

Intrinsic Stabilization of Output Rates by Spike-Based Hebbian Learning

Richard Kempter

kempter@phy.ucsf.edu

Keck Center for Integrative Neuroscience, University of California at San Francisco, San Francisco, CA 94143-0732, U.S.A.

Wulfram Gerstner

Wulfram.Gerstner@epfl.ch

Swiss Federal Institute of Technology Lausanne, Laboratory of Computational Neuroscience, DI-LCN, CH-1015 Lausanne EPFL, Switzerland

J. Leo van Hemmen

Leo.van.Hemmen@ph.tum.de

Physik Department, Technische Universität München, D-85747 Garching bei München, Germany

We study analytically a model of long-term synaptic plasticity where synaptic changes are triggered by presynaptic spikes, postsynaptic spikes, and the time differences between presynaptic and postsynaptic spikes. The changes due to correlated input and output spikes are quantified by means of a learning window. We show that plasticity can lead to an intrinsic stabilization of the mean firing rate of the postsynaptic neuron. Subtractive normalization of the synaptic weights (summed over all presynaptic inputs converging on a postsynaptic neuron) follows if, in addition, the mean input rates and the mean input correlations are identical at all synapses. If the integral over the learning window is positive, firing-rate stabilization requires a non-Hebbian component, whereas such a component is not needed if the integral of the learning window is negative. A negative integral corresponds to anti-Hebbian learning in a model with slowly varying firing rates. For spike-based learning, a strict distinction between Hebbian and anti-Hebbian rules is questionable since learning is driven by correlations on the timescale of the learning window. The correlations between presynaptic and postsynaptic firing are evaluated for a piecewise-linear Poisson model and for a noisy spiking neuron model with refractoriness. While a negative integral over the learning window leads to intrinsic rate stabilization, the positive part of the learning window picks up spatial and temporal correlations in the input.

1 Introduction

“Hebbian” learning (Hebb, 1949), that is, synaptic plasticity driven by correlations between pre- and postsynaptic activity, is thought to be an important mechanism for the tuning of neuronal connections during development and thereafter, in particular, for the development of receptive fields and computational maps (see, e.g. von der Malsburg, 1973; Sejnowski, 1977; Kohonen, 1984; Linsker, 1986; Sejnowski & Tesauro, 1989; MacKay & Miller, 1990; Wimbauer, Gerstner, & van Hemmen, 1994; Shouval & Perrone, 1995; Miller, 1996a; for a review, see Wiskott & Sejnowski, 1998). It is well known that simple Hebbian rules may lead to diverging synaptic weights so that weight normalization turns out to be an important topic (Kohonen, 1984; Oja, 1982; Miller & MacKay, 1994; Miller, 1996b). In practice, normalization of the weights w_i is imposed by either an explicit rescaling of all weights after a learning step or a constraint on the summed weights (e.g., $\sum_i w_i = c$, or $\sum_i w_i^2 = c$ for some constant c), or an explicit decay term proportional to the weight itself.

In recent simulation studies of spike-based learning (Gerstner, Kempter, van Hemmen, & Wagner, 1996; Kempter, Gerstner, & van Hemmen, 1999a; Song, Miller, & Abbott, 2000; Xie & Seung, 2000; Kempter, Leibold, Wagner, & van Hemmen, 2001), however, intrinsic normalization properties of synaptic plasticity have been found. Neither an explicit normalization step nor a constraint on the summed weights was needed. So we face the question: How can we understand these findings? Some preliminary arguments as to why intrinsic normalization occurs in spike-based learning rules have been given (Gerstner et al., 1998; Kempter et al., 1999a; Song et al., 2000). In this article, we study normalization properties in more detail. In particular, we show that in spike-based plasticity with realistic learning windows, the procedure of continuously strengthening and weakening the synapses can automatically lead to a normalization of the total input strength to the postsynaptic neuron in a competitive self-organized process and, hence, to a stabilization of the output firing rate.

In order to phrase the problem of normalization in a general framework, let us start with a rate-based learning rule of the form

$$\begin{aligned} \tau_w \frac{d}{dt} w_i(t) = & a_0 + a_1^{\text{in}} \lambda_i^{\text{in}}(t) + a_1^{\text{out}} \lambda^{\text{out}}(t) + a_2^{\text{corr}} \lambda_i^{\text{in}}(t) \lambda^{\text{out}}(t) \\ & + a_2^{\text{in}} [\lambda_i^{\text{in}}(t)]^2 + a_2^{\text{out}} [\lambda^{\text{out}}(t)]^2, \end{aligned} \quad (1.1)$$

which can be seen as an expansion of local adaptation rules up to second order in the input rates λ_i^{in} ($1 \leq i \leq N$) and the output rate λ^{out} (see Bienenstock, Cooper, & Munro, 1982; Linsker, 1986; Sejnowski & Tesauro, 1989). (A definition of *rate* will be given below.) The timescale of learning is set by τ_w . The coefficients a_0 , a_1^{in} , a_1^{out} , a_2^{corr} , a_2^{in} , and a_2^{out} may, and in general will, depend on the weights w_i . For example, in Oja’s rule (Oja, 1982), we have

$a_2^{\text{corr}} = 1$ and $a_2^{\text{out}}(w_i) = -w_i$, while all other coefficients vanish; for linear neurons, the weight vector converges to a unit vector. Here we study learning rules of the form of equation 1.1, where the coefficients do *not* depend on the weights; we do, however, introduce upper and lower bounds for the weights (e.g., $dw_i/dt = 0$ if $w_i \geq w^{\text{max}}$ or $w_i \leq 0$ for an excitatory synapse) so as to exclude runaway of individual weight values.

As a first issue, we argue that the focus on normalization of the weights, be it in the linear form $\sum_i w_i$ or in the quadratic form $\sum_i w_i^2$, is too narrow. A broad class of learning rules will naturally stabilize the output rate rather than the weights. As an example, let us consider a learning rule of the form

$$\tau_w \frac{d}{dt} w_i(t) = -[\lambda^{\text{out}}(t) - \lambda^c] \lambda_i^{\text{in}} \quad (1.2)$$

with constant rates λ_i^{in} , ($1 \leq i \leq N$). Equation 1.2 is a special case of equation 1.1. For excitatory synapses, λ^{out} increases with w_i ; more precisely, $\lambda^{\text{out}}(w_1, \dots, w_N)$ is an increasing function of each of the w_i . Hence learning stops if λ^{out} approaches $\lambda^c > 0$, and the fixed point $\lambda^{\text{out}} = \lambda^c$ is stable. In the special case that all input lines i have an equal rate $\lambda_i^{\text{in}} = \lambda^{\text{in}}$ independent of i , stabilization of the mean output rate implies the normalization of the sum $\sum_i w_i$. Section 3 will make this argument more precise.

A disadvantage of the learning rule in equation 1.2 is that it has a correlation coefficient $a_2^{\text{corr}} < 0$. Thus, in a pure rate description, the learning rule would be classified as anti-Hebbian rather than Hebbian. As a second issue, we show in section 4 that for spike-based learning rules, the strict distinction between Hebbian and anti-Hebbian learning rules becomes questionable. A learning rule with a realistic time window (Levy & Stewart, 1983; Markram, Lübke, Frotscher, & Sakmann, 1997; Zhang, Tao, Holt, Harris, & Poo, 1998; Debanne, Gähwiler, & Thompson, 1998; Bi & Poo, 1998; Feldman, 2000) can be anti-Hebbian in a time-averaged sense and still pick up positive “Hebbian” correlations between input and output spikes (Gerstner, Kempter, et al., 1996; Kempter et al., 1999a; Kempter, Gerstner, & van Hemmen, 1999b; Kempter, Gerstner, van Hemmen, & Wagner, 1996). These correlations between input and output spikes that enter the learning dynamics will be calculated in this article for a linear Poisson neuron and for a noisy spiking neuron with refractoriness.

Spike-based rules open the possibility of a direct comparison of model parameters with experiments (Senn, Tsodyks, & Markram, 1997; Senn, Markram, & Tsodyks, 2001). A theory of spike-based learning rules has been developed in Gerstner, Ritz, and van Hemmen (1993), Häfliger, Mahowald, and Watts (1997), Ruf and Schmitt (1997), Senn et al. (1997, 2001), Eurich, Pawelzik, Ernst, Cowan, and Milton (1999), Kempter et al. (1999a), Roberts (1999), Kistler and van Hemmen (2000), and Xie and Seung (2000); for a review, see, van Hemmen (2000). Spike-based learning rules are closely related to rules for sequence learning (Herz, Sulzer, Kühn, & van Hemmen,

1988, 1989; van Hemmen et al., 1990; Gerstner et al., 1993; Abbott & Blum, 1996; Gerstner & Abbott, 1997), where the idea of asymmetric learning windows is exploited. In section 4 a theory of spike-based learning is outlined and applied to the problem of weight normalization by rate stabilization.

2 Input Scenario

We consider a single neuron that receives input from N synapses with weights w_i where $1 \leq i \leq N$ is the index of the synapse. At each synapse i , input spikes arrive stochastically. In the spike-based description of section 4, we will model the input spike train at synapse i as an inhomogeneous Poisson process with rate $\lambda_i^{\text{in}}(t)$. In the rate description of section 3, the input spike train is replaced by the continuous-rate variable $\lambda_i^{\text{in}}(t)$. Throughout the article, we focus on two different input scenarios—the static-pattern scenario and the translation-invariant scenario. In both scenarios, input rates are defined as statistical ensembles that we now describe.

2.1 Static-Pattern Scenario. The static-pattern scenario is the standard framework for the theory of unsupervised Hebbian learning (Hertz, Krogh, & Palmer, 1991). The input ensemble contains p patterns defined as vectors $\mathbf{x}^\mu \in \mathbb{R}^N$ where $\mathbf{x}^\mu = (x_1^\mu, x_2^\mu, \dots, x_N^\mu)^T$ with pattern label $1 \leq \mu \leq p$. Time is discretized in intervals Δ_t . In each time interval Δ_t , a pattern is chosen randomly from the ensemble of patterns and applied at the input; during application of pattern μ , the input rate at synapse i is $\lambda_i^{\text{in}} = x_i^\mu$. Temporal averaging over a time T yields the average input rate at each synapse:

$$\overline{\lambda_i^{\text{in}}(t)} := \frac{1}{T} \int_{t-T}^t dt' \lambda_i^{\text{in}}(t'). \quad (2.1)$$

For long intervals $T \gg p \Delta_t$, the average input rate is constant, $\overline{\lambda_i^{\text{in}}(t)} = \overline{\lambda_i^{\text{in}}}$, and we find

$$\overline{\lambda_i^{\text{in}}} = \frac{1}{p} \sum_{\mu=1}^p x_i^\mu. \quad (2.2)$$

An overbar will always denote temporal averaging over long intervals T . The normalized correlation coefficient of the input ensemble is

$$C_{ij}^{\text{static}} = \frac{1}{p} \sum_{\mu=1}^p \frac{(x_i^\mu - \overline{\lambda_i^{\text{in}}})(x_j^\mu - \overline{\lambda_j^{\text{in}}})}{\overline{\lambda_i^{\text{in}}} \cdot \overline{\lambda_j^{\text{in}}}} = C_{ji}^{\text{static}}. \quad (2.3)$$

An explicit example of the static-pattern scenario is given in the appendix. Input correlations play a major role in the theory of unsupervised Hebbian learning (Hertz et al., 1991).

2.2 Translation-Invariant Scenario. In developmental learning, it is often natural to assume a translation-invariant scenario (Miller, 1996a). In this scenario, the mean input rates are the same at all synapses and (again) constant in time:

$$\overline{\lambda_i^{\text{in}}(t)} = \overline{\lambda^{\text{in}}} \quad \text{for all } i. \quad (2.4)$$

At each synapse i , the actual rate $\lambda_i^{\text{in}}(t)$ is time dependent and fluctuates on a timescale Δ_{corr} around its mean. In other words, the correlation function, defined as

$$C_{ij}(s) := \frac{\overline{\Delta\lambda_i^{\text{in}}(t) \Delta\lambda_j^{\text{in}}(t-s)}}{\overline{\lambda_i^{\text{in}}} \cdot \overline{\lambda_j^{\text{in}}}} = C_{ji}(-s), \quad (2.5)$$

is supposed to vanish for $|s| \gg \Delta_{\text{corr}}$. Here $\Delta\lambda_i^{\text{in}}(t) := \lambda_i^{\text{in}}(t) - \overline{\lambda_i^{\text{in}}}$ is the deviation from the mean.

