

The Upstart Algorithm: A Method for Constructing and Training Feedforward Neural Networks

Marcus Frean

*Department of Physics and Centre for Cognitive Science, Edinburgh University,
The Kings Buildings, Mayfield Road, Edinburgh, Scotland*

A general method for building and training multilayer perceptrons composed of linear threshold units is proposed. A simple recursive rule is used to build the structure of the network by adding units as they are needed, while a modified perceptron algorithm is used to learn the connection strengths. Convergence to zero errors is guaranteed for any boolean classification on patterns of binary variables. Simulations suggest that this method is efficient in terms of the numbers of units constructed, and the networks it builds can generalize over patterns not in the training set.

1 Introduction

The perceptron learning algorithm (Rosenblatt 1962) offers a powerful but restricted method for learning binary classifications. All classifications that can in theory be learned by the perceptron architecture *will* be learned; however, the number of classifications it can learn is only a tiny subset (*linearly separable* patterns) of those that are possible (Minsky and Papert 1969). To perform any arbitrary classification successfully, "hidden" units and/or feedback between units are required. The problem is to train such networks, and recently quite powerful methods have become available, most notably "backpropagation" in its various forms (e.g., Rumelhart, et al. 1986). However, because many of these methods are based on hill climbing, which has the perennial problem of becoming stuck in local optima, they cannot guarantee that the classification will be learned. Another problem is that *a priori* no realistic estimate can be made of the number of hidden units that are required. Recently, methods such as the tiling algorithm (Mezard and Nadal 1989) and others (Gallant 1986a; Nadal 1989) have been proposed that get around both these problems. In these, the hidden units are constructed in layers one by one as they are needed. By showing that at least one unit in each layer makes fewer errors than a corresponding unit in the previous layer, eventual convergence to zero errors is guaranteed.

The method described here also constructs units as it goes, but in a simple and quite different way. Instead of building layers from the input outward until convergence, new units are interpolated between the input layer and the output. The role of these units is to correct mistakes made by the output unit.

2 Basics

Suppose we are given a binary classification to be learned. Each input pattern of N binary values has an associated target output that the network must learn to produce. The units are all *linear threshold* units connected by variable weights to the inputs, with output o given by

$$o = \begin{cases} 1 & \text{if } \phi \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

where

$$\phi = \sum_{i=0}^N W_i \xi_i$$

The W s are the weights, and ξ_i is the value of the i th input in the given pattern. The necessary threshold or "bias" is included by having an extra input that is set to 1 for all the input patterns. On presentation of pattern μ , perceptron learning (Rosenblatt 1962) alters the weights if the target t^μ differs from the output:

$$\Delta W_i^\mu = (t^\mu - o^\mu) \xi_i^\mu \quad (2.1)$$

Using this method any linearly separable class will be learned, but when the patterns are not linearly separable the values of the weights never stabilize. However, a simple extension called the *pocket algorithm* (Gallant 1986b) suffices to make the system well behaved. This consists of running a perceptron exactly as above with a random presentation of patterns, but also keeping a copy of the set of weights that has lasted longest without being changed. This set of weights will give the minimum possible number of errors with a probability approaching unity as the training time increases. That is, if a solution giving say p or fewer errors exists then the pocket algorithm can be used to find it (although unfortunately there is no bound known for the training time actually required). I make use of this algorithm to demonstrate convergence to zero errors.

3 Rationale

The basic idea is that a unit builds other units to correct its mistakes. Any unit (say Z) can make two kinds of mistake:

“wrongly ON” ($\sigma_Z^\mu = 1$, but $t_Z^\mu = 0$)

“wrongly OFF” ($\sigma_Z^\mu = 0$, but $t_Z^\mu = 1$)

Consider patterns for which Z is wrongly ON: Z could be corrected by a large negative weight from a new unit (say X), which is active only for those patterns. Likewise when Z is wrongly OFF it could be corrected by a large positive weight from another unit (say Y), which is active at the right time. Hence X and Y (also connected to the input layer by variable weights) can be trained with targets which depend on Z 's response. These new units might be called “daughters” since they are generated¹ by the established “parent” unit, Z .

Consider, for example, the targets we should assign to X , the unit whose role is to inhibit Z . We would like X to be active if Z was wrongly ON, and silent if Z was correctly ON. Similarly X should be silent if Z was wrongly OFF (to avoid further inhibition of Z). Finally, X could be silent if Z was correctly OFF, although if X is active in this case, the effect is merely to reinforce Z 's response when it was *already correct*. This does not itself cause an error, so in practice we can eliminate these patterns from X 's training set. This elimination makes the problem easier and faster to solve, but is not essential for the error-correcting property described below. Targets for Y are similarly derived. These target assignments are summarized in Figure 1.

An important point is that the “raw” output of unit Z is used to set the daughter's targets, rather than the value of Z after the daughters have exerted any effect, since this would introduce feedback.

Two useful results follow immediately from this training method, because it essentially gives daughters (X or Y) an easier problem to solve than their parent (Z). First, daughters can always make fewer errors than their parent, and second, connecting daughter to parent with the appropriate weight will always reduce the errors made by the parent.

