

A Biologically Supported Error-Correcting Learning Rule

Peter J. B. Hancock

Leslie S. Smith

William A. Phillips

Centre for Cognitive and Computational Neuroscience,

Departments of Computing Science and Psychology,

University of Stirling, Stirling, Scotland FK9 4LA

We show that a form of synaptic plasticity recently discovered in slices of the rat visual cortex (Artola *et al.* 1990) can support an error-correcting learning rule. The rule increases weights when both pre- and postsynaptic units are highly active, and decreases them when pre-synaptic activity is high and postsynaptic activation is less than the threshold for weight increment but greater than a lower threshold. We show that this rule corrects false positive outputs in feedforward associative memory, that in an appropriate opponent-unit architecture it corrects misses, and that it performs better than the optimal Hebbian learning rule reported by Willshaw and Dayan (1990).

1 Introduction

Learning rules that correct errors are most often used in cognitive simulations and in the technological applications of neural nets. The Delta rule (Widrow and Hoff 1960) is a typical example. Three terms are required to specify the weight change: presynaptic activity, the postsynaptic activity produced by the net, and the postsynaptic activity specified by the training signal. Performance improves gradually with repeated presentation of the whole training set. There is psychological evidence for such a rule (e.g., Sutton and Barto 1981), but no biological evidence has yet been presented for a rule of this kind. Learning rules based on biological evidence typically use just two terms to specify weight change: presynaptic activity and postsynaptic activity. They do not require multiple presentations of the training set to reach their optimum performance. The many forms of this kind of learning are collectively called Hebbian rules. It is well established that the computational power of error-correcting rules exceeds that of the Hebbian rules.

Recently Artola *et al.* (1990) reported a new form of synaptic plasticity in slices of adult rat visual cortex. They show that tetanic presynaptic input produces long-term potentiation (LTP) if postsynaptic depolarization exceeds a high threshold, and long-term depression (LTD) if it does not

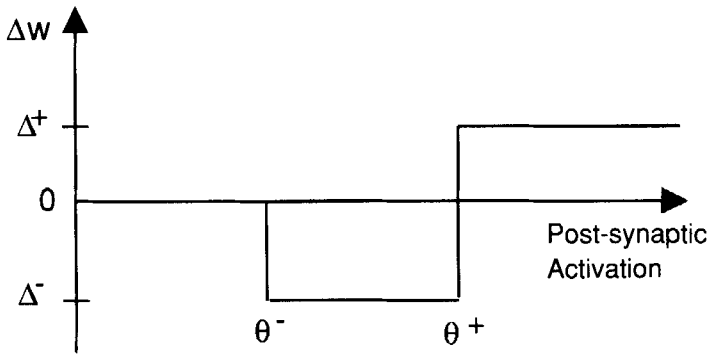


Figure 1: Simple version of the ABS rule, showing weight change for a synapse from an active presynaptic unit.

exceed the high threshold but does exceed a lower threshold. The high threshold is related to NMDA receptor-gated conductances. At first sight this seems to be just another Hebbian rule, but it is unusual because LTD occurs when the postsynaptic unit is moderately active but not when it is less active. This nonmonotonic relationship of weight change to postsynaptic activation is the critical difference.

A simple form of this rule is shown in Figure 1. We shall refer to it as the ABS rule. It resembles the proposal of Bienenstock *et al.* (1982), but it does not use the time averages of unit activity to specify weight change thresholds.

To see a possible rationale for this rule consider the development of a feedforward associative net learning a random set of pairs of binary patterns. The net consists of M input units fully connected to N output units. These output units compute a weighted sum of their inputs (including the training signal) and give a binary output determined by whether the activation is above or below their output threshold. Initially, all the weights are assumed to be sufficiently small that the only output units to be active are those driven by the training signal. We assume that this signal reaches the threshold for weight increase. As the loading increases, some of the units that should, according to the training signal, be OFF start to become active. This activity triggers the weight decrement: the rule thus reduces specifically the weights that are causing problems. This is a simple form of error correction. A few high weights from active input units to inactive output units can be tolerated, and indeed should be because of the other patterns that have been learned. Reducing all

such weights, as a simple two term rule would, is likely to lead to other errors.

Here we begin the computational study of this learning rule, and compare it with the Hebbian rule that Willshaw and Dayan (1990) have shown to be optimal of that class. They demonstrate the requirement for decreases in synaptic efficacy that on average match the increases. The optimal rule is the covariance rule (Sejnowski 1977), which they call Rule 1. Two simpler cases (Rules 2 and 3) are shown to give good but slightly less than optimal results. There is biological evidence for both of the simpler rules (Rauschecker and Singer 1979; Stanton and Sejnowski 1989).