The essential hypothesis is that the correlation function is translation invariant, that is, $C_{ij} = C_{|i-j|}$. Let us suppose, for example, that the ensemble of stimuli consists of objects or patterns of some typical size spanning l input pixels so that the correlation function has a spatial width of order l . If the objects are moved randomly with equal probability in an arbitrary direction across the stimulus array, then the temporal part of the correlation function will be symmetric in that $C_{|i-j|}(s) = C_{|i-j|}(-s)$.

Using the static-pattern scenario or the translation-invariant scenario, we are able to simplify the analysis of Hebbian learning rules considerably (rate-based learning in section 3 and spike-based learning in section 4).

3 Plasticity in Rate Description

In this section, we start with some preliminary considerations formulated on the level of rates and then analyze the rate description proper. Our aim is to show that the learning dynamics defined in equation 1.1 can yield a stable fixed point of the output rate. In order to keep the discussion as simple as possible, we focus mainly on a linearized rate neuron model and use the input scenarios of section 2. The learning rule is that of equation 1.1 with $a_2^{\text{in}} = a_2^{\text{out}} = 0$. All synapses are assumed to be of the same type so that the coefficients a_1^{in} , a_1^{out} , and a_2^{corr} do not depend on the synapse index i .

As is usual in the theory of Hebbian learning, we assume that the slow timescale τ_w of learning and the fast timescale of fluctuations in the input and neuronal output can be separated by the averaging time T . For instance, we assume $\tau_w \gg T \gg p \Delta_i$ in the static-pattern scenario or $\tau_w \gg T \gg \Delta_{\text{corr}}$ in the translation-invariant scenario. A large τ_w implies that weight changes within a time interval of duration T are small compared to typical values of weights. The right-hand side of equation 1.1 is then “self-averaging”

(Kempster et al., 1999a) with respect to both randomness and time; that is, synaptic changes are driven by statistically and temporally averaged quantities,

$$\tau_w \frac{d}{dt} w_i(t) = a_0 + a_1^{\text{in}} \overline{\lambda_i^{\text{in}}} + a_1^{\text{out}} \overline{\lambda^{\text{out}}(t)} + a_2^{\text{corr}} \overline{\lambda_i^{\text{in}} \cdot \lambda^{\text{out}}(t)} + a_2^{\text{corr}} \overline{\Delta \lambda_i^{\text{in}}(t) \Delta \lambda^{\text{out}}(t)}. \quad (3.1)$$

“Hebbian” learning of a synapse should be driven by the correlations of pre- and postsynaptic neuron, that is, the last term on the right-hand side of equation 3.1 (cf. Sejnowski, 1977; Sejnowski & Tesauro, 1989; Hertz et al., 1991). This term is of second order in $\Delta \lambda^{\text{in}}$, and it might therefore seem that it is small compared to the term proportional to $\overline{\lambda_i^{\text{in}} \cdot \lambda^{\text{out}}}$ —but is it really? Let us speculate for the moment that in an ideal neuron, the mean output rate $\overline{\lambda^{\text{out}}}$ is attracted toward an operating point,

$$\lambda_{\text{FP}}^{\text{out}} = - \frac{a_0 + a_1^{\text{in}} \overline{\lambda^{\text{in}}}}{a_1^{\text{out}} + a_2^{\text{corr}} \overline{\lambda^{\text{in}}}}, \quad (3.2)$$

where $\overline{\lambda^{\text{in}}}$ is the typical mean input rate that is identical at all synapses. Then the dominant term on the right-hand side of equation 3.1 would indeed be the covariance term $\propto \overline{\Delta \lambda_i^{\text{in}} \Delta \lambda^{\text{out}}}$ because all other terms on the right-hand side of equation 3.1 cancel each other. The idea of a fixed point of the rate $\overline{\lambda^{\text{out}}} = \lambda_{\text{FP}}^{\text{out}}$ will be made more precise in the following two sections.

3.1 Piecewise-Linear Neuron Model. As a first example, we study the piecewise-linear neuron model with rate function

$$\lambda^{\text{out}}(t) = \left[\lambda_0 + \gamma_0 \frac{1}{N} \sum_{i=1}^N w_i(t) \lambda_i^{\text{in}}(t) \right]_+. \quad (3.3)$$

Here, λ_0 is a constant, and $\gamma_0 > 0$ is the slope of the activation function. The normalization by $1/N$ ensures that the mean output rate stays the same if the number N of synaptic inputs (with identical input rate λ^{in}) is increased. In case $\lambda_0 < 0$, we must ensure explicitly that λ^{out} is nonnegative. To do so, we have introduced the notation $[\cdot]_+$ with $[x]_+ = x$ for $x > 0$ and $[x]_+ = 0$ for $x \leq 0$. In the following, we will always assume that the argument inside the square brackets is positive, drop the brackets, and treat the model as strictly linear. Only at the end of the calculation do we check for consistency—for $\lambda^{\text{out}} \geq 0$. At the end of learning, we also check for lower and upper bounds, $0 \leq w_i \leq w^{\text{max}}$ for all i . (For a discussion of the influence of lower and upper bounds for individual weights on the stabilization of the output rate, see section 5.3.)

Since learning is assumed to be slow ($\tau_w \gg \Delta_{\text{corr}}$), the weights w_i do not change on the fast timescale of input fluctuations. Hence the correlation term $\overline{\Delta \lambda_i^{\text{in}} \Delta \lambda^{\text{out}}}$ in equation 3.1 can be evaluated for constant w_i . Using equation 3.3, the definition $\Delta \lambda^{\text{out}}(t) := \lambda^{\text{out}}(t) - \overline{\lambda^{\text{out}}(t)}$, and the definition of the input correlations C_{ij} as given by equation 2.5, we obtain

$$\overline{\Delta \lambda_i^{\text{in}}(t) \Delta \lambda^{\text{out}}(t)} = \overline{\lambda_i^{\text{in}}} \gamma_0 \frac{1}{N} \sum_{j=1}^N w_j(t) \overline{\lambda_j^{\text{in}}} C_{ij}(0). \quad (3.4)$$

We now want to derive the conditions under which the mean output rate $\overline{\lambda^{\text{out}}}$ approaches a fixed point $\lambda_{\text{FP}}^{\text{out}}$. Let us define, whatever j , the ‘‘average correlation,’’

$$C := \frac{\sum_{i=1}^N (\overline{\lambda_i^{\text{in}}})^2 C_{ij}(0)}{\sum_{i=1}^N (\overline{\lambda_i^{\text{in}}})^2}; \quad (3.5)$$

that is, we require that the average correlation be independent of the index j . This condition may look artificial, but it is, in fact, a rather natural assumption. In particular, for the translation-invariant scenario, equation 3.5 always holds. We use equations 3.3 through 3.5 in equation 3.1, multiply by $\gamma_0 N^{-1} \overline{\lambda_i^{\text{in}}}$, and sum over i . Rearranging the result, we obtain

$$\tau \frac{d}{dt} \overline{\lambda^{\text{out}}(t)} = \lambda_{\text{FP}}^{\text{out}} - \overline{\lambda^{\text{out}}(t)} \quad (3.6)$$

with fixed point

$$\lambda_{\text{FP}}^{\text{out}} := \frac{\tau}{\tau_w} \frac{\gamma_0}{N} \left[a_0 \sum_{i=1}^N \overline{\lambda_i^{\text{in}}} + a_1^{\text{in}} \sum_{i=1}^N (\overline{\lambda_i^{\text{in}}})^2 - C a_2^{\text{corr}} \lambda_0 \sum_{i=1}^N (\overline{\lambda_i^{\text{in}}})^2 \right], \quad (3.7)$$

and time constant

$$\tau := -\tau_w \frac{N}{\gamma_0} \left[a_1^{\text{out}} \sum_{i=1}^N \overline{\lambda_i^{\text{in}}} + (1 + C) a_2^{\text{corr}} \sum_{i=1}^N (\overline{\lambda_i^{\text{in}}})^2 \right]^{-1}. \quad (3.8)$$

The fixed point $\lambda_{\text{FP}}^{\text{out}}$ is asymptotically stable if and only if $\tau > 0$. For the translation-invariant scenario, the mean input rates are the same at all synapses so that the fixed point 3.7 reduces to

$$\lambda_{\text{FP}}^{\text{out}} = -\frac{a_0 + a_1^{\text{in}} \overline{\lambda^{\text{in}}} - C a_2^{\text{corr}} \lambda_0 \overline{\lambda^{\text{in}}}}{a_1^{\text{out}} + a_2^{\text{corr}} \overline{\lambda^{\text{in}}} (1 + C)}. \quad (3.9)$$

This is the generalization of equation 3.2 to the case of nonvanishing correlations.

We require $\lambda_{\text{FP}}^{\text{out}} > 0$ so as to ensure that we are in the linear regime of our piecewise-linear neuron model. For excitatory synapses, we require in addition that the weights w_i are positive. Hence, a realizable fixed point should have a rate $\lambda_{\text{FP}}^{\text{out}} > \max\{0, \lambda_0\}$. On the other hand, the weight values being bounded by w^{max} , the fixed point must be smaller than $\lambda_0 + \gamma_0 \overline{\lambda^{\text{in}}}$ since otherwise it cannot be reached. The learning dynamics would stop once all weights have reached their upper bound. Finally, the stability of the fixed point requires $\tau > 0$. The above requirements define conditions on the parameters $a_0, a_1^{\text{in}}, a_1^{\text{out}}, a_2^{\text{corr}}$, and λ_0 that we now explore in two examples. We note that a fixed point of the output rate implies a fixed point of the average weight $w := \sum_j \overline{\lambda_j^{\text{in}}} w_j / \sum_j \overline{\lambda_j^{\text{in}}}$ (cf. equation 3.3).

Example 1 (no linear terms). In standard correlation-based learning, there are usually no linear terms. Let us therefore set $a_0 = a_1^{\text{in}} = a_1^{\text{out}} = 0$. We assume that the average correlation is nonnegative, $C \geq 0$. Since we have $\gamma_0 > 0$, the term inside the square brackets of equation 3.8 must be negative in order to guarantee asymptotic stability of the fixed point. Since $\overline{\lambda_i^{\text{in}}} > 0$ for all i , a stable fixed point $\lambda_{\text{FP}}^{\text{out}}$ can therefore be achieved only if $a_2^{\text{corr}} < 0$, so that we end up with anti-Hebbian learning. The fixed point is

$$\lambda_{\text{FP}}^{\text{out}} = \lambda_0 \frac{C}{1 + C}. \quad (3.10)$$

As it should occur for positive rates, we must require $\lambda_0 > 0$. In addition, because $\lambda_{\text{FP}}^{\text{out}} < \lambda_0$, some of the weights w_i must be negative at the fixed point, which contradicts the assumption of excitatory synapses (see equation 3.3). We note that for $C \rightarrow 0$, the output rate approaches zero; that is, the neuron becomes silent. In summary, without linear terms, a rate model of correlation-based learning cannot show rate stabilization unless some weights are negative and $a_2^{\text{corr}} < 0$ (anti-Hebbian learning).

Example 2 (Hebbian learning). In standard Hebbian learning, the correlation term should have positive sign. We therefore set $a_2^{\text{corr}} > 0$. As before, we assume $C > 0$. From $\tau > 0$ in equation 3.8, we then obtain the condition that a_1^{out} be sufficiently negative in order to make the fixed point stable. In particular, we must have $a_1^{\text{out}} < 0$. Hence, firing-rate stabilization in rate-based Hebbian learning requires a linear “non-Hebbian” term $a_1^{\text{out}} < 0$.

3.2 Nonlinear Neuron Model. So far we have focused on a piecewise-linear rate model. Can we generalize the above arguments to a nonlinear neuron model? To keep the arguments transparent, we restrict our analysis to the translation-invariant scenario and assume identical mean input rates, $\overline{\lambda_i^{\text{in}}} = \overline{\lambda^{\text{in}}}$ for all i . As a specific example, we consider a neuron model with a sigmoidal activation (or gain) function $\lambda^{\text{out}} = g(u)$ where

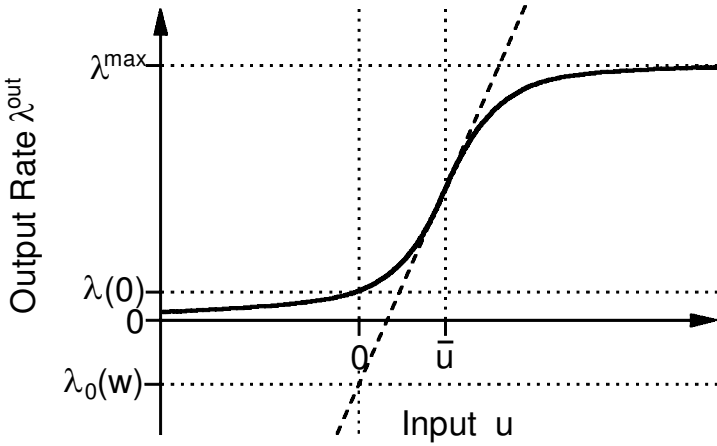


Figure 1: Sigmoidal activation function g . We plot the output rate λ^{out} (solid line) as a function of the neuron's input $u := N^{-1} \sum_i w_i \lambda_i^{\text{in}}$. The spontaneous output rate at $u = 0$ is $\lambda(0)$, and the maximum output rate is λ^{max} . The dashed line is a linearization of g around a mean input $\bar{u} > 0$. We note that here $\lambda_0(w) < 0$; cf. equation 3.11 and examples 1 and 4.