Proof. Z 's errors are

$$e(Z) = e(Z)_{\text{ON}} + e(Z)_{\text{OFF}}$$

where $e(Z)_{\text{ON}}$ is the number of patterns for which Z is wrongly ON.

¹Note, however, that the *activity* proceeds from daughter to parent.

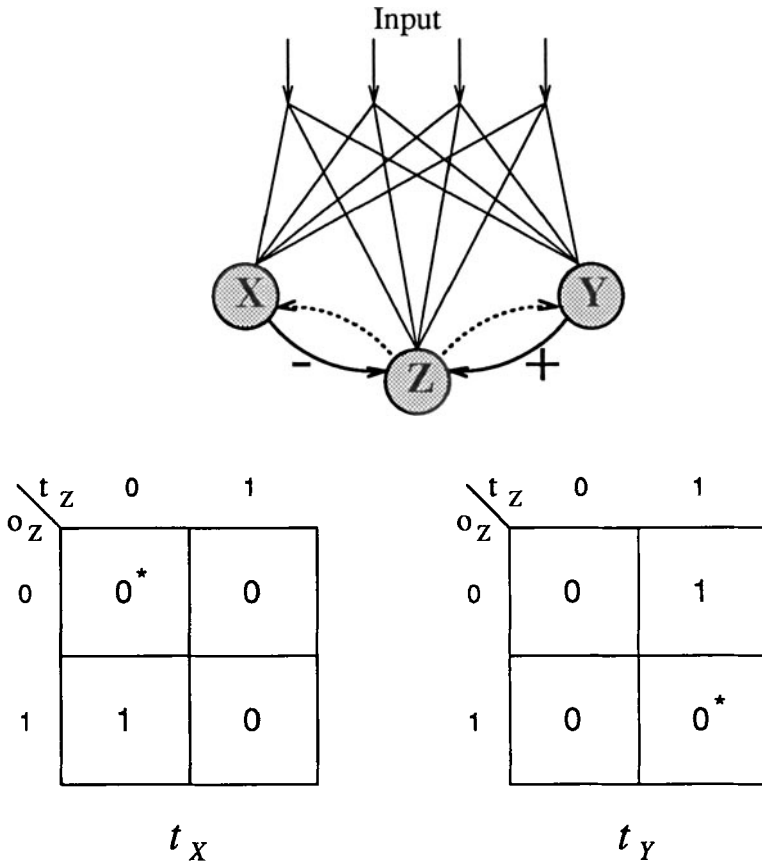


Figure 1: Correcting a parent unit: the left-hand table gives the targets, t_X , for the daughter unit X for each combination of (o_Z, t_Z) . For example, the lower left-hand entry assigns t_X to be 1 when $o_Z = 1$ and $t_Z = 0$: the “wrongly ON” case. Similarly the right-hand table gives the values of t_Y for the daughter unit Y . The dotted line represents the flow of this target information. The “starred” entries correspond to cases where the pattern could be eliminated from the daughter’s training set.

If X responded OFF to every pattern, it would make as many errors as there were patterns of target $t_X = 1$. However, X can always do better than this. In particular, it can always be ON for at least one input pattern

whose target is 1 and OFF for all other patterns. For example if the input weights are

$$W_i = 2\xi_i^\mu - 1$$

with a bias weight

$$W_0 = -\sum_j \xi_j^\mu (2\xi_j^\mu - 1)$$

only the μ th pattern turns the unit ON. Given that the pocket algorithm can find the optimal weights visited by a perceptron with any given probability, at worst X could find the above weights. Therefore

$$e(X) < e(Z)_{\text{ON}} \leq e(Z) \quad (3.1)$$

A similar argument applies to Y . It also follows that Z 's errors are reduced by X , since

$$\begin{aligned} e(Z \text{ with } X) &= e(X)_{\text{ON}} + e(X)_{\text{OFF}} + e(Z)_{\text{OFF}} \\ &= e(X) + e(Z)_{\text{OFF}} \\ &< e(Z) \end{aligned} \quad (3.2)$$

and similarly for Y on its own. When the joint action of X and Y is considered, the same result holds, that is, $e(Z \text{ with } X, Y) < e(Z) - 1$. In the next section an algorithm that uses the first of these results is described. Other possibilities are discussed in Section 6.

4 The Upstart Algorithm

Assume we already have a unit Z that sees input patterns ξ_i^μ : $i = 1, \dots, N$ and has associated targets t_Z^μ . The weights from the input layer to Z are trained to minimize the discrepancies between Z 's output and target and, once trained, these weights remain frozen. This "first" unit is actually the eventual output unit, and its targets are the classification to be learned. The following two steps are then applied recursively, generating a binary branching tree of units. Thus daughter units behave just as Z did, constructing daughter units themselves if they are required.

Step 1. If Z makes any "wrongly ON" mistakes, it builds a new unit X , using the targets given in Figure 1. Similarly if Z is ever "wrongly OFF" it builds a unit Y . Apart from the different targets, these units are trained and then frozen just as Z was.

