

3. Results and Discussion

3.1. Quantitative comparison of the model and the cubic approximation

In the present section we test the prediction of equation (13) for the dependence of the weight size on the standard deviation σ_j and kurtosis K_j of the input distribution. Given that the input distribution has a finite width, the integrated input x cannot fall into the minima of the objective function for every point in the distribution, but rather the cloud of x points generated will tend to spread around these minima. The discrepancies in the rule from the cubic approximation in the vicinity and away from the minima are then expected to affect the final result of the learning procedure.

In order to test (13), we use as an input for the direction of the first principal component (FPC, which is chosen to be along y_1 , without loss of generality), the sum

$$\frac{1}{2} \left[N \left(x - \frac{1+2d}{2}, \sigma_s \right) + N \left(x - \frac{1-2d}{2}, \sigma_s \right) \right]$$

of two normal distributions $N(x, \sigma_s)$ with individual standard deviations σ_s , whose peaks are at a distance $\pm d$ from the center of the input range (0.5).

- σ_s is adjusted, changing d , such that the overall standard deviation σ_1 remains constant. In this way, one can select with d different kurtosis levels, while retaining a constant standard deviation. For $d = 0$, one gets a bound (since $y_1 \in [0, 1]$) normal distribution with $K_1 \approx 0$ (slightly negative since the distributions are bound). In this way, we can evaluate the size of w_1 after training for a varying $K_1 \in [-2, 0)$ for any given σ_1 .
- For the other $N_w - 1$ directions, we use bound normal distributions with standard deviations $\sigma_i = \sigma_1/2$ as in [32].

We tested training the neuron using both the original learning rule (6), derived for the Fermi transfer function, and the cubic rule (20) for the error function.

In Figure 4 (b), the value of w_1 after training is presented together with the prediction (13) from the cubic approximation, as a function of K_1 (the kurtosis in the y_1 direction), for a constant $b = 0$. In this case we have used $\sigma_1 = 0.1$, and $\sigma_{i \neq 1} = \sigma_1/2$.

The prediction of the cubic approximation is indeed practically exact for the error transfer function, as expected, since the input distributions are symmetric and, if one sets the axes parallel to the

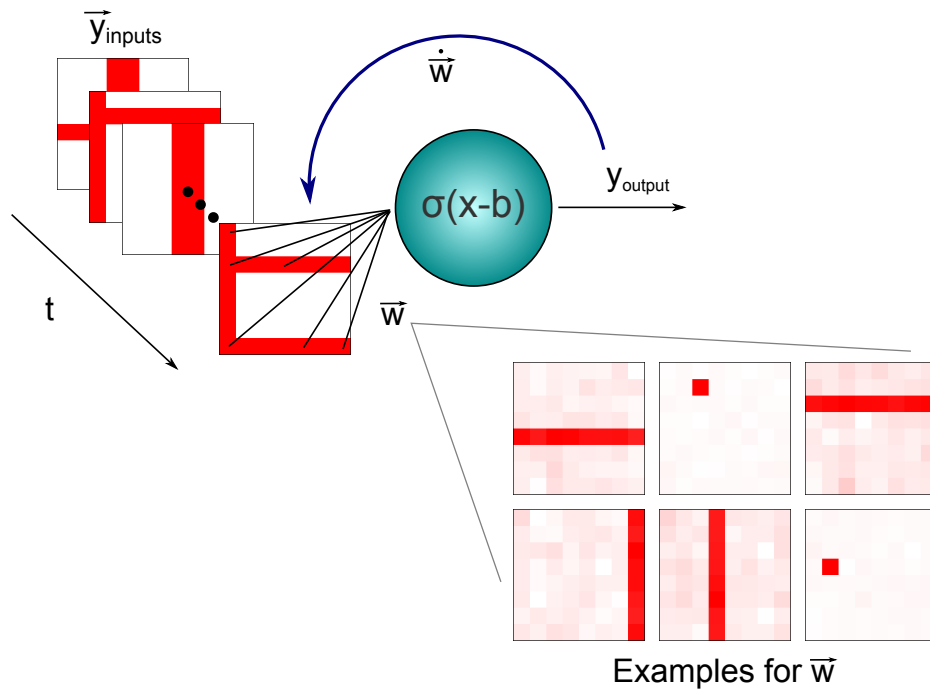


Figure 5. A single neuron, whose synaptic weights evolve according to (6) is presented with a set of input images consisting of the non-linear superposition of a random set of bars. We find that, on subsequent iterations, the neuron becomes selective to either single bars (the independent components of the input distribution), or to points.

principal components, as we did in this case, the correlation terms indeed vanish, therefore fulfilling the assumptions made during the averaging procedure. Apart from this procedure no other approximations were made in this case, since the rule for the error function is identical to the cubic approximation in the $b = 0$ case.