$u(t) = N^{-1} \sum_{i=1}^N w_i(t) \lambda_i^{\text{in}}(t)$ and the derivative $g'(u)$ is positive for all u . For vanishing input, $u = 0$, we assume a spontaneous output rate $\lambda^{\text{out}} = \lambda(0) > 0$. For $u \rightarrow -\infty$, the rate vanishes. For $u \rightarrow \infty$, the function g approaches the maximum output rate $\lambda^{\text{out}} = \lambda^{\text{max}}$ (see Figure 1).

Let us set $w = N^{-1} \sum_j w_j$. For a given mean weight w , we obtain a mean input $\bar{u} = w \bar{\lambda}^{\text{in}}$. If the fluctuations $\Delta \lambda_i^{\text{in}}(t) := \lambda_i^{\text{in}}(t) - \bar{\lambda}^{\text{in}}$ are small, we can linearize $g(u)$ around the mean input \bar{u} , that is, $g(u) = g(\bar{u}) + g'(\bar{u})(u - \bar{u})$ (cf. Figure 1). The linearization leads us back to equation 3.3,

$$\lambda^{\text{out}} = \lambda_0(w) + \gamma_0(w) \frac{1}{N} \sum_{i=1}^N w_i \lambda_i^{\text{in}}, \quad (3.11)$$

with $\lambda_0(w) = g(\bar{u}) - g'(\bar{u}) \bar{u}$ and $\gamma_0(w) = g'(\bar{u})$.

The discussion of output rate stabilization proceeds as in the linear case. Does a fixed point exist? The condition $d\lambda^{\text{out}}/dt = 0$ yields a fixed point that is given by equation 3.9 with λ_0 replaced by $\lambda_0(w)$. Thus, equation 3.9 becomes an implicit equation for $\lambda_{\text{FP}}^{\text{out}}$. A sigmoidal neuron can reach the fixed point if the latter lies in the range $0 < \lambda_{\text{FP}}^{\text{out}} < \lambda^{\text{max}}$. The local stability analysis does not change. As before, a fixed point of λ^{out} implies a fixed point of the mean weight w . The statement follows directly from equation 3.11 since $d\bar{\lambda}^{\text{out}}/dt = g'(u) \bar{\lambda}^{\text{in}} N^{-1} \sum_i dw_i/dt$ and $g' > 0$. That is, stabilization of the output rate implies a stabilization of the weights.

Let us summarize this section. The value of the fixed point $\lambda_{\text{FP}}^{\text{out}}$ is identical for both linear and the linearized rate model and given by equation 3.9. It is stable if τ in equation 3.8 is positive. If we kept higher-order terms in the expansion of equation 3.11, the fixed point may shift slightly, but as long as the nonlinear corrections are small, the fixed point remains stable (due to the continuity of the differential equation, 1.1). This highlights the fact that the learning rule, equation 3.1, with an appropriate choice of parameters achieves a stabilization of the output rate rather than a true normalization of the synaptic weights.

4 Plasticity Driven by Spikes

In this section, we generalize the concept of intrinsic rate stabilization to spike-time-dependent plasticity. We start by reviewing a learning rule where synaptic modifications are driven by the relative timing of pre- and post-synaptic spikes and defined by means of a learning window (Gerstner, Kempter, et al., 1996; Kempter et al., 1996); extensive experimental evidence supports this idea (Levy & Stewart, 1983; Bell, Han, Sugawara, & Grant, 1997; Markram et al., 1997; Zhang et al., 1998; Debanne et al., 1998; Bi & Poo, 1998, 1999; Egger, Feldmeyer, & Sakmann, 1999; Abbott & Munro, 1999; Feldman, 2000). The comparison of the spike-based learning rule developed below with a rate-based learning rule discussed above allows us to give the coefficients a_1^{in} , a_1^{out} , and a_2^{corr} in equation 3.1 a more precise meaning. In particular, we show that a_2^{corr} corresponds to the integral over the learning window. Anti-Hebbian learning in the above rate model (see example 1) can thus be identified with a negative integral over the learning window (Gerstner, Kempter, et al. 1996; Gerstner, Kempter, van Hemmen, & Wagner, 1997; Kempter et al., 1999a, 1999b).

4.1 Learning Rule. We are going to generalize equation 1.1 to spike-based learning. Changes in synaptic weights w_i are triggered by input spikes, output spikes, and the time differences between input and output spikes. To simplify the notation, we write the sequence $\{t_i^1, t_i^2, \dots\}$ of spike arrival times at synapse i in the form of a “spike train” $S_i^{\text{in}}(t) := \sum_m \delta(t - t_i^m)$ consisting of δ functions (real spikes of course have a finite width, but what counts here is the events “spike”). Similarly, the sequence of output spikes is denoted by $S^{\text{out}}(t) := \sum_n \delta(t - t^n)$. These notations allow us to formulate the spike-based learning rule as (Kempter, et al., 1999a; Kistler & van Hemmen, 2000; van Hemmen, 2000),

$$\begin{aligned} \tau_w \frac{d}{dt} w_i(t) = & a_0 + a_1^{\text{in}} S_i^{\text{in}}(t) + a_1^{\text{out}} S^{\text{out}}(t) + S_i^{\text{in}}(t) \int_{-\infty}^t dt' W(t - t') S^{\text{out}}(t') \\ & + S^{\text{out}}(t) \int_{-\infty}^t dt' W(-t + t') S_i^{\text{in}}(t'). \end{aligned} \quad (4.1)$$

As in equation 1.1, the coefficients a_0 , a_1^{in} , a_1^{out} , and the learning window W in general depend on the current weight value w_i . They may also depend on other local variables, such as the membrane potential or the calcium concentration. For the sake of clarity of presentation, we drop these dependencies and assume constant coefficients. We can, however, assume upper and lower bounds for the synaptic efficacies w_i ; that is, weight changes are zero if $w_i > w^{\text{max}}$ or $w_i < 0$ (see also the discussion in section 5.3). Again for the sake of clarity, in the remainder of this section we presume that all weights are within the bounds.

How can we interpret the terms on the right in equation 4.1? The coefficient a_0 is a simple decay or growth term that we will drop. S_i^{in} and S^{out} are sums of δ functions. We recall that integration of a differential equation $dx/dt = f(x, t) + a\delta(t)$ with arbitrary f yields a discontinuity at zero: $x(0_+) - x(0_-) = a$. Thus, an input spike arriving at synapse i at time t changes w_i at that time by a constant amount $\tau_w^{-1} a_1^{\text{in}}$ and a variable amount $\tau_w^{-1} \int_{-\infty}^t dt' W(t-t') S^{\text{out}}(t')$, which depends on the sequence of output spikes occurring at times earlier than t . Similarly, an output spike at time t results in a weight change $\tau_w^{-1} [a_1^{\text{out}} + \int_{-\infty}^t dt' W(t'-t) S_i^{\text{in}}(t')]$.

Our formulation of learning assumes discontinuous and instantaneous weight changes at the times of the occurrence of spikes, which may not seem very realistic. The formulation can be generalized to continuous and delayed weight changes that are triggered by pre- and postsynaptic spikes. Provided the gradual weight change terminates after a time that is small compared to the timescale of learning, all results derived below would basically be the same.

Throughout what follows, we assume that discontinuous weight changes are small compared to typical values of synaptic weights, meaning that learning is by small increments, which can be achieved for large τ_w . If τ_w is large, we can separate the timescale of learning from that of the neuronal dynamics, and the right-hand side of equation 4.1 is “self-averaging” (Kempster et al., 1999a; van Hemmen, 2000). The evolution of the weight vector (see equation 4.1) is then

$$\tau_w \frac{d}{dt} w_i(t) = a_1^{\text{in}} \overline{\langle S_i^{\text{in}} \rangle(t)} + a_1^{\text{out}} \overline{\langle S^{\text{out}} \rangle(t)} + \int_{-\infty}^{\infty} ds W(s) \overline{\langle S_i^{\text{in}}(t) S^{\text{out}}(t-s) \rangle}. \quad (4.2)$$

Our notation with angular brackets plus horizontal bar is intended to emphasize that we average over both the spike statistics given the rates (angular brackets) and the ensemble of rates (horizontal bar). An example will be given in section 4.4. Because weights change slowly and the statistical input ensemble is assumed to have stationary correlations, the two integrals in equation 4.1 have been replaced by a single integral in equation 4.2.

The formulation of the learning rule in equation 4.2 allows a direct link to rate-based learning as introduced in equations 1.1 and 3.1. Indeed, averaged spike trains $\overline{\langle S_i^{\text{in}} \rangle(t)}$ and $\overline{\langle S_i^{\text{out}} \rangle(t)}$ are identical to the mean firing rates $\overline{\lambda_i^{\text{in}}(t)}$ and $\overline{\lambda_i^{\text{out}}(t)}$, respectively, as defined in sections 2 and 3. In order to interpret the correlation terms, we have to be more careful. The terms $\overline{\langle S_i^{\text{in}} S_i^{\text{out}} \rangle}$ describe the correlation between input and output on the level of spikes. In section 4.3, we will give a detailed treatment of the correlation term. In the next section, we simplify the correlation term under the assumption of slowly changing firing rates.

4.2 Correlations in Rate-Based Learning. In this section, we relate the correlation term a_2^{corr} in equation 3.1 to the integral over the learning window W in equation 4.2. In order to make the transition from equation 4.2 to equation 3.1, two approximations are necessary (Kempter et al., 1999a, 1999b). First, we have to neglect correlations between input and output spikes apart from the correlations contained in the rates. That is, we make the approximation $\overline{\langle S_i^{\text{in}}(t) S_i^{\text{out}}(t-s) \rangle} \approx \overline{\lambda_i^{\text{in}}(t) \lambda_i^{\text{out}}(t-s)}$, the horizontal overbar denoting an average over the learning time. Second, we assume that these rates change slowly as compared to the width of the learning window W (cf. Figure 2a). Consequently we set $\lambda_i^{\text{out}}(t-s) \approx \lambda_i^{\text{out}}(t) - s [d/dt \lambda_i^{\text{out}}(t)]$ in the integrand of equation 4.2. With the above two assumptions, we obtain

$$\begin{aligned} \tau_w \frac{d}{dt} w_i(t) = & a_1^{\text{in}} \overline{\lambda_i^{\text{in}}} + a_1^{\text{out}} \overline{\lambda_i^{\text{out}}(t)} + \overline{\lambda_i^{\text{in}}(t) \lambda_i^{\text{out}}(t)} \int_{-\infty}^{\infty} ds W(s) \\ & - \overline{\lambda_i^{\text{in}}(t)} \frac{d}{dt} \overline{\lambda_i^{\text{out}}(t)} \int_{-\infty}^{\infty} ds s W(s). \end{aligned} \quad (4.3)$$

The last term on the right in equation 4.3 has been termed “differential-Hebbian” (Xie & Seung, 2000). Equation 4.3 is equivalent to equation 3.1, if we neglect the differential-Hebbian term and set $a_2^{\text{corr}} = \int ds W(s)$ and $a_0 = 0$. We note, however, that $\overline{\lambda_i^{\text{in}}(t) \lambda_i^{\text{out}}(t)} = \overline{\lambda_i^{\text{in}} \cdot \lambda_i^{\text{out}}(t)} + \overline{\Delta \lambda_i^{\text{in}}(t) \Delta \lambda_i^{\text{out}}(t)}$ still contains the correlations between the fluctuations in the input and output rates (but we had to assume that these correlations are slow compared to the width of the learning window). The condition of slowly changing rates will be dropped in the next section.