Step 2. The outputs of X and Y are connected as inputs to Z . The weight from X is large and negative while that from Y is large and positive. The size of the weight from X [Y] needs to exceed the sum of Z 's positive [negative] input weights, which could be done explicitly or by perceptron learning.

New units are only generated if the parent makes errors, and the number of errors decreases at every branching. It follows that eventually none of the terminal daughters makes any mistakes, so neither do their parents, and neither do *their* parents, and so on. Therefore every unit in the whole tree produces its target output, including Z , the output unit. Hence the classification is learned.

5 Simulations

In all the simulations shown here the “starred” entries in Figure 1 were not included in a daughter’s training set.² To decrease training times, a fast and well behaved version of perceptron learning (Frean 1990) was used to train the weights, in preference to the pocket algorithm. While this method is not guaranteed to find the optimal weights, in practice it produces substantially fewer errors in a given time than the pocket algorithm. The weight changes given by the usual perceptron learning rule (equation 2.1) were simply multiplied by

$$\frac{T}{T_0} \exp\left(\frac{-|\phi^\mu|}{T}\right)$$

This factor decreases with $|\phi|$, which measures how “serious” the error is. The rationale behind this approach is that an error where $|\phi|$ is large is difficult to correct without causing other errors and should therefore be weighted as less significant than those where $|\phi|$ is small. The “temperature” T controls how strongly this weighting is biased to small $|\phi|$. T was reduced linearly from T_0 to zero over the entire training period. At high T the perceptron rule is recovered, but as T decreases the weights are “frozen.” Unless otherwise stated, T_0 was set to 1, and 1000 passes were made through the training set.

5.1 Parity. In this problem the output should be ON if the number of active inputs is odd, and OFF if it is even. Parity is often cited as a difficult problem for neural networks to learn, being a predicate of order N (Minsky and Papert 1969); that is, at least one hidden unit must sample all of the N inputs. It is also of interest because there is a known solution consisting of a single layer of N hidden units projecting to an output unit. It is easy to see how the upstart algorithm tackles parity (see Fig. 2). Essentially the same structure as that shown for $N = 3$ would arise for any N , although for large problems the optimal weights become much harder to find. I have tried parity for up to $N = 10$, and in all cases N units are constructed, including the output unit. In all cases except

²If the whole training set is used in every case, the number of units produced is relatively unaffected for the problems investigated here, but the training time (a combination of the time per epoch and number of epochs required to generate a comparable network) is approximately doubled for the problems discussed.

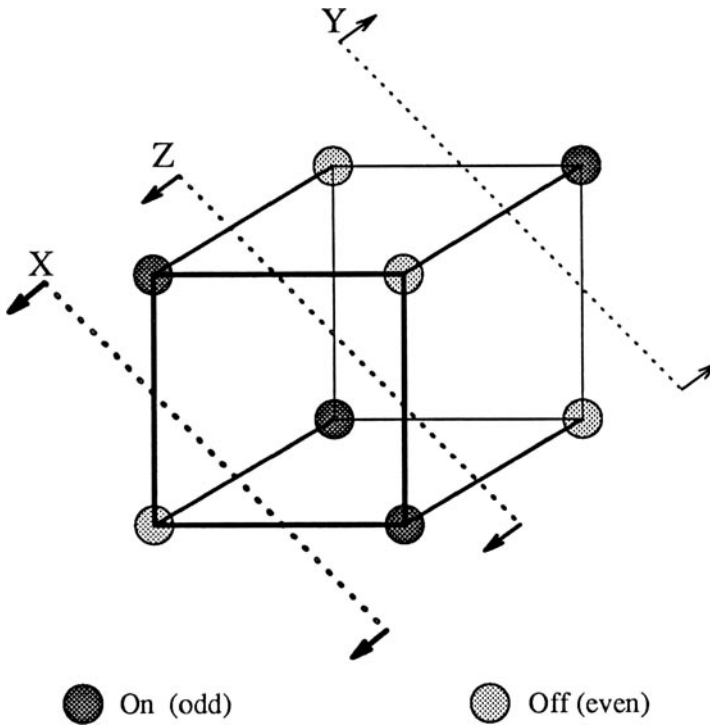


Figure 2: Solution for 3-bit parity. The output unit Z on its own can clearly make a minimum of two mistakes, when the plane defined by its weights cuts the cube as shown. X corrects the wrongly ON pattern by responding to it alone, and similarly Y corrects the wrongly OFF pattern.

$N = 10$, 1000 passes were sufficient to generate the minimal tree. For 10-bit parity, the figure was 10,000.

5.2 Random Mappings. In this problem the binary classification is defined by assigning each of the 2^N patterns its target 0 or 1 with 50% probability. Again this is a difficult problem, due to the absence of correlations and structure in the input for the network to exploit. The networks obtained for N up to 10 are summarized in Figure 3, with comparisons to the tiling algorithm.

5.3 Generalization. Neural networks are often ascribed the property of generalization: the ability to perform well on all patterns taken from