2 ABS Rule Definition

This study of the ABS rule is designed for direct comparison with the results of Willshaw and Dayan (1990). They considered the storage in a single-layer feedforward associative net of Ω pattern pairs, each consisting of an input vector $A(\omega)$ and an output vector $B(\omega)$. The components of each $A(\omega)$ are set to 1 with a probability of s , and to a low value c with probability $(1 - s)$ (we are substituting s for their p to avoid confusion with probability p later). Part of their conclusion is that the value of c is not important, given appropriate output thresholds and their rules, so we always set it to 0. Components of a $B(\omega)$ are set to 1 with a probability of r and 0 with a probability of $1 - r$. The activation of an output unit, X_j , is given by the weighted sum of its inputs:

$$X_j = \sum_{i=1}^M A_i(\omega) W_{ij}$$

If the activation is above the unit's threshold θ_j , its output O_j is set to 1, otherwise to the low value c (0):

$$O_j = \begin{cases} 1 & \text{if } X_j > \theta_j \\ 0 & \text{otherwise} \end{cases}$$

The simple form of the learning rule shown in Figure 1 may be defined by

$$\Delta W_{ij} = \begin{cases} \Delta^+ & \text{if } X_j \geq \theta^+ \text{ and } A_i(\omega) = 1 \\ \Delta^- & \text{if } \theta^- < X_j < \theta^+ \text{ and } A_i(\omega) = 1 \\ 0 & \text{otherwise} \end{cases}$$

We do not need a specific value for θ^+ in our simulations. We assume that the target output signal is strong enough to drive units into weight increment and that the signals from the adaptive weights are not.¹ Here

¹Artola *et al.* show that, under bicuculline disinhibition, the internal signals from the adaptive weights can drive the cell sufficiently to cause weight increment. Since

we also assume that the training inputs consist of binary signals. With these two assumptions we can reformulate the ABS rule as follows:

$$\begin{aligned}\Delta W_{ij} = & \Delta^+ \text{ if } B_j(\omega) = 1 \text{ and } A_i(\omega) = 1 \\ & \Delta^- \text{ if } B_j(\omega) = 0 \text{ and } X_j > \theta^- \text{ and } A_i(\omega) = 1 \\ & 0 \text{ otherwise}\end{aligned}$$

These two forms of the ABS rule are equivalent if each unit is seen as having two different inputs, the modifiable connections of the associative memory and the training signal with which the associations are being made:

$$X_j = \sum_{i=1}^M A_i(\omega) W_{ij} + B_j(\omega) d$$

Here d is the strength of the training signal, set such that only it can drive the unit sufficiently to reach the weight increment threshold.

This specification of the rule allows weights to change between being positive and being negative, which is not biologically plausible (Crick and Asanuma 1986). One of the implications of Willshaw and Dayan's work is that negative weights are required for optimal storage. To allow direct comparison with their rules we have allowed negative weights in the first experiments reported below. This is then corrected in the section discussing an architecture with opponent units.

Note that $\theta^- < \theta_j$, in order to prevent falsely high outputs. The unit is active above θ^- , but not sufficiently active to be counted as ON in binary terms. However, if the difference is too large, the rule resembles the simple binary rule and there is a danger of overcorrection.

3 The Effects of Error Correction

Hebbian rules with binary signals lead to a distribution of activation levels after learning illustrated by Figure 2a. The overlapping tails of the desired high and desired low distribution are where the errors occur. The ABS rule is able to cut the tail off the high end of the desired low distribution, while the full three-term Delta rule is able to correct errors in both directions, (Fig. 2b, c). Obviously there comes a point where the Delta rule will also fail, but it occurs at higher loadings than for two-term rules. Note that we are using threshold logic units, which may not be biologically plausible. However it is clear that no form of output function could prevent errors if the two distributions overlap.

With additional circuitry the ABS rule is also able to correct misses. The requirement is to replace the single output units with mutually

this would lead to a runaway self-association, with strong weights getting stronger, we assume that the threshold was reached because of the disinhibition, and that normally other inputs would also be required.