As shown in Figure 2 (b), even though the cubic approximation is able to reproduce the roots of the learning rule derived for the Fermi transfer function, the cubic approximation grows faster in the outer region, restricting the growth of the synaptic weight more than the original rule would. One therefore expects the cubic approximation to underestimate the final value of w_1 , as indeed is observed in Figure 4 (b). When $K = -2$, the input distribution becomes the sum of two deltas and the rule is able to assign each delta to a root. The prediction is of course once again exact in this case. Otherwise, while the quantitative result differs from one rule to another, the qualitative behavior remains unchanged.

3.2. Independent component analysis: an application to the nonlinear bars problem

Incoming signals may result from the sum of non-Gaussian independent sources and the process of extracting these sources is denoted independent component analysis (ICA) [48].

The non-Gaussianness of the independent sources may be characterized in principle by any quantity which is zero for the normal distribution and non-zero otherwise. Common measures include the kurtosis K and the skewness S , or the negentropy in general. For a comprehensive summary of the principles behind typical ICA procedures, as well as a description of independent components in terms of the cumulative moments of the probability distributions, see [49]. Our learning rule is functionally

261 dependent both on the kurtosis in general, as evident within the cubic approximation (10), and on the
262 skewness for non-zero bias b (25), and hence prone to perform an ICA. We note, however, that this
263 dependency does not result from maximizing a given measure of non-Gaussianness, as performed in
264 the past by several groups [50,51]. The resulting preference for non-Gaussianness is in the case of the
265 present work a by-product of the stationarity principle of statistical learning. In [27], the author shows
266 how, under certain conditions, a neural network evolving under a nonlinear principal component analysis
267 learning rule, is also capable of performing ICA.

268 The classical ICA is, strictly speaking, defined only for linear superpositions of sources. One can
269 generalize this concept to non-linear tasks and we test our multiplicative learning rule (6) using a
270 classical example for a non-linear ICA: the non-linear bars problem [35], in which a neuron, or a network
271 of neurons, is trained with a set of inputs, each representing an image consisting on the non-linear
272 superposition of horizontal and vertical bars.

273 In a grid of N_w inputs where $N_w = L \times L$, each horizontal and vertical bar has a constant probability
274 of being present of $p = 1/L$. Each input or pixel can take only two values: a low intensity and a high
275 intensity value. Each bar then corresponds to a whole row or a column of high intensity pixels, where
276 at the intersection of two bars the pixel has the same value (high) as in the rest of the bar, making the
277 problem non-linear.

278 The here examined synaptic plasticity rules (6) are able, as illustrated in Fig. 5, to discriminate
279 individual bars, the independent components of the input patterns, or points [37]. One might argue
280 that, given that in the original training set of [35] single bars will occur in the training set with a finite
281 probability, the neuron is simply selecting this particular input. To rule out this possibility, we trained
282 the neuron also with sets having at least of one horizontal and one vertical bar each time, so that no bar
283 is ever presented in isolation. No appreciable change in the performance was observed.

284 4. Conclusions

285 We presented guiding principles for deriving plasticity rules for the synaptic weight of interconnected
286 neurons which are motivated by considering an information theoretical measure, namely the Fisher
287 information, for the implementation of the stationarity principle of statistical learning. We showed
288 how, in the case of ellipsoidal input distributions, the resulting plasticity rules find, as usual for
289 Hebbian learning, the dominant principal component in the data input stream, when present, being
290 selective otherwise for non-Gaussianness in terms of the excess kurtosis and the skewness of the input
291 activities. The plasticity rules are hence also prone to perform an independent component analysis, as
292 we demonstrated by considering the non-linear bars problem.