We saw in example 1 that with rate-based learning, intrinsic normalization of the output rate is achievable once $a_2^{\text{corr}} < 0$. We now argue that $a_2^{\text{corr}} < 0$, which has been coined “anti-Hebbian” in the framework of rate-based learning, can still be “Hebbian” in terms of spike-based learning. The intuitive reason is that the learning window $W(s)$ as sketched in Figure 2a has two parts, a positive and a negative one. Even if the integral over the learning window is negative, the size of weight changes in the positive part

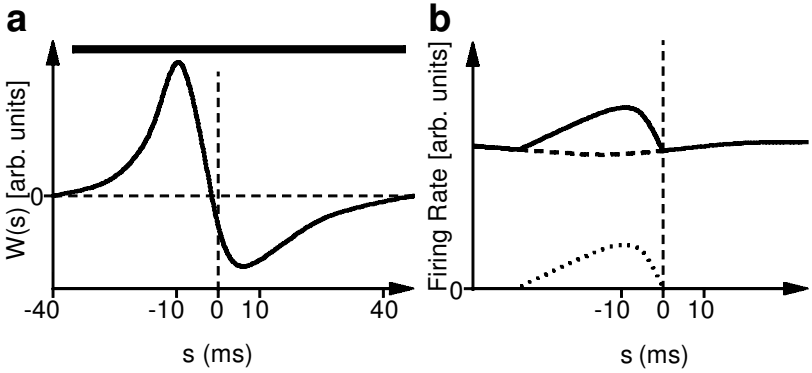


Figure 2: (a) Learning window W (arbitrary units) as a function of the delay $s = t_i^m - t^n$ between presynaptic spike arrival time t_i^m and postsynaptic firing time t^n (schematic). If $W(s)$ is positive (negative) for some s , the synaptic efficacy w_i is increased (decreased). The increase of w_i is most efficient if a presynaptic spike arrives a few milliseconds before the postsynaptic neuron starts firing. For $|s| \rightarrow \infty$, we have $W(s) \rightarrow 0$. The bar denotes the width of the learning window. The qualitative time course of the learning window is supported by experimental results (Levy & Stewart, 1983; Markram et al., 1997; Zhang et al., 1998; Debanne et al., 1998; Bi & Poo, 1998; Feldman, 2000; see also the reviews by Brown & Chattarji, 1994; Linden, 1999; Paulsen & Sejnowski, 2000; and Bi & Poo, 2001). The learning window can be described theoretically by a phenomenological model with microscopic variables (Senn et al., 1997, 2001; Gerstner, Kempter, van Hemmen, & Wagner, 1998). (b) Spike-spike correlations (schematic). The correlation function $\text{Corr}_i(s, \mathbf{w}) / \lambda_i^{\text{in}}$ (full line) defined in equation 4.10 is the sum of two terms: (1) the causal contribution of an input spike at synapse i at time s to the output firing rate at time $s = 0$ is given by the time-reverted EPSP, that is, $\gamma_0 N^{-1} w_i \epsilon(-s)$ (dotted line); (2) the correlations between the mean firing rates $\lambda_i^{\text{in}}(t) \lambda^{\text{out}}(t - s) / \lambda_i^{\text{in}}$ are indicated by the dashed line. For independent inputs with stationary mean, the rate contribution would be constant and equal to λ^{out} . In the figure, we have sketched the situation where the rates themselves vary with time.

can be large enough to pick up correlations. In other words, the two approximations necessary to make the transition from equation 4.2 to equation 4.3 are, in general, not justified. An example is given in the the appendix. In order to make our argument more precise, we have to return to spike-based learning.

4.3 Spike-Spike Correlations. In the previous section, we replaced $\langle S_i^{\text{in}}(t) S^{\text{out}}(t - s) \rangle$ by the temporal correlation between slowly changing rates. This replacement leads to an oversimplification of the situation be-

cause it does not take into account the correlations on the level of spikes or fast changes of the rates. Here we present the correlation term in its general form.

Just as in section 3, we require learning to be a slow process so that $\tau_w \gg p\Delta_t$ or $\tau_w \gg \Delta_{\text{corr}}$ (see section 2 for a definition of $p\Delta_t$ and Δ_{corr}). The correlation term can then be evaluated for constant weights w_i , $1 \leq i \leq N$. For input S_i^{in} with some given statistical properties that do not change on the slow timescale of learning, the correlations between $S_i^{\text{in}}(t)$ and $S^{\text{out}}(t-s)$ depend on only the time difference s and the current weight vector $\mathbf{w}(t)$,

$$\text{Corr}_i[s, \mathbf{w}(t)] := \overline{\langle S_i^{\text{in}}(t) S^{\text{out}}(t-s) \rangle}. \quad (4.4)$$

The horizontal bar refers, as before, to the temporal average introduced through the separation of timescales (cf. equation 2.1). Substituting equation 4.4 into equation 4.2, we obtain the dynamics of weight changes,

$$\tau_w \frac{d}{dt} w_i(t) = a_1^{\text{in}} \overline{\lambda_i^{\text{in}}} + a_1^{\text{out}} \overline{\lambda^{\text{out}}(t)} + \int_{-\infty}^{\infty} ds W(s) \text{Corr}_i[s, \mathbf{w}(t)]. \quad (4.5)$$

Equation 4.5 is completely analogous to equation 3.1 except that the correlations between rates have been replaced by those between spikes. We may summarize equation 4.5 by saying that learning is driven by correlations on the timescale of the learning window.

4.4 Poisson Neuron. To proceed with the analysis of equation 4.5, we need to determine the correlations Corr_i between input spikes at synapse i and output spikes. The correlations depend strongly on the neuron model under consideration. As a generalization to the linear rate-based neuron model in section 3.1, we study an inhomogeneous Poisson model that generates spikes at a rate $\lambda_{\text{sin}}^{\text{out}}(t)$. For a membrane potential $u(t) \leq \mathcal{G}$, the model neuron is quiescent. For $u > \mathcal{G}$, the rate $\lambda_{\text{sin}}^{\text{out}}(t)$ is proportional to $u - \mathcal{G}$. We call this model the piecewise-linear Poisson neuron or, for short, the Poisson neuron.

4.4.1 Definition. The input to the neuron consists of N Poissonian spike trains with time-dependent intensities $\langle S_i^{\text{in}} \rangle(t) = \lambda_i^{\text{in}}(t)$, $1 \leq i \leq N$ (for the mathematics of an inhomogeneous Poisson process, we refer to appendix A of Kempster, Gerstner, van Hemmen, & Wagner, 1998). As in the rate model of section 3, the normalized correlation function for the rate fluctuations at synapse i and j is defined by equation 2.5. The input scenarios are the same as in section 2—static-pattern scenario or translation-invariant scenario.

A spike arriving at t_i^f at synapse i evokes a postsynaptic potential (PSP) with time course $w_i \epsilon(t - t_i^f)$, which we assume to be excitatory (EPSP). We

impose the condition $\int_0^\infty ds \epsilon(s) = 1$ and require causality, that is, $\epsilon(s) = 0$ for $s < 0$. The amplitude of the EPSP is given by the synaptic efficacy $w_i \geq 0$. The membrane potential u of the neuron is the linear superposition of all contributions,

$$\begin{aligned}
 u(t) &= \frac{1}{N} \sum_{i=1}^N w_i(t) \sum_f \epsilon(t - t_i^f) \\
 &= \frac{1}{N} \sum_{i=1}^N w_i(t) \int_{-\infty}^{+\infty} ds \epsilon(s) S_i^{\text{in}}(t - s).
 \end{aligned}
 \tag{4.6}$$

Since ϵ is causal, only spike times $t_i^f < t$ count. The sums run over all spike arrival times t_i^f of all synapses. $S_i^{\text{in}}(t) = \sum_f \delta(t - t_i^f)$ is the spike train at synapse i .

In the Poisson neuron, output spikes are generated stochastically with a time-dependent rate $\lambda_{\text{Sin}}^{\text{out}}$ that depends linearly on the membrane potential whenever u is above the threshold $\vartheta = -\lambda_0/\gamma_0$,

$$\lambda_{\text{Sin}}^{\text{out}}(t) = [\lambda_0 + \gamma_0 u(t)]_+ = \left[\lambda_0 + \frac{\gamma_0}{N} \sum_{i=1}^N w_i(t) \sum_f \epsilon(t - t_i^f) \right]_+. \tag{4.7}$$

Here $\gamma_0 > 0$ and λ_0 are constants. As before, the notation with square brackets $[\cdot]_+$ ensures that $\lambda_{\text{Sin}}^{\text{out}}$ be nonnegative. The brackets will be dropped in the following. The subscript Sin of $\lambda_{\text{Sin}}^{\text{out}}$ indicates that the output rate depends on the specific realization of the input spike trains. We note that the spike-generation process is independent of previous output spikes. In particular, the Poisson neuron model does not show refractoriness. We drop this restriction and include refractoriness in section 4.6.

4.4.2 Expectation Values. In order to evaluate equations 4.4 and 4.5, we first have to calculate the expectation values $\langle \cdot \rangle$ (i.e., perform averages over the spike statistics of both input *and* output given the rates) and then average over time. By definition of the rate of an inhomogeneous Poisson process, the expected input is $\langle S_i^{\text{in}} \rangle = \lambda_i^{\text{in}}$. Taking advantage of equations 4.6 and 4.7, we find that the expected output rate $\langle S^{\text{out}} \rangle = \langle \lambda_{\text{Sin}}^{\text{out}} \rangle = \lambda^{\text{out}}$, which is the rate contribution of all synapses to an output spike at time t , is given by

$$\langle S^{\text{out}} \rangle(t) = \lambda^{\text{out}}(t) = \lambda_0 + \frac{\gamma_0}{N} \sum_{i=1}^N w_i(t) \int ds \epsilon(s) \lambda_i^{\text{in}}(t - s), \tag{4.8}$$

where the final term is a convolution of ϵ with the input rate.

Next we consider correlations between input and output, $\langle S_i^{\text{in}}(t) S^{\text{out}}(t - s) \rangle$ —the expectation of finding an input spike at synapse i at time t and an output spike at $t - s$. Since this expectation is defined as a joint probability density, it equals the probability density $\lambda_i^{\text{in}}(t)$ for an input spike at time t times the conditional probability density $\langle S^{\text{out}}(t - s) \rangle_{\{i,t\}}$ of observing an output spike at time $t - s$ given the input spike at synapse i at t ; that is, $\langle S_i^{\text{in}}(t) S^{\text{out}}(t - s) \rangle = \lambda_i^{\text{in}}(t) \langle S^{\text{out}}(t - s) \rangle_{\{i,t\}}$. Within the framework of the linear Poisson neuron, the term $\langle S^{\text{out}}(t - s) \rangle_{\{i,t\}}$ equals the sum of the expected output rate $\lambda^{\text{out}}(t - s)$ in equation 4.8, and the specific contribution $\gamma_0 N^{-1} w_i \epsilon(-s)$ of a single input spike at synapse i (cf. equation 4.7). To summarize, we get (Kempster et al., 1999a),

$$\langle S_i^{\text{in}}(t) S^{\text{out}}(t - s) \rangle = \lambda_i^{\text{in}}(t) \left[\frac{\gamma_0}{N} w_i(t) \epsilon(-s) + \lambda^{\text{out}}(t - s) \right]. \quad (4.9)$$

Due to causality of ϵ , the first term on the right, inside the square brackets of equation 4.9, must vanish for $s > 0$.

4.4.3 Temporal Average. In order to evaluate the correlation term, equation 4.4, we need to take the temporal average of equation 4.9, that is,

$$\text{Corr}_i[s, \mathbf{w}(t)] = \overline{\lambda_i^{\text{in}}} \frac{\gamma_0}{N} w_i(t) \epsilon(-s) + \overline{\lambda_i^{\text{in}}(t) \lambda^{\text{out}}(t - s)}. \quad (4.10)$$

For excitatory synapses, the first term gives for $s < 0$ a positive contribution to the correlation function, as it should be (see Figure 2b). We recall that $s < 0$ means that a presynaptic spike precedes postsynaptic firing. The second term on the right in equation 4.10 describes the correlations between the input and output rates. These rates and correlations between them can, in principle, vary on an arbitrary fast timescale (for modeling papers, see Gerstner, Kempster, et al., 1996; Xie & Seung, 2000). We assume, however, for reasons of transparency that the mean input rates $\overline{\lambda_i^{\text{in}}(t)} = \overline{\lambda_i^{\text{in}}}$ are constant for all i (see also the input scenarios in section 2). The mean output rate $\overline{\lambda^{\text{out}}(t)}$ and the weights $w_i(t)$ may vary on the slow timescale of learning.