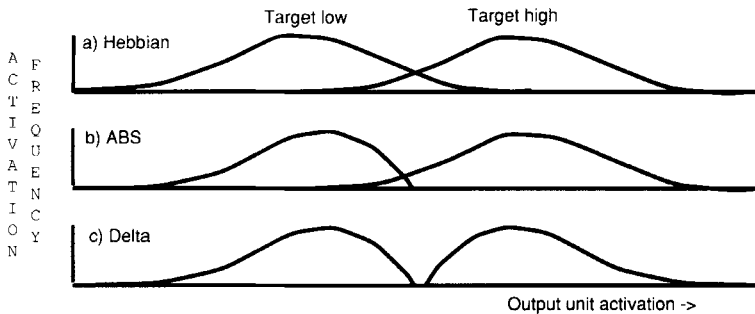


Figure 2: Idealized activation frequency distributions after learning for a single output unit, plotted separately for when it should have low and high outputs. (a) Result of simple binary two term rules: the region of overlap indicates that errors will be made, wherever the threshold is put. (b) The ABS rule corrects false positive errors, reducing the region of overlap. (c) The Delta rule corrects errors in both directions.

inhibitory opponent pairs. This may be regarded as a simplification of the local inhibition that is common in cortex. Whenever a unit is trained to be ON, its opponent is trained to be OFF, and vice versa. False positive outputs are corrected as before. The way that misses are corrected can be seen by considering why a unit that is below its output threshold has too little activation. Part of the reason is that its own weights are too low, but it is also being inhibited by its opponent cell. By symmetry, this unit is responding too strongly and will reach the weight decrement threshold. Its activity will be decreased, reducing the inhibition and allowing the other unit to give a higher output.

The simple two-term rules learn in a single pass: unless weights are limited in some way further presentations of the training set will not affect the result, only the size of the weights and activation. As with the Delta rule, the ABS rule gives improved performance with additional presentations.

4 Simulation Experiments

4.1 Single Unit Architecture. The performance of the ABS rule was tested by repeating the experiments of Willshaw and Dayan (1990), who measured signal-to-noise ratios and the number of bit errors for a feedforward associative net. Their computations of signal to noise ratio assume

the distributions are gaussian. The ABS rule distribution (Fig. 2b) is not gaussian, so that the figures produced are not directly comparable. We therefore report only the actual numbers of errors produced, since minimizing this is more important.

The bit errors were counted by two methods. Initially we used the threshold set by Willshaw and Dayan's method, which is designed for gaussian distributions. As might be expected from Figure 2, the ideal threshold for the ABS rule is rather lower than for the simpler Hebbian rules. An optimal threshold can be found by searching through the actual responses of each unit in the region where desired low and desired high outputs overlap to minimize the number that are wrong. This procedure produced significantly better results. We have not yet looked for a method of setting something like the optimal threshold for each unit without recourse to such serial search procedures.

The optimal Hebbian rules of Willshaw and Dayan (1990) specify the sizes of the weight change parameters for given bit probabilities. The ABS rule decrements the weight only when an error occurs, so that if Willshaw and Dayan's conclusion that the expected value of the weight should be zero still holds, the size of the decrement has to be larger than is given by their binary homosynaptic depression Rule 3. This rule gives a weight increment of $1 - r$ and a decrement of r (r is the output bit probability). We therefore fixed the increment size at $1 - r$, and experimented with a range of values for the decrement size. The initial value of all the weights was zero, there being no need here for the symmetry breaking required by some other methods. The results are shown in Table 1.

A number of things are apparent from Table 1. The absolute level of performance is good, and improves as bit probability decreases. It learns 200 patterns when the bit probability is 0.1 with on average only 0.05 bits in error out of 20, so at least 190 of the 200 output patterns will be completely correct. As predicted, for both bit probabilities, the optimal size of Δ^- is larger than the value specified by Willshaw and Dayan's Rule 3. Near the optimum, the precise value of Δ^- is not critical. In both cases the average value of the weights is near zero at optimal performance. These results were used to set the sizes of the weight changes to their optimal value in the following experiments.

We next compared the ABS rule with the optimal Hebbian rule (Rule 1) of Willshaw and Dayan (W&D) over a range of bit probabilities. The results are given in Table 2. The ABS rule does better at all bit probabilities and in contrast to normal Hebbian rules, its performance improves with training. However, there is little room for improvement at low bit probabilities and the limit is quickly reached.