293 The here examined adaption rules are self-limiting. This is a natural consequence of implementing
294 the stationarity principle of statistical learning, which states that the statistics of the postsynaptic neural
295 firing will become stationary whenever learning is complete and when the statistics of the input activity is
296 itself stationary. The self-limitation is achieved through a multiplicative factor to the usual Hebbian-type
297 plasticity function, in contrast to other approaches, where runaway growth of the synaptic weights is
298 avoided by performing either an overall renormalization of the synaptic weights, or by adding an explicit
299 weight decay term.

300 In previous work [32], a numerical comparison between the learning rules here proposed and the
301 traditionally employed Oja's rule was performed, showing differences in the sensitivity to higher
302 moments of the input distribution (Oja's rule is tailored to be sensitive to the second moment of the
303 input distribution only), as well as a stark contrast in terms of transient dynamics. While Oja's rule
304 predicts the neuron to learn and unlearn the direction of the PCA within the same timescale when a
305 new interesting direction is presented, a fading memory effect is observed when the present rules are
306 employed. We believe then, that depending on the application at hand, and the level of noise in the
307 environment, one or the other might prove more suitable.

308 The objective function (16) and the learning rules discussed here depend on the specific form
309 $y(x)$ of the transfer function, becoming a cubic polynomial in the membrane potential x when the
310 transfer function is a rescaled error function. This cubic plasticity rule (20) is, at the same time, an
311 excellent approximation for the update rule (6) valid for sigmoidal transfer functions, allowing to derive
312 analytically the sensibility of our learning rules to the excess kurtosis, as discussed in Sect. 2.1.2. The
313 polynomial update rules allows also to study, given its polynomial character, the stability of the learning
314 dynamics quite general in terms of the moments of the input distribution.

315 Finally, we have shown here and in [37], how neurons operating under several transfer functions, and
316 with learning rules which are only qualitatively equivalent to the cubic function, are able to perform
317 identical computational tasks. We have also tested whether a neuron defined by one particular transfer
318 function can be trained using the learning rule derived for another choice of sigmoidal function, finding
319 no major changes to the results. The procedure is then very robust to quantitative deviations from the
320 derived rules, as long as the plasticity rule remains qualitatively similar to a cubic polynomial in the
321 membrane potential, an important requirement for biological plausibility.

322 Acknowledgments

323 We thank Bulcsú Sándor for his valuable input on gradient systems. The support of the German
324 Science Foundation (DFG) and the German Academic Exchange Service (DAAD) are acknowledged.

325 Conflicts of Interest

326 The authors declare no conflict of interest.

327 References

- 328 1. Attwell, D.; Laughlin, S.B. An energy budget for signaling in the grey matter of the brain.
329 *Journal of Cerebral Blood Flow & Metabolism* **2001**, *21*, 1133–1145.
- 330 2. Mink, J.W.; Blumenshine, R.J.; Adams, D.B. Ratio of central nervous system to body
331 metabolism in vertebrates: its constancy and functional basis. *American Journal of*
332 *Physiology-Regulatory, Integrative and Comparative Physiology* **1981**, *241*, R203–R212.
- 333 3. Niven, J.E.; Laughlin, S.B. Energy limitation as a selective pressure on the evolution of sensory
334 systems. *Journal of Experimental Biology* **2008**, *211*, 1792–1804.
- 335 4. Bullmore, E.; Sporns, O. The economy of brain network organization. *Nature Reviews*
336 *Neuroscience* **2012**, *13*, 336–349.