4.4.4 Learning Equation. Substituting equations 2.5, 4.8, and 4.10 into 4.5, we find

$$\begin{aligned} \tau_w \frac{d}{dt} w_i(t) &= a_1^{\text{in}} \overline{\lambda_i^{\text{in}}} + a_1^{\text{out}} \overline{\lambda^{\text{out}}(t)} + \overline{\lambda_i^{\text{in}}} \cdot \overline{\lambda^{\text{out}}(t)} \int_{-\infty}^{\infty} ds W(s) + \frac{\gamma_0}{N} \sum_{j=1}^N \\ &\times w_j(t) \left[Q_{ij} \overline{\lambda_i^{\text{in}}} \cdot \overline{\lambda_j^{\text{in}}} + \delta_{ij} \overline{\lambda_i^{\text{in}}} \int_{-\infty}^0 ds W(s) \epsilon(-s) \right], \quad (4.11) \end{aligned}$$

where

$$Q_{ij} := \int_{-\infty}^{\infty} ds W(s) \int_0^{\infty} ds' \epsilon(s') C_{ij}(s + s'). \quad (4.12)$$

The important factor in equation 4.12 is the input correlation function C_{ij} (as it appears in equation 2.5) (low-pass) filtered by learning window W and the postsynaptic potential ϵ . For all input ensembles with $C_{ij}(\tau) = C_{ij}(-\tau)$ (viz. temporal symmetry), we have $Q_{ij} = Q_{ji}$. Hence, the matrix (Q_{ij}) can be diagonalized and has real eigenvalues. In particular, this remark applies to the static-pattern scenario and the translation-invariant scenario introduced in section 2. It does not apply, for example, to sequence learning; cf. (Herz et al., 1988, 1989; van Hemmen et al., 1990; Gerstner et al., 1993; Abbott & Blum, 1996; Gerstner & Abbott, 1997).

To get some preliminary insight into the nature of solutions to equation 4.11, let us suppose that all inputs have the same mean $\overline{\lambda_i^{\text{in}}} = \overline{\lambda^{\text{in}}}$. As we will see, for a broad parameter range, the mean output rate will be attracted toward a fixed point so that the first three terms on the right-hand side of equation 4.11 almost cancel each other. The dynamics of pattern formation in the weights w_i is then dominated by the largest positive eigenvalue of the matrix $[Q_{ij} + \delta_{ij} (\overline{\lambda^{\text{in}}})^{-1} \int ds W(s) \epsilon(-s)]$. The eigenvectors of this matrix are identical to those of the matrix (Q_{ij}) . The eigenvalues of (Q_{ij}) can be positive even though $a_2^{\text{corr}} = \int ds W(s)$ is negative. A simple example is given in the appendix.

To summarize this section, we have solved the dynamics of spike-time-dependent plasticity defined by equation 4.2 for the Poisson neuron, a (piecewise) linear spiking neuron model. In particular, we have succeeded in evaluating the spike-spike correlations between presynaptic input and postsynaptic firing.

4.5 Stabilization in Spike-Based Learning. We are aiming at a stabilization of the output rate in analogy to equations 3.7 and 3.8. To arrive at a differential equation for $\overline{\lambda^{\text{out}}}$, we multiply equation 4.11 by $\overline{\lambda_i^{\text{in}}} \gamma_0 / N$ and sum over i . Similarly to equation 3.5, we assume that

$$Q := \frac{1}{N} \sum_{i=1}^N (\overline{\lambda_i^{\text{in}}})^2 Q_{ij} \quad (4.13)$$

is independent of the index j . To keep the arguments simple, we assume in addition that all inputs have the same mean $\overline{\lambda_i^{\text{in}}} = \overline{\lambda^{\text{in}}}$. For the translation-invariant scenario, both assumptions can be verified. With these simplifications, we arrive at a differential equation for the mean output rate $\overline{\lambda^{\text{out}}}$,

which is identical to equation 3.6, that is, $\tau \frac{d}{dt} \overline{\lambda^{\text{out}}} = \lambda_{\text{FP}}^{\text{out}} - \overline{\lambda^{\text{out}}}$ with fixed point

$$\lambda_{\text{FP}}^{\text{out}} = \frac{\tau}{\tau_w} \gamma_0 \left[a_1^{\text{in}} \left(\overline{\lambda^{\text{in}}} \right)^2 - (Q + \beta) \lambda_0 \right] \quad (4.14)$$

and time constant

$$\tau = -\frac{\tau_w}{\gamma_0} \left[a_1^{\text{out}} \overline{\lambda^{\text{in}}} + \left(\overline{\lambda^{\text{in}}} \right)^2 \int_{-\infty}^{\infty} ds W(s) + Q + \beta \right]^{-1}, \quad (4.15)$$

where

$$\beta := \frac{\overline{\lambda^{\text{in}}}}{N} \int_{-\infty}^0 ds W(s) \epsilon(-s) \quad (4.16)$$

is the contribution that is due to the spike-spike correlations between pre- and postsynaptic neuron. A comparison with equations 3.7 and 3.8 shows that we may set $a_0 = 0$ and $a_2^{\text{corr}} = \int ds W(s)$. The average correlation is to be replaced by

$$C = \frac{Q + \beta}{\left(\overline{\lambda^{\text{in}}} \right)^2 \int ds W(s)}. \quad (4.17)$$

All previous arguments regarding intrinsic normalization apply. The only difference is the new definition of the factor C in equation 4.17 as compared to equation 3.5. This difference is, however, important. First, the sign of $\int ds W(s)$ enters the definition of C in equation 4.17. Furthermore, an additional term β appears. It is due to the extra spike-spike correlations between presynaptic input and postsynaptic firing. Note that β is of order N^{-1} . Thus, it is small compared to the first term in equation 4.17 whenever the number of nonzero synapses is large (Kempter et al., 1999a). If we neglect the second term ($\beta = 0$) and if we set $\epsilon(s) = \delta(s)$ and $W(s) = a_2^{\text{corr}} \delta(s)$, equation 4.17 is identical to equation 3.5. On the other hand, β is of order $\overline{\lambda^{\text{in}}}$ in the input rates, whereas Q is of order $(\overline{\lambda^{\text{in}}})^2$. Thus, β becomes important whenever the mean input rates are low. As shown in the appendix, β is, for a realistic set of parameters, of the same order of magnitude as Q .

Let us now study the intrinsic normalization properties for two combinations of parameters.

Example 3 (positive integral). We consider the case $\int_{-\infty}^{+\infty} ds W(s) > 0$ and assume that correlations in the input on the timescale of the learning window are positive so that $Q > 0$. Furthermore, for excitatory synapses, it is

natural to assume $\beta > 0$. In order to guarantee the stability of the fixed point, we must require $\tau > 0$, which yields the condition

$$a_1^{\text{out}} < -\overline{\lambda^{\text{in}}} \int_{-\infty}^{\infty} ds W(s) - \frac{Q + \beta}{\lambda^{\text{in}}} < 0. \quad (4.18)$$

Hence, if a_1^{out} is sufficiently negative, output rate stabilization is possible even if the integral over the learning window has a positive sign. In other words, for a positive integral over the learning window, a non-Hebbian “linear” term is necessary.

Example 4 (no linear terms). In standard Hebbian learning, all linear terms vanish. Let us therefore set $a_1^{\text{in}} = a_1^{\text{out}} = 0$. As in example 3, we assume $Q + \beta > 0$. To guarantee the stability of the fixed point, we must now require

$$\int_{-\infty}^{\infty} ds W(s) < -\frac{Q + \beta}{(\lambda^{\text{in}})^2} < 0. \quad (4.19)$$

Thus, a stable fixed point of the output rate is possible if the integral over the learning window is sufficiently negative. The fixed point is

$$\lambda_{\text{FP}}^{\text{out}} = \lambda_0 \frac{Q + \beta}{Q + \beta + (\lambda^{\text{in}})^2 \int ds W(s)}. \quad (4.20)$$

Note that the denominator is negative because of equation 4.19. For a fixed point to exist, we need $\lambda_{\text{FP}}^{\text{out}} > 0$; hence, $\lambda_0 < 0$ (cf. Figure 1). We emphasize that in contrast to example 1, all weights at the fixed point can be positive, as should be the case for excitatory synapses.

To summarize the two examples, a stabilization of the output rate is possible for positive integral with a non-Hebbian term $a_1^{\text{out}} < 0$ or for a negative integral and vanishing non-Hebbian terms.

4.6 Spiking Model with Refractoriness. In this section, we generalize the results to a more realistic model of a neuron: the spike response model (SRM) (Gerstner & van Hemmen, 1992; Gerstner, van Hemmen, & Cowan, 1996; Gerstner, 2000). Two aspects change with respect to the Poisson neuron. First, we include arbitrary refractoriness into the description of the neuron in equation 4.6. In contrast to a Poisson process, events in disjoint intervals are not independent. Second, we replace the linear firing intensity in equation 4.7 by a nonlinear one (and linearize only later on).

Let us start with the neuronal membrane potential. As before, each input spike evokes an excitatory postsynaptic potential described by a response kernel $\epsilon(s)$ with $\int_0^\infty ds \epsilon(s) = 1$. After each output spike, the neuron undergoes a phase of refractoriness, which is described by a further response kernel η . The total membrane potential is

$$u(t | \hat{t}) = \eta(t - \hat{t}) + \frac{1}{N} \sum_{i=1}^N w_i(t) \sum_f \epsilon(t - t_i^f), \quad (4.21)$$

where \hat{t} is the last output spike of the postsynaptic neuron. As an example, let us consider the refractory kernel

$$\eta(s) = -\eta_0 \exp\left(-\frac{s}{\tau_\eta}\right) \Theta(s), \quad (4.22)$$

where $-\eta(0) = \eta_0 > 0$ and $\tau_\eta > 0$ are parameters and $\Theta(s)$ is the Heaviside function, that is, $\Theta(s) = 0$ for $s \leq 0$ and $\Theta(s) = 1$ for $s > 0$. With this definition of η , the SRM is related to the standard integrate-and-fire model. A difference is that in the integrate-and-fire model, the voltage is reset after each firing to a fixed value, whereas in equations 4.21 and 4.22, the membrane potential is reset at time $t = \hat{t}$ by an amount $-\eta_0 - \eta(t - t')$, which depends on the time t' of the previous spike.

In analogy to equation 4.7, the probability of spike firing depends on the momentary value of the membrane potential. More precisely, the stochastic intensity of spike firing is a (nonlinear) function of $u(t | \hat{t})$,

$$\lambda^{\text{out}}(t | \hat{t}) = f[u(t | \hat{t})], \quad (4.23)$$

with $f(u) \rightarrow 0$ for $u \rightarrow -\infty$. For example, we may take $f(u) = \lambda_0 \exp[\gamma_0 (u - \vartheta)/\lambda_0]$ with parameters $\lambda_0, \gamma_0, \vartheta > 0$. For $\gamma_0 \rightarrow \infty$, spike firing becomes a threshold process and occurs whenever u reaches the threshold ϑ from below. For $\vartheta = 0$ and $|\gamma_0 u/\lambda_0| \ll 1$, we are led back to the linear model of equation 4.7. A systematic study of escape functions f can be found in Plesser and Gerstner (2000).

We now turn to the calculation of the correlation function Corr_i as defined in equation 4.4. More precisely, we are interested in finding the generalization of equation 4.10 (for Poisson neurons) to spiking neurons with refractoriness. Due to the dependence of the membrane potential in equation 4.21 upon the last firing time \hat{t} , the first term on the right-hand side of equation 4.10 will become more complicated. Intuitively, we have to average over the firing times \hat{t} of the postsynaptic neuron. The correct averaging can be found from the theory of population dynamics studied in Gerstner (2000). To keep the formulas as simple as possible, we consider constant input rates $\lambda_i^{\text{in}}(t) \equiv \lambda^{\text{in}}$ for all i . If many input spikes arrive on average

within the postsynaptic integration time set by the kernel ϵ , the membrane potential fluctuations are small, and we can linearize the dynamics about the mean trajectory of the membrane potential:

$$\tilde{u}(t | \hat{t}) = \eta(t - \hat{t}) + \overline{\lambda^{\text{in}}} \frac{1}{N} \sum_{i=1}^N w_i. \tag{4.24}$$

To linear order in the membrane potential fluctuations, the correlations are (Gerstner, 2000),

$$\text{Corr}_i(s, \mathbf{w}) = \overline{\lambda^{\text{in}}} P_i(-s) + \overline{\lambda^{\text{in}} \lambda^{\text{out}}}, \tag{4.25}$$

where

$$P_i(s) = \overline{\lambda^{\text{out}}} \frac{w_i}{N} \frac{d}{ds} \int_0^\infty dx \mathcal{L}(x) \epsilon(s - x) + \int_0^s d\hat{t} f[u(s | \hat{t})] S(s | \hat{t}) P_i(\hat{t}). \tag{4.26}$$

The survivor function S is defined as

$$S(s | \hat{t}) = \exp \left\{ - \int_{\hat{t}}^s dt' f[u(t' | \hat{t})] \right\}, \tag{4.27}$$

the mean output rate is

$$\overline{\lambda^{\text{out}}(t)} = \left[\int_0^\infty dt S(t | 0) \right]^{-1}, \tag{4.28}$$

and the “filter” \mathcal{L} is

$$\mathcal{L}(x) = \int_x^\infty d\xi \frac{d}{d\xi} f[u(\xi | x)] S(\xi | 0). \tag{4.29}$$

In equations 4.26 through 4.29, the membrane potential u is given by the mean trajectory \tilde{u} in equation 4.24.