4.2 Opponent Architecture. In the simple architecture of the preceding experiments the ABS rule corrects false positives. In an architecture with twice as many output units arranged in mutually inhibitory pairs it also corrects misses. The internal activations for each unit are calculated

Table 1: Results from a net with 512 input units and 20 output units, with 200 patterns, averaged over 10 runs, with 5 training cycles. Weight increment is Δ^+ , weight decrement Δ^- . Avg weight is the average of all weights to all 20 output units. Bit errs is the average number of errors per 20-bit pattern, counted using the threshold used by W&D. Min errs is the number of bit errors given by an optimal threshold for each unit.

Bit prob $s,r = 0.1, \Delta^+ = 0.9$				Bit prob $s,r = 0.5, \Delta^+ = 0.5$		
Δ^-	Bit errs	Min errs	Avg weight	Bit errs	Min errs	Avg weight
0.1	0.135	0.065	1.8			
0.2	0.06	0.018	0.43			
0.3	0.05	0.0125	0.19	6.38	5.97	50.3
0.4	0.05	0.0055	0.06	3.86	3.58	26.4
0.5	0.05	0.003	-0.02	1.6	1.41	9.3
0.6	0.05	0.0025	-0.03	1.04	0.85	3.26
0.7	0.05	0.002	-0.09	1.01	0.81	1.61
0.8	0.052	0.001	-0.16	1.01	0.82	1.32
0.9	0.053	0.0015	-0.19	1.05	0.86	1.09
1.0				1.1	0.94	0.82

Table 2: Results from a 512 input, 20 output net with 200 random input-output patterns, averaged over 10 runs with different pattern sets. Bit errors refers to the average number of errors per pattern, counted using the threshold used by W&D. Min errors is the number of bit errors given by an optimal threshold for each unit.

s,r	W & D Rule 1		ABS 5 epochs		ABS 10 epochs		ABS 20 epochs	
	Bit errors	Min errors	Bit errors	Min errors	Bit errors	Min errors	Bit errors	Min errors
0.5	1.07	0.89	1.03	0.86	0.77	0.50	0.71	0.34
0.4	0.97	0.82	0.82	0.63	0.64	0.29	0.61	0.13
0.3	0.72	0.56	0.54	0.32	0.48	0.12	0.47	0.044
0.2	0.35	0.25	0.24	0.06	0.24	0.015	0.23	0.005
0.1	0.08	0.027	0.05	0.004	0.05	0.004	0.05	0.004
0.05	0.03	0.003	0.02	0.0	0.02	0.0	0.02	0.0

as before, the units then inhibit each other by subtracting some fraction κ of the opponent unit's activation:

$$X_j = \sum_{i=1}^M A_i(\omega)W_{ij} - \kappa \sum_{i=1}^M A_i(\omega)W_{ik}$$

During training, for each unit where the target is 1, its opponent unit is set to 0, and vice versa. The weight change procedure for each unit is the same as for the single-sided architecture.

We are not suggesting that such an orderly arrangement of pairs of units is biologically plausible. This design is a simplification that matches the assumption of binary training signals. However, local mutual inhibition is widespread in the cortex and a more realistic simulation might contain a layer of units such as that suggested by von der Malsburg (1973). Here we only wish to demonstrate the possibilities of the learning rule and have kept the architecture as simple as possible.

The opponent architecture also allows the problem of negative weights to be addressed. Effectively, we are simply splitting each weight in two, and putting the inhibitory part on a separate unit. For this to work requires only that the weight decrement threshold θ^- be above zero. The value is not critical, since the weights and activations are automatically adjusted appropriately. A value of 50 proved satisfactory for the weight change parameters in use here.

Results are given in Table 3. This system can learn 200 patterns without errors, though convergence to this accuracy is quite slow, requiring

Table 3: Results from a net with 512 input units and 20×2 output units trained with 200 random input-output patterns, for a variety of parameters. In all cases Δ^+ is 0.02, and there are 30 training cycles.

Bit probability	Δ^-	κ	Min bit errors per pattern
0.5	0.1	0.5	1.28
0.5	0.1	0.8	0.03
0.5	0.1	0.9	0
0.5	0.1	1.0	0.37
0.5	0.15	0.9	0
0.3	0.1	0.9	0
0.2	0.1	0.9	0.007
0.1	0.1	0.9	0.048

30 or 40 training epochs. Performance is distinctly better than the single-sided architecture, which still makes about 0.1 errors per pattern after 40 epochs.