- 337 5. Lee, H.; Battle, A.; Raina, R.; Ng, A.Y. Efficient sparse coding algorithms. *Advances in neural*
338 *information processing systems*, 2006, pp. 801–808.
- 339 6. Stemmler, M.; Koch, C. How voltage-dependent conductances can adapt to maximize the
340 information encoded by neuronal firing rate. *Nature neuroscience* **1999**, *2*, 521–527.
- 341 7. Gros, C. Generating functionals for guided self-organization. In *Guided Self-Organization:*
342 *Inception*; Prokopenko, M., Ed.; Springer, 2014; pp. 53–66.
- 343 8. MacKay, D. Information-based objective functions for active data selection. *Neural computation*
344 **1992**, *4*, 590–604.
- 345 9. Marler, R.T.; Arora, J.S. Survey of multi-objective optimization methods for engineering.
346 *Structural and multidisciplinary optimization* **2004**, *26*, 369–395.
- 347 10. Intrator, N.; Cooper, L.N. Objective function formulation of the BCM theory of visual cortical
348 plasticity: Statistical connections, stability conditions. *Neural Networks* **1992**, *5*, 3–17.
- 349 11. Kay, J.W.; Phillips, W. Coherent infomax as a computational goal for neural systems. *Bulletin of*
350 *mathematical biology* **2011**, *73*, 344–372.
- 351 12. Polani, D. Information: currency of life? *HFSP journal* **2009**, *3*, 307–316.
- 352 13. Zahedi, K.; Ay, N.; Der, R. Higher coordination with less control - a result of information
353 maximization in the sensorimotor loop. *Adaptive Behavior* **2010**, *18*, 338–355.
- 354 14. Polani, D.; Prokopenko, M.; Yaeger, L.S. Information and self-organization of behavior.
355 *Advances in Complex Systems* **2013**, *16*.
- 356 15. Prokopenko, M.; Gershenson, C. Entropy Methods in Guided Self-Organisation. *Entropy* **2014**,
357 *16*, 5232–5241.
- 358 16. Der, R.; Martius, G.; *The Playful Machine: Theoretical Foundation and Practical Realization of*
359 *Self-Organizing Robots*; Vol. 15, Springer-Verlag Berlin Heidelberg, 2012
- 360 17. Markovic, D.; Gros, C. Self-organized chaos through polyhomeostatic optimization. *Physical*
361 *Review Letters* **2010**, *105*, 068702.
- 362 18. Marković, D.; Gros, C. Intrinsic adaptation in autonomous recurrent neural networks. *Neural*
363 *Computation* **2012**, *24*, 523–540.
- 364 19. Triesch, J. Synergies between intrinsic and synaptic plasticity mechanisms. *Neural Computation*
365 **2007**, *19*, 885–909.
- 366 20. Linsker, R. Local synaptic learning rules suffice to maximize mutual information in a linear
367 network. *Neural Computation* **1992**, *4*, 691–702.
- 368 21. Chechik, G. Spike-timing-dependent plasticity and relevant mutual information maximization.
369 *Neural computation* **2003**, *15*, 1481–1510.
- 370 22. Toyozumi, T.; Pfister, J.P.; Aihara, K.; Gerstner, W. Generalized Bienenstock–Cooper–Munro
371 rule for spiking neurons that maximizes information transmission. *Proceedings of the National*
372 *Academy of Sciences of the United States of America* **2005**, *102*, 5239–5244.
- 373 23. Friston, K. The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*
374 **2010**, *11*, 127–138.
- 375 24. Mozzachiodi, R.; Byrne, J.H. More than synaptic plasticity: role of nonsynaptic plasticity in
376 learning and memory. *Trends in neurosciences* **2010**, *33*, 17–26.