Equations 4.25 and 4.26 are our main result. All considerations regarding normalization as discussed in the preceding subsection are valid with the function Corr_i defined in equations 4.25 and 4.26 and an output rate $\overline{\lambda^{\text{out}}}$

defined in equation 4.28. In passing, we note that equation 4.28 defines the gain function of the corresponding rate model, $\overline{\lambda^{\text{out}}} = g[\overline{\lambda^{\text{in}}} N^{-1} \sum_i w_i]$.

Let us now discuss equation 4.26 in more detail. The first term on its right-hand side describes the immediate influence of the EPSP ϵ on the firing probability of the postsynaptic neuron. The second term describes the reverberation of the primary effect that occurs one interspike interval later. If the interspike interval distribution is broad, the second term on the right is small and may be neglected (Gerstner, 2000). In order to understand the relation of equations 4.25 and 4.26 to equation 4.10, let us study two examples.

Example 5 (no refractoriness). We want to show that in the limit of no refractoriness, the spiking neuron model becomes identical with a Poisson model. If we neglect refractoriness ($\eta \equiv 0$), the mean membrane potential in equation 4.24 is $\tilde{u} = \overline{\lambda^{\text{in}}} N^{-1} \sum_i w_i$. The output rate is $\overline{\lambda^{\text{out}}} = f(\tilde{u})$ (cf. equation 4.23), and we have $S(s | 0) = \exp[-\overline{\lambda^{\text{out}}} s]$ (see equation 4.27). In equation 4.29, we set $\gamma_0 = df/du$ evaluated at \tilde{u} and find $P_i(s) = \epsilon(s) w_i \gamma_0 / N$. Hence, equation 4.25 reduces to equation 4.10, as it should be.

Example 6 (high- and low-noise limit). If we take an exponential escape rate $f(u) = \lambda_0 \exp[\gamma_0(u - \vartheta)/\lambda_0]$, the high-noise limit is found for $|\gamma_0(u - \vartheta)/\lambda_0| \ll 1$ and the low-noise limit for $|\gamma_0(u - \vartheta)/\lambda_0| \gg 1$. In the high-noise limit, it can be shown that the filter \mathcal{L} is relatively broad. More precisely, it starts with an initial value $\mathcal{L}(0) > 0$ and decays slowly, $\mathcal{L}(x) \rightarrow 0$ for $x \rightarrow \infty$ (Gerstner, 2000). Since decay is small, we set $\frac{d}{dx} \mathcal{L}(x) \approx 0$ for $x > 0$. This yields $P_i(s) = w_i N^{-1} \overline{\lambda^{\text{out}}} \mathcal{L}(0) \epsilon(s) + \int_0^s d\hat{t} f[u(s | \hat{t})] S(s | \hat{t}) P_i(\hat{t})$. Thus, the correlation function Corr_i contains a term proportional to $\epsilon(-s)$ as in the linear model. In the low-noise limit ($\gamma_0 \rightarrow \infty$), we retrieve the threshold process, and the filter \mathcal{L} approaches a δ -function— $\mathcal{L}(x) = d \delta(x)$ with some constant $d > 0$ (Gerstner, 2000). In this case, the correlation function Corr_i contains a term proportional to the *derivative* of the EPSP: $d\epsilon(s)/ds$.

To summarize this section, we have shown that it is possible to calculate the correlation function $\text{Corr}_i(s, \mathbf{w})$ for a spiking neuron model with refractoriness. Once we have obtained the correlation function, we can use it in the learning dynamics, equation 4.5. We have seen that the correlation function depends on the noise level. This theoretical result is in agreement with experimental measurements on motoneurons (Fetz & Gustafsson, 1983; Poliakov, Powers, Sawczuk, & Binder, 1996). It was found that for high noise levels, the correlation function contained a peak with a time course roughly proportional to the postsynaptic potential. For low noise, however, the time course of the peak was similar to the derivative of the postsynaptic potential. Thus, the (piecewise-linear) Poisson neuron introduced in section 4.4 is a valid model in the high-noise limit.

5 Discussion

In this article we have compared learning rules at the level of spikes with those at the level of rates. The learning rules can be applied to both linear and nonlinear neuron models. We discuss our results in the context of the existing experimental and theoretical literature.

5.1 Experimental Results. Rate normalization has been found in experiments performed by Turrigiano, Leslie, Desai, Rutherford, and Nelson (1998). They blocked GABA-mediated inhibition in a cortical culture, which initially raised activity. After about two days, the firing rate returned close to the control value, and at the same time, all synaptic strengths decreased. Conversely, a pharmacologically induced firing-rate reduction leads to synaptic strengthening. Again, the output rate was normalized. Turrigiano et al. (1998) suggest that rate normalization is achieved by multiplying all weights by the same factor. In our ansatz, however, weights are normalized subtractively, that is, by addition or subtraction of a fixed amount from all weights. The discrepancy between the experiment and our model can be resolved by assuming—beyond upper and lower bounds for individual weights—that the learning coefficients themselves depend on the weights too. The extended ansatz exceeds, however, the scope of this article. For additional mechanisms contributing to the activity-dependent stabilization of firing rates we refer to Desai, Rutherford, and Turrigiano (1999), Turrigiano and Nelson (2000), and Turrigiano (1999).

In our model, several scenarios for intrinsic rate normalization are possible, depending on the choice of parameters for the Hebbian and non-Hebbian terms. Let us discuss each of them in turn.

5.1.1 Sign of Correlation Term. Rate normalization in our model is most easily achieved if the integral of the learning window W is negative, even though this is not necessary (see examples 1 and 4). On the basis of neurophysiological data known at present (Markram et al., 1997; Zhang et al., 1998; Debanne et al., 1998; Bi & Poo, 1998; Feldman, 2000), the sign of the integral over the learning window cannot be decided (see also the reviews by Linden, 1999; Bi & Poo, 2001). We have shown that a negative integral corresponds, in terms of rate coding, to a negative coefficient a_2^{corr} and hence to an anti-Hebbian rule. Nevertheless, the learning rule can pick up positive correlations if the input contains temporal correlations C_{ij} on the timescale of the learning window. An example is given in the appendix.

5.1.2 Role of Non-Hebbian Terms. We saw in examples 1 and 4 that the linear terms in the learning rule are not necessary for rate normalization. Nevertheless, we saw from example 2 or 3 that a negative coefficient a_1^{out} helps. A value $a_1^{\text{out}} < 0$ means that in the absence of presynaptic activation, postsynaptic spiking alone induces heterosynaptic long-term depression.

This has been seen, for example, in hippocampal slices (Pockett, Brookes, & Bindman, 1990; Christofi, Nowicky, Bolsover, & Bindman, 1993; see also the reviews of Artola & Singer, 1993; Brown & Chattarji, 1994; Linden, 1999). For cortical slices, both $a_1^{\text{out}} > 0$ and $a_1^{\text{in}} < 0$ have been reported (Volgushev, Voronin, Chistiakova, & Singer, 1994).

On the other hand, a positive effect of presynaptic spikes on the synapses in the absence of postsynaptic spiking ($a_1^{\text{in}} > 0$) helps to keep the fixed point in the range of positive rates $\lambda_{\text{FP}}^{\text{out}} > 0$ (see equation 3.7). Some experiments indeed suggest that presynaptic activity alone results in homosynaptic long-term potentiation ($a_1^{\text{in}} > 0$) (Bliss & Collingridge, 1993; Urban & Barrionuevo, 1996; Bell et al., 1997; see also Brown & Chattarji, 1994).

5.1.3 Zero-Order Term. As we have seen from equation 3.7, spontaneous weight growth $a_0 > 0$ helps keep the fixed point of the rate in the positive range. Turrigiano et al. (1998) chronically blocked cortical culture activity and found an increase of synaptic weights, supporting $a_0 > 0$.

5.2 Spike-Time-Dependent Learning Window. In Figure 2a we indicated a learning window with two parts: a positive part (potentiation) and a negative one (depression). Such a learning window is in accordance with experimental results (Markram et al., 1997; Zhang et al., 1998; Debanne et al., 1998; Bi & Poo, 1998; Feldman, 2000). We may ask, however, whether there are theoretical reasons for this time dependence of the learning window? For excitatory synapses, a presynaptic input spike that precedes postsynaptic firing may be the cause of the postsynaptic activity or, at least, “takes part in firing it” (Hebb, 1949). Thus, a literal understanding of the Hebb rule suggests that for excitatory synapses, the learning window $W(s)$ is positive for $s < 0$, if s is the difference between the times of occurrence of an input and an output spike. (Recall that $s < 0$ means that a presynaptic spike precedes postsynaptic spiking; see also Figure 2b.) In fact, a positive learning phase (potentiation) for $s < 0$ has been found to be an important ingredient in models of sequence learning (Herz et al., 1988, 1989; van Hemmen et al., 1990; Gerstner et al., 1993; Abbott & Blum, 1996; Gerstner & Abbott, 1997). If the positive phase is followed by a negative phase (depression) for $s > 0$ as in Figure 2a, the learning rule acts as a temporal contrast filter enhancing the detection of temporal structure in the input (Gerstner, Kempster, et al., 1996; Kempster et al., 1996, 1999a; Song et al., 2000; Xie & Seung, 2000). A two-phase learning rule with both potentiation and depression was, to a certain degree, a theoretical prediction. The advantages of such a learning rule have been realized in models (Gerstner, Kempster, et al., 1996; Kempster et al., 1996) even before experimental results on the millisecond timescale have become available (Markram et al., 1997; Zhang et al., 1998; Debanne et al., 1998; Bi & Poo, 1998; Feldman, 2000).

5.3 Structural Stability of Output Rates. We have seen that synaptic plasticity with an appropriate combination of parameters pushes the output rate $\bar{\lambda}^{\text{out}}$ toward a fixed point. To explicitly show the existence and stability of the fixed point, we had to require that the factor C defined by equation 3.5 is independent of the synapse index j . The standard example that we have used throughout the article is an input scenario that has the same mean rate at each synapse ($\bar{\lambda}_i^{\text{in}} = \bar{\lambda}^{\text{in}}$ for all i) and is translation invariant $C_{ij} = C_{|i-j|}$. In this case, rate normalization is equivalent to a stable fixed point of the mean weight $N^{-1} \sum_i w_i$ (see section 3.2).

We may wonder what happens if equation 3.5 is not exactly true but holds only approximately. To answer this question, let us write the linear model of equation 4.11 in vector notation $d\mathbf{w}/dt = k_1 \mathbf{e}_1 + M\mathbf{w}$ with some matrix M and the unit vector $\mathbf{e}_1 = N^{-1/2} (1, 1, \dots, 1)^T$. A stable fixed point of the mean weight implies that $\mathbf{w} = \mathbf{e}_1$ is an eigenvector of M with negative eigenvalue. If such a negative eigenvalue $\lambda^M < 0$ exists under the condition 3.5, then, due to (local) continuity of solutions of the differential equation 4.11 in dependence on parameters, the eigenvalue will stay negative in some regime where equation 3.5 holds only approximately. The eigenvector corresponding to this negative eigenvalue will be close to but not identical with \mathbf{e}_1 . In other words, the output rate remains stable, but the weight vector \mathbf{w} is no longer exactly normalized.