The system is sensitive to the value of κ (the strength of the mutual inhibition), with 1 giving distinctly worse performance than slightly lower values. The value of Δ^- is less critical, a good value being five times the size of Δ^+ . Performance tails off as the bit probability decreases: precisely the opposite of the simpler architecture, although with bit probabilities as low as 0.2 this system still does better than the optimal Hebbian rule. The reason for the effect of bit probability is clear: if one unit is ON with a probability of 0.1, then its opponent is ON with a probability of 0.9. The required values of Δ^+ and Δ^- are very different for the two opponent units. Choosing appropriate values does allow all the patterns to be learned. Although a mechanism for adjusting weight change sizes to suit the measured bit probability is possible, we prefer a more biological solution with small groups of mutually inhibitory units (like a winner-takes-all cluster), each of which responds approximately equally often.

5 Stability

An important question concerning any learning rule is its stability and convergence, both in terms of errors and synaptic weights. Consider an individual weight W_{ij} . It will be incremented for those patterns where $A_i(\omega) = B_j(\omega) = 1$. Assuming input and output bit probabilities s and r are equal, weight increment would be expected on Ω/r^2 patterns. Weight decrement is expected for some of the cases where $A_i(\omega) = 1$ and $B_j(\omega) = 0$, specifically those when the unit activation $X_j > \theta^-$. Zero weight change is achieved if

$$\frac{\Omega}{r^2} \Delta^+ = \frac{\Omega}{r(1-r)} \Delta^- p [X_j > \theta^- | B_j(\omega) = 0]$$

This rearranges to give

$$\frac{\Delta^+}{\Delta^-} = \frac{r}{(1-r)} p [X_j > \theta^- | B_j(\omega) = 0]$$

That this is at least moderately stable may be seen by considering the situation where the value of Δ^+ is too high. Weights will tend to increase, leading to an increased probability of exceeding θ^- and provoking a weight decrement. Conversely, an overlarge value for Δ^- will reduce the probability of exceeding θ^- , allowing the weights to build up. Exceeding θ^- does not necessarily imply registering an error, provided there is a gap between θ^- and the binary output threshold θ_j . As with many systems a suitably pathological input sequence will break it; in practice with the runs reported here we saw just an occasional single bit error in the epochs

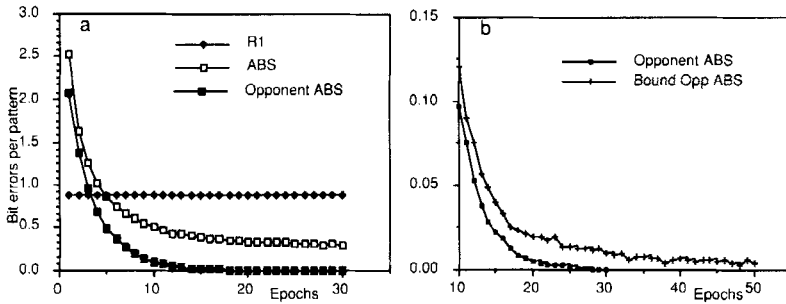


Figure 3: (a) Course of learning for optimal Hebbian and simple and opponent ABS rules, 512 input units, 20 output units/pairs, 200 random patterns of bit probability 0.5, average of 10 experiments. Both ABS rules had $\Delta^+ = 0.02$ and $\Delta^- = 0.1$. θ^- was 0 for simple, 50 for opponent ABS. Latter had weights constrained to be positive and initial weights 4.0. (b) Expanded scale, showing effects of adding an upper weight bound of 6 to opponent ABS. Other details as (a).

following initial convergence. Overall convergence is smooth, as shown by Figure 3.

As the error rate is reduced, the ratio of Δ^+/Δ^- required for weight stability will also decrease. Work in progress with adjustment of the ratio indicates that performance is indeed improved.

As the ABS rule contains a purely Hebbian increment component, it is clear that there is no upper limit on the weights: a single input repeatedly applied would cause all the active weights to grow indefinitely. Although frequently ignored in simulation work, any real synapse (in brains or silicon) will clearly have an upper limit on its strength. So the behavior of the ABS rule with an upper weight bound is important. It was checked by simply clipping any weight that exceeded a limit. This was arbitrarily set at 6, an intentionally very tight constraint given that the weights start at 4, and that some normally reach around 15 while learning 200 patterns in 30 epochs. As would be expected, the performance deteriorated noticeably, but the system still converges well and approaches zero errors (Fig. 3b). The weights were followed beyond 50 epochs and do not change significantly. In practice, therefore, the ABS rule is stable and tolerant of constraints.

6 Discussion

We have shown that a learning rule based on the form of synaptic plasticity reported by Artola *et al.* (1990) can correct false positives and misses.

It can learn more random paired associates than the optimal classical Hebbian rule, and its performance continues to improve with repeated presentations of the training set.