- 377 25. Strogatz, S.H. *Nonlinear dynamics and chaos: with applications to physics, biology and*
 378 *chemistry*; Perseus publishing, 2001.
- 379 26. Hebb, D.O. *The organization of behavior: A neuropsychological theory*; Psychology Press, 2002.
- 380 27. Oja, E. The nonlinear PCA learning rule in independent component analysis. *Neurocomputing*
 381 **1997**, *17*, 25–45.
- 382 28. Bi, G.q.; Poo, M.m. Synaptic modifications in cultured hippocampal neurons: dependence on
 383 spike timing, synaptic strength, and postsynaptic cell type. *The Journal of neuroscience* **1998**,
 384 *18*, 10464–10472.
- 385 29. Froemke, R.C.; Dan, Y. Spike-timing-dependent synaptic modification induced by natural spike
 386 trains. *Nature* **2002**, *416*, 433–438.
- 387 30. Izhikevich, E.M.; Desai, N.S. Relating stdp to bcm. *Neural computation* **2003**, *15*, 1511–1523.
- 388 31. Echeveste, R.; Gros, C. Two-trace model for spike-timing-dependent synaptic plasticity. *Neural*
 389 *computation* **2015**, *27*, 672–698.
- 390 32. Echeveste, R.; Gros, C. Generating functionals for computational intelligence: The Fisher
 391 information as an objective function for self-limiting Hebbian learning rules. *Frontiers in*
 392 *Robotics and AI* **2014**, *1*.
- 393 33. Bell, A.J.; Sejnowski, T.J. An information-maximization approach to blind separation and blind
 394 deconvolution. *Neural computation* **1995**, *7*, 1129–1159.
- 395 34. Martius, G.; Der, R.; Ay, N. Information driven self-organization of complex robotic behaviors.
 396 *PloS one* **2013**, *8*, e63400.
- 397 35. Földiak, P. Forming sparse representations by local anti-Hebbian learning. *Biological cybernetics*
 398 **1990**, *64*, 165–170.
- 399 36. Brunel, N.; Nadal, J.P. Mutual information, Fisher information, and population coding. *Neural*
 400 *Computation* **1998**, *10*, 1731–1757.
- 401 37. Echeveste, R.; Gros, C. An objective function for self-limiting neural plasticity rules.
 402 *Proceedings of the European Symposium on Artificial Neural Networks, Computational*
 403 *Intelligence and Machine Learning (ESANN) (In press: to appear in April 2015)*.
- 404 38. Hyvärinen, A.; Karhunen, J.; Oja, E. *Independent component analysis*; Vol. 46, John Wiley &
 405 Sons, 2004.
- 406 39. Bell, A.J.; Sejnowski, T.J. The “independent components” of natural scenes are edge filters.
 407 *Vision research* **1997**, *37*, 3327–3338.
- 408 40. Paradiso, M. A theory for the use of visual orientation information which exploits the columnar
 409 structure of striate cortex. *Biological cybernetics* **1988**, *58*, 35–49.
- 410 41. Seung, H.; Sompolinsky, H. Simple models for reading neuronal population codes. *Proceedings*
 411 *of the National Academy of Sciences* **1993**, *90*, 10749–10753.
- 412 42. Gutnisky, D.A.; Dragoi, V. Adaptive coding of visual information in neural populations. *Nature*
 413 **2008**, *452*, 220–224.
- 414 43. Bethge, M.; Rotermund, D.; Pawelzik, K. Optimal neural rate coding leads to bimodal firing rate
 415 distributions. *Network: Computation in Neural Systems* **2003**, *14*, 303–319.
- 416 44. Lansky, P.; Greenwood, P.E. Optimal signal in sensory neurons under an extended rate coding
 417 concept. *BioSystems* **2007**, *89*, 10–15.

- 418 45. Ecker, A.S.; Berens, P.; Tolias, A.S.; Bethge, M. The effect of noise correlations in populations
419 of diversely tuned neurons. *The Journal of Neuroscience* **2011**, *31*, 14272–14283.
- 420 46. Reginatto, M. Derivation of the equations of nonrelativistic quantum mechanics using the
421 principle of minimum Fisher information. *Physical Review A* **1998**, *58*, 1775–1778.
- 422 47. DeCarlo, L.T. On the meaning and use of kurtosis. *Psychological methods* **1997**, *2*, 292.
- 423 48. Comon, P. Independent component analysis, a new concept? *Signal processing* **1994**,
424 *36*, 287–314.
- 425 49. Hyvärinen, A.; Oja, E. Independent component analysis: algorithms and applications. *Neural*
426 *networks* **2000**, *13*, 411–430.
- 427 50. Girolami, M.; Fyfe, C. Negentropy and kurtosis as projection pursuit indices provide generalised
428 ICA algorithms. A. C, Back A (eds.), NIPS-96 Blind Signal Separation Workshop, 1996, Vol. 8.
- 429 51. Li, H.; Adali, T. A class of complex ICA algorithms based on the kurtosis cost function. *Neural*
430 *Networks, IEEE Transactions on* **2008**, *19*, 408–420.

431 © June 5, 2015 by the authors; submitted to *Entropy* for possible open access
432 publication under the terms and conditions of the Creative Commons Attribution license
433 <http://creativecommons.org/licenses/by/4.0/>.