The normalization properties of our learning rule are thus akin to, but not identical with, subtractive normalization (Miller & MacKay, 1994; Miller, 1996b). From our point of view, the basic property of such a rule is a stabilization of the output rate. The (approximate) normalization of the mean weight $N^{-1} \sum_i w_i = c$ with some constant c is a consequence. The normalization is exact if C in equation 3.5 is independent of the synapse index j , and mean input rates $\bar{\lambda}_j^{\text{in}}$ are the same at all synapses.

Finally, the learning rule 1.1 with constant coefficients is typically unstable because weights receiving strongest enhancement grow without bounds at the expense of synapses receiving less or no enhancement. Unlimited growth can be avoided by explicitly introducing upper and lower bounds for individual weights. As a consequence, most of the weights saturate at these bounds. Saturated weights no longer participate in the learning dynamics, and the stabilization arguments of sections 3 and 4 can be applied to the set of remaining weights. For certain models, it can be shown explicitly that the fixed point of the output rate remains unchanged compared to the case without bounds (Kempster et al., 1999a). We simply have to check that the fixed point lies within the parameter range allowed by the bounds on the weights.

5.4 Principal Components and Covariance. Whereasthe output rate converges to a fixed point so that the mean weight adapts to a constant value $N^{-1} \sum_i w_i = c$, some synapses may grow and others decay. It is this synapse-

specific change that leads to learning in its proper sense (in contrast to mere adaptation). Spike-based learning is dominated by the eigenvector of (Q_{ij}) with the largest eigenvalue (Kempster et al., 1999a, 1999b). Thus, learning “detects” the principal component of the matrix (Q_{ij}) defined in equation 4.12.

A few remarks are in order. First, we emphasize that for spike-time-dependent learning, a clear-cut distinction between Hebbian and anti-Hebbian rules is difficult. Since (Q_{ij}) is sensitive to the covariance of the input on the timescale of the learning window, the same rule can be considered as anti-Hebbian for one type of stimulus and as Hebbian for another one (see the appendix).

Second, in contrast to $\overline{\lambda}_i^{\text{in}}$ Oja’s rule (Oja, 1982), it is not necessary to require that the mean input $\overline{\lambda}_i^{\text{in}}$ vanishes at each synapse. Intrinsic rate normalization implies that the mean input level does not play a role. After an initial adaptation phase, neuronal plasticity has automatically “subtracted” the mean input and becomes sensitive to the covariance of the input (see also the discussion around equations 3.2 and 3.9). Thus, after convergence to the fixed point, the learning rate becomes equivalent to Sejnowski’s covariance rule (Sejnowski, 1977; Sejnowski & Tesauro, 1989; Stanton & Sejnowski, 1989).

5.5 Subthreshold Regime and Coincidence Detection. We have previously exploited the normalization properties of learning rules with negative integral of the learning window in a study of the barn owl auditory system (Gerstner, Kempster, et al., 1996, 1997). Learning led to a structured delay distribution with submillisecond time resolution, albeit the width of the learning window was in the millisecond range. Intrinsic normalization of the output rate was used to keep the postsynaptic neuron in the subthreshold regime where the neuron functions as a coincidence detector (Kempster et al., 1998). Song et al. (2000) used the same mechanism to keep the neuron in the subthreshold regime where the neuron may show large output fluctuations. In this regime, rate stabilization induces stabilization of the value of the coefficient of variation (CV) over a broad input regime.

Due to stabilization of output rates, we speculate that a coherent picture of neural signal transmission may emerge, where each neuron in a chain of several processing layers is equally active once averaged over the complete stimulus ensemble. In that case, a set of normalized input rates $\overline{\lambda}_i^{\text{in}}$ in one “layer” is transformed, on average, into a set of equal rates in the next layer. Ideally the target value of the rates to which the plasticity rule converges would be in the regime where the neuron is most sensitive: the subthreshold regime. This is the regime where signal transmission is maximal (Kempster et al., 1998). The plasticity rule could thus be an important ingredient to optimize the neuronal transmission properties (Brenner, Agam, Bialek, & de Ruyter van Steveninck, 1998; Kempster et al., 1998; Stemmler & Koch, 1999; Song et al., 2000).

Appendix

We illustrate the static-pattern scenario for the case of a learning window with negative integral. At each time step of length Δ_t , we present a pattern vector \mathbf{x}^μ chosen stochastically from some database $\{\mathbf{x}^\mu; 1 \leq \mu \leq p\}$. The input rates are

$$\lambda_i^{\text{in}}(t) = x_i^\mu \quad \text{for } \mu \Delta_t < t < (\mu + 1) \Delta_t. \tag{A.1}$$

The “static” correlation between patterns is defined to be

$$C_{ij}^{\text{static}} = \frac{1}{p} \sum_{\mu=1}^p \frac{(x_i^\mu - \overline{\lambda_i^{\text{in}}})(x_j^\mu - \overline{\lambda_j^{\text{in}}})}{\overline{\lambda_i^{\text{in}}} \cdot \overline{\lambda_j^{\text{in}}}}, \tag{A.2}$$

where $\overline{\lambda_i^{\text{in}}} := p^{-1} \sum_{\mu=1}^p x_i^\mu$ now denotes the average over all input patterns. We assume that there is no intrinsic order in the sequence of presentation of patterns. The “time-dependent” correlation in the input is then (see Figure 3)

$$C_{ij}(s) = \begin{cases} C_{ij}^{\text{static}} (1 - |s|/\Delta_t) & \text{for } -\Delta_t \leq s \leq \Delta_t \\ 0 & \text{else.} \end{cases}$$

Let us consider a postsynaptic potential $\epsilon(s) = \delta(s - \Delta_t)$, that is, the form of the potential is approximated by a delayed localized pulse, a (Dirac) delta function. As a learning window we take (see Figure 3)

$$W(s) = \begin{cases} 1/\Delta_t & \text{for } -2\Delta_t < s \leq 0 \\ -1/\Delta_t & \text{for } 0 < s \leq 3\Delta_t \end{cases}$$

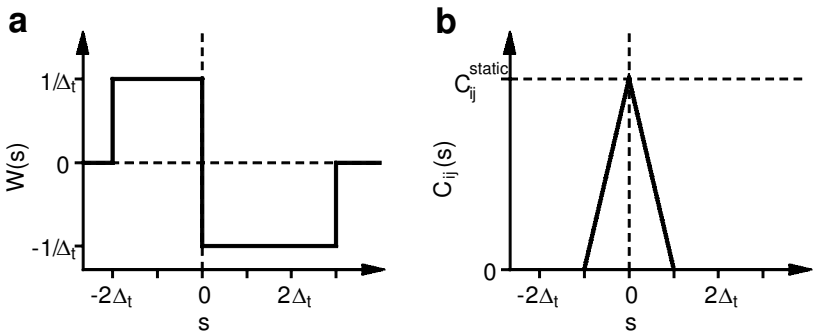


Figure 3: (a) Example of a learning window W and (b) a correlation function C_{ij} . Both are plotted as a function of the time difference s .

with $\int ds W(s) = -1 =: a_2^{\text{corr}}$. Hence, the learning rule could be classified as “anti-Hebbian.” On the other hand,

$$Q_{ij} = \int ds W(s) \int ds' \epsilon(s') C_{ij}(s + s') = C_{ij}^{\text{static}} \quad (\text{A.3})$$

gives the static correlations. Since $C_{ij} = C_{ji}$, also $Q_{ij} = Q_{ji}$ holds, and the matrix (Q_{ij}) is Hermitian. Furthermore, $\sum_{ij} C_{ij}^{\text{static}} > 0$ and $\int ds W(s) (1 - |s + \Delta_t|/\Delta_t) > 0$, so that $\sum_{ij} Q_{ij} c_i^* c_j > 0$ for all vectors $\mathbf{c} = (c_i) \neq 0$. Hence, (Q_{ij}) is a positive-definite matrix whose eigenvalues are positive (> 0). The dynamics of learning is dominated by the eigenvector of (C_{ij}^{static}) with the largest eigenvalue just as in standard Hebbian learning (Hertz et al., 1991). We note that in contrast to standard Hebbian learning, we need not impose $\overline{\lambda_i^{\text{in}}} = 0$; the spike-time-dependent learning rule automatically subtracts the mean.

Let us now assume that both $\overline{\lambda_j^{\text{in}}}$ and $Q := (\overline{\lambda^{\text{in}}})^2 N^{-1} \sum_i Q_{ij}$ are independent of j (see also section 4.5). We want to compare the value of Q with that of $\beta = N^{-1} \overline{\lambda^{\text{in}}} \int_{-\infty}^0 ds W(s) \epsilon(-s)$ in equation 4.17. We find

$$\frac{Q}{\beta} = N \overline{\lambda^{\text{in}}} \Delta_t \left[\frac{1}{N} \sum_{i=1}^N C_{ij}^{\text{static}} \right]. \quad (\text{A.4})$$

Let us substitute some numbers so as to estimate the order of magnitude of this fraction. The learning window has a duration (width) of, say, the order of $\Delta_t = 10$ ms. The mean rate, averaged over the whole stimulus ensemble, is about $\overline{\lambda^{\text{in}}} = 10$ Hz. A typical number of synapses is $N = 1000$. Let us assume that each input channel is correlated with 10% of the others with a value of $C_{ij}^{\text{static}} = 0.1$ and uncorrelated with the remaining 90%. Thus, $N^{-1} \sum_i C_{ij}^{\text{static}} = 0.01$. This yields $Q/\beta = 1$. Hence, β has the same order of magnitude as Q . For a lower value of the mean rate, Q/β will decrease; for a larger number of synapses, it will increase.

Acknowledgments

We thank Werner Kistler for a critical reading of a first version of the manuscript for this article. R. K. has been supported by the Deutsche Forschungsgemeinschaft (FG Hörobjekte and Ke 788/1-1).

References

Abbott, L. F., & Blum, K. I. (1996). Functional significance of long-term potentiation for sequence learning and prediction. *Cereb. Cortex*, *6*, 406–416.

- Abbott, L. F., & Munro, P. (1999). *Neural information processing systems (NIPS)*. Workshop on spike timing and synaptic plasticity, Breckenridge, CO. Available on-line at: http://www.pitt.edu/~pwm/LTP_LTD_99/.
- Artola, A., & Singer, W. (1993). Long-term depression of excitatory synaptic transmission and its relationship to long-term potentiation. *Trends Neurosci.*, *16*, 480–487.
- Bell, C. C., Han, V. Z., Sugawara, Y., & Grant, K. (1997). Synaptic plasticity in a cerebellum-like structure depends on temporal order. *Nature*, *387*, 278–281.
- Bi, G.-q., & Poo, M.-m. (1998). Synaptic modifications in cultured hippocampal neurons: Dependence on spike timing, synaptic strength, and postsynaptic cell type. *J. Neurosci.*, *18*, 10464–10472.
- Bi, G.-q., & Poo, M.-m. (1999). Distributed synaptic modification in neural networks induced by patterned stimulation. *Nature*, *401*, 792–796.
- Bi, G.-q., & Poo, M.-m. (2001). Synaptic modification by correlated activity: Hebb's postulate revisited. *Annu. Rev. Neurosci.*, *24*, 139–166.
- Bienenstock, E. L., Cooper, L. N., & Munro, P. W. (1982). Theory for the development of neuron selectivity: Orientation specificity and binocular interaction in visual cortex. *J. Neurosci.*, *2*, 32–48.
- Bliss, T. V. P., & Collingridge, G. L. (1993). A synaptic model of memory: Long-term potentiation in the hippocampus. *Nature*, *361*, 31–39.
- Brenner, N., Agam, O., Bialek, W., & de Ruyter van Steveninck, R. R. (1998). Universal statistical behavior of neuronal spike trains. *Phys. Rev. Lett.*, *81*(18), 4000–4003.
- Brown, T. H., & Chattarji, S. (1994). Hebbian synaptic plasticity: Evolution of the contemporary concept. In E. Domany, J. L. van Hemmen, & K. Schulten (Eds.), *Models of neural networks II* (pp. 287–314). New York: Springer-Verlag.
- Christofi, G., Nowicky, A. V., Bolsover, S. R., & Bindman, L. J. (1993). The postsynaptic induction of nonassociative long-term depression of excitatory synaptic transmission in rat hippocampal slices. *J. Neurophysiol.*, *69*, 219–229.
- Debanne, D., Gähwiler, B. H., & Thompson, S. M. (1998). Long-term synaptic plasticity between pairs of individual CA3 pyramidal cells in rat hippocampal slice cultures. *J. Physiol.*, *507*, 237–247.
- Desai, N. S., Rutherford, L. C., & Turrigiano, G. G. (1999). Plasticity in the intrinsic excitability of cortical pyramidal neurons. *Nat. Neurosci.*, *2*, 515–520.
- Egger, V., Feldmeyer, D., & Sakmann, B. (1999). Coincidence detection and changes of synaptic efficacy in spiny stellate neurons in rat barrel cortex. *Nat. Neurosci.*, *2*, 1098–1105.
- Eurich, C. W., Pawelzik, K., Ernst, U., Cowan, J. D., & Milton, J. G. (1999). Dynamics of self-organized delay adaption. *Phys. Rev. Lett.*, *82*, 1594–1597.
- Feldman, D. E. (2000). Timing-based LTP and LTD at vertical inputs to layer II/III pyramidal cells in rat barrel cortex. *Neuron*, *27*, 45–56.
- Fetz, E. E., & Gustafsson, B. (1983). Relation between shapes of postsynaptic potentials and changes in the firing probability of cat motoneurons. *J. Physiol.*, *341*, 387–410.
- Gerstner, W. (2000). Population dynamics of spiking neurons: Fast transients, asynchronous states and locking. *Neural Computation*, *12*, 43–89.