In essence the rule assumes that during training the required outputs are signaled by distinctively high levels of postsynaptic activation. Lower levels of postsynaptic activation within a specified range can thus be treated as false positives and the weights from active input lines reduced. This entails two further assumptions. (1) The maximum sizes of the weights produced by this rule must be limited so that they cannot produce levels of postsynaptic activation that mimic the training signals. (2) The functionally effective level of postsynaptic activation must be less than the level required for weight increment, otherwise the learning would not be effective at test. Both assumptions are biologically plausible because the weight must be limited, and it is known that neural output activity can be functionally effective at activation levels well below the NMDA threshold.

In our formulation of the rule for the opponent architecture we calculate the internal activations according to the modifiable weights, then use the output signal to decide the weight change. In reality, one of the units will be being driven hard on by the output signal. This should thoroughly inhibit the other unit, which would therefore never reach the decrement threshold, causing the rule to revert to a simple Hebbian. A plausible solution to this problem is given by the possibility of dendritic processing. The patch of the dendrite receiving the modifiable input may then reach decrement threshold, while the inhibition prevents the cell from firing.

A simple biological mechanism could provide the predicted change in ratio of weight increment to decrement as error rate declines. The size of the decrement could be controlled by the concentration of some substance, an enzyme perhaps, at the synapse. Frequent weight decrement events would use up the stock of enzyme, reducing the size of the change. A low error rate would result in occasional, larger decrements.

Further work on the ABS rule to be reported elsewhere (Hancock *et al.* 1991) shows that it compares favorably with the classical perceptron learning rule (PLR) in the early stages of learning. The PLR does not perform particularly well on the first pass of training data and there has traditionally been a divide between single-pass Hebb-like rules and multipass error-correcting rules. The ABS rule thus raises the possibility of obtaining the benefits of both, with a relatively good performance on a single pass, but continuing to improve with further training. The rule also works well in autoassociative architectures.

Important unresolved issues on which we are currently working include the extension of the rule to nonbinary signals, and its role in multilayer architectures when combined with other biologically supported learning rules.

Acknowledgments

This work was funded by the SERC of the UK, and by the BRAIN initiative of the EEC. We are very grateful to Wolf Singer and Alain Artola for helpful discussions and to Peter Dayan for help with setting up the simulations. Roland Baddeley and Peter Cahusac made helpful comments on earlier drafts of this paper.

References

- Artola, A., Bröcher, S., and Singer, W. 1990. Different voltage-dependent thresholds for the induction of long-term depression and long-term potentiation in slices of the rat visual cortex. *Nature (London)* **347**, 69–72.
- Bienenstock, E. L., Cooper, L. N., and Munro, P. W. 1982. Theory for the development of neuron selectivity: Orientation specificity and binocular interaction in visual cortex. *J. Neurosci.* **2**, 32–48.
- Crick, F. H. C., and Asanuma, C. 1986. Certain aspects of the anatomy and physiology of the cerebral cortex. In *Parallel Distributed Processing*, J. L. McClelland and D. E. Rumelhart, eds., Vol. 2, pp. 333–371. Bradford Books, MIT Press.
- Hancock, P. J. B., Smith, L. S., and Phillips, W. A. 1991. Error correcting capabilities of a recently discovered form of cortical synaptic plasticity. In preparation.
- Rauschecker, J. P., and Singer, W. 1979. Changes in the circuitry of the kitten's visual cortex are gated by post-synaptic activity. *Nature (London)* **280**, 58–60.
- Sejnowski, T. J. 1977. Storing covariance with nonlinearly interacting neurons. *J. Math Biol.* **4**, 303–321.
- Stanton, P., and Sejnowski, T. J. 1989. Associative long-term depression in the hippocampus: Induction of synaptic plasticity by Hebbian covariance. *Nature (London)* **339**, 215–218.
- Sutton, R. S., and Barto, A. G. 1981. Toward a modern theory of adaptive networks: Expectation and prediction. *Psychol. Rev.* **88-2**, 135–170.
- von der Malsburg, C. 1973. Self-organization of orientation sensitive cells in the striate cortex. *Kybernetik* **14**, 85–100.
- Widrow, B., and Hoff, M. E. 1960. Adaptive switching circuits. IRE WESCON Convention Record, New York: IRE, 96–104.
- Willshaw, D. J., and Dayan, P. 1990. Optimal plasticity from matrix memories: What goes up must come down. *Neural Comp.* **2**, 85–93.