- Gerstner, W., & Abbott, L. F. (1997). Learning navigational maps through potentiation and modulation of hippocampal place cells. *J. Comput. Neurosci.*, *4*, 79–94.
- Gerstner, W., & van Hemmen, J. L. (1992). Associative memory in a network of “spiking” neurons. *Network*, *3*, 139–164.
- Gerstner, W., van Hemmen, J. L., & Cowan, J. D. (1996). What matters in neuronal locking. *Neural Comput.*, *8*, 1689–1712.
- Gerstner, W., Kempster, R., van Hemmen, J. L., & Wagner, H. (1996). A neuronal learning rule for sub-millisecond temporal coding. *Nature*, *383*, 76–78.
- Gerstner, W., Kempster, R., van Hemmen, J. L., & Wagner, H. (1997). A developmental learning rule for coincidence tuning in the barn owl auditory system. In J. Bower (Ed.), *Computational neuroscience: Trends in research 1997* (pp. 665–669). New York: Plenum Press.
- Gerstner, W., Kempster, R., van Hemmen, J. L., & Wagner, H. (1998). Hebbian learning of pulse timing in the barn owl auditory system. In W. Maass, & C. M. Bishop (Eds.), *Pulsed neural networks* (pp. 353–377). Cambridge, MA: MIT Press.
- Gerstner, W., Ritz, R., & van Hemmen, J. L. (1993). Why spikes? Hebbian learning and retrieval of time-resolved excitation patterns. *Biol. Cybern.*, *69*, 503–515.
- Häfliger, P., Mahowald, M., & Watts, L. (1997). A spike based learning neuron in analog VLSI. In M. C. Mozer, M. I. Jordan, & T. Petsche (Eds.), *Advances in neural information processing systems*, *9* (pp. 692–698). Cambridge, MA: MIT Press.
- Hebb, D. O. (1949). *The organization of behavior*. New York: Wiley.
- Hertz, J., Krogh, A., & Palmer, R. G. (1991). *Introduction to the theory of neural computation*. Redwood City, CA: Addison-Wesley.
- Herz, A. V. M., Sulzer, B., Kühn, R., & van Hemmen, J. L. (1988). The Hebb rule: Representation of static and dynamic objects in neural nets. *Europhys. Lett.*, *7*, 663–669.
- Herz, A. V. M., Sulzer, B., Kühn, R., & van Hemmen, J. L. (1989). Hebbian learning reconsidered: Representation of static and dynamic objects in associative neural nets. *Biol. Cybern.*, *60*, 457–467.
- Kempster, R., Gerstner, W., & van Hemmen, J. L. (1999a). Hebbian learning and spiking neurons. *Phys. Rev. E*, *59*, 4498–4514.
- Kempster, R., Gerstner, W., & van Hemmen, J. L. (1999b). Spike-based compared to rate-based Hebbian learning. In M. S. Kearns, S. A. Solla, & D. A. Cohn (Eds.), *Advances in neural information processing systems*, *11* (pp. 125–131). Cambridge, MA: MIT Press.
- Kempster, R., Gerstner, W., van Hemmen, J. L., & Wagner, H. (1996). Temporal coding in the sub-millisecond range: Model of barn owl auditory pathway. In D. S. Touretzky, M. C. Mozer, & M. E. Hasselmo (Eds.), *Advances in neural information processing systems*, *8* (pp. 124–130). Cambridge, MA: MIT Press.
- Kempster, R., Gerstner, W., van Hemmen, J. L., & Wagner, H. (1998). Extracting oscillations: Neuronal coincidence detection with noisy periodic spike input. *Neural Comput.*, *10*, 1987–2017.
- Kempster, R., Leibold, C., Wagner, H., & van Hemmen, J. L. (2001). Formation of temporal-feature maps by axonal propagation of synaptic learning. *Proc. Natl. Acad. Sci. U.S.A.*, *98*, 4166–4171.

- Kistler, W. M., & van Hemmen, J. L. (2000). Modeling synaptic plasticity in conjunction with the timing of pre- and postsynaptic action potentials *Neural Comput.*, *12*, 385–405.
- Kohonen, T. (1984). *Self-organization and associative memory*. Berlin: Springer-Verlag.
- Levy, W. B., & Stewart, D. (1983). Temporal contiguity requirements for long-term associative potentiation/depression in the hippocampus. *Neurosci.*, *8*, 791–797.
- Linden, D. J. (1999). The return of the spike: Postsynaptic action potentials and the induction of LTP and LTD. *Neuron*, *22*, 661–666.
- Linsker, R. (1986). From basic network principles to neural architecture: Emergence of spatial-opponent cells. *Proc. Natl. Acad. Sci. USA*, *83*, 7508–7512.
- MacKay, D. J. C., & Miller, K. D. (1990). Analysis of Linsker's application of Hebbian rules to linear networks. *Network*, *1*, 257–297.
- Markram, H., Lübke, J., Frotscher, M., & Sakmann, B. (1997). Regulation of synaptic efficacy by coincidence of postsynaptic APs and EPSPs. *Science*, *275*, 213–215.
- Miller, K. D. (1996a). Receptive fields and maps in the visual cortex: Models of ocular dominance and orientation columns. In E. Domany, J. L. van Hemmen, & K. Schulten (Eds.), *Models of neural networks III* (pp. 55–78). New York: Springer-Verlag.
- Miller, K. D. (1996b). Synaptic economics: Competition and cooperation in correlation-based synaptic plasticity. *Neuron*, *17*, 371–374.
- Miller, K. D., & MacKay, D. J. C. (1994). The role of constraints in Hebbian learning. *Neural Comput.*, *6*, 100–126.
- Oja, E. (1982). A simplified neuron model as a principal component analyzer. *J. Math. Biol.*, *15*, 267–273.
- Paulsen, O., & Sejnowski, T. J. (2000). Natural patterns of activity and long-term synaptic plasticity. *Curr. Opin. Neurobiol.*, *10*, 172–179.
- Plesser, H. E., & Gerstner, W. (2000). Noise in integrate-and-fire models: From stochastic input to escape rates. *Neural Computation*, *12*, 367–384.
- Pockett, S., Brookes, N. H., & Bindman, L. J. (1990). Long-term depression at synapses in slices of rat hippocampus can be induced by bursts of postsynaptic activity. *Exp. Brain Res.*, *80*, 196–200.
- Poliakov, A. V., Powers, R. K., Sawczuk, A., & Binder, M. D. (1996). Effects of background noise on the response of rat and cat motoneurons to excitatory current transients. *J. Physiol.*, *495*, 143–157.
- Roberts, P. D. (1999). Computational consequences of temporally asymmetric learning rules: I. Differential Hebbian learning. *J. Comput. Neurosci.*, *7*, 235–246.
- Ruf, B., & Schmitt, M. (1997). Unsupervised learning in networks of spiking neurons using temporal coding. In W. Gerstner, A. Germond, M. Hasler, & J.-D. Nicoud (Eds.), *Proc. 7th Int. Conf. Artificial Neural Networks (ICANN'97)* (pp. 361–366). Heidelberg: Springer-Verlag.
- Sejnowski, T. J. (1977). Storing covariance with nonlinearly interacting neurons *J. Mathematical Biology*, *4*, 303–321.

- Sejnowski, T. J., & Tesauero, G. (1989). The Hebb rule for synaptic plasticity: Algorithms and implementations. In J. H. Byrne, & W. O. Berry (Eds.), *Neural models of plasticity: Experimental and theoretical approaches* (pp. 94–103). San Diego: Academic Press.
- Senn, W., Tsodyks, M., & Markram, H. (1997). An algorithm for synaptic modification based on exact timing of pre- and postsynaptic action potentials. In W. Gerstner, A. Germond, M. Hasler, & J.-D. Nicoud (Eds.), *Proc. 7th Int. Conf. Artificial Neural Networks (ICANN'97)* (pp. 121–126). Heidelberg: Springer-Verlag.
- Senn, W., Markram, H., & Tsodyks, M. (2001). An algorithm for modifying neurotransmitter release probability based on pre- and postsynaptic spike timing. *Neural Comput.*, *13*, 35–67.
- Shouval, H. Z., & Perrone, M. P. (1995). Post-Hebbian learning rules. In M. A. Arbib (Ed.), *The handbook of brain theory and neural networks* (pp. 745–748). Cambridge, MA: MIT Press.
- Song, S., Miller, K. D., & Abbott, L. F. (2000). Competitive Hebbian learning through spike-timing-dependent synaptic plasticity. *Nat. Neurosci.*, *3*, 919–926.
- Stanton, P. K., & Sejnowski, T. J. (1989). Associative long-term depression in the hippocampus induced by Hebbian covariance. *Nature*, *339*, 215–218.
- Stemmler, M., & Koch, C. (1999). How voltage-dependent conductances can adapt to maximize the information encoded by neurons. *Nat. Neurosci.*, *2*, 521–527.
- Turrigiano, G. G. (1999). Homeostatic plasticity in neuronal networks: The more things change, the more they stay the same. *Trends Neurosci.*, *22*, 221–227.
- Turrigiano, G. G., Leslie, K. R., Desai, N. S., Rutherford, L. C., & Nelson, S. B. (1998). Activity-dependent scaling of quantal amplitude in neocortical neurons. *Nature*, *391*, 892–896.
- Turrigiano, G. G., & Nelson, S. B. (2000). Hebb and homeostasis in neuronal plasticity. *Curr. Opin. Neurobiol.*, *10*, 358–364.
- Urban, N. N., & Barrionuevo, G. (1996). Induction of Hebbian and non-Hebbian mossy fiber long-term potentiation by distinct patterns of high-frequency stimulation. *J. Neurosci.*, *16*, 4293–4299.
- van Hemmen, J. L. (2000). Theory of synaptic plasticity. In F. Moss & S. Gielen (Eds.), *Handbook of biological physics, Vol. 4: Neuro-informatics, neural modelling* (pp. 749–801). Amsterdam: Elsevier.
- van Hemmen, J. L., Gerstner, W., Herz, A., Kühn, R., Sulzer, B., & Vaas, M. (1990). Encoding and decoding of patterns which are correlated in space and time. In G. Dorffner (Ed.), *Konnektionismus in artificial intelligence und kognitionsforschung* (pp. 153–162). Berlin: Springer-Verlag.
- Volgushev, M., Voronin, L. L., Chistiakova, M., & Singer, W. (1994). Induction of LTP and LTD in visual cortex neurons by intracellular tetanization. *NeuroReport*, *5*, 2069–2072.
- von der Malsburg, C. (1973). Self-organization of orientation sensitive cells in the striate cortex. *Kybernetik*, *14*, 85–100.

- Wimbauer, S., Gerstner, W., & van Hemmen, J. L. (1994). Emergence of spatio-temporal receptive fields and its application to motion detection. *Biol. Cybern.*, *72*, 81–92.
- Wiskott, L., & Sejnowski, T. (1998). Constrained optimization for neural map formation: A unifying framework for weight growth and normalization. *Neural Comput.*, *10*, 671–716.
- Xie, X., & Seung, S. H. (2000). Spike-based learning rules and stabilization of persistent neural activity. In S. A. Solla, T. K. Leen, & K.-R. Müller (Eds.), *Advances in neural information processing systems*, *12* (pp. 199–205). Cambridge, MA: MIT Press.
- Zhang, L. I., Tao, H. W., Holt, C. E., Harris, W. A., & Poo, M.-m. (1998). A critical window for cooperation and competition among developing retinotectal synapses. *Nature*, *395*, 37–44.

Received January 4, 2000; accepted March 1, 2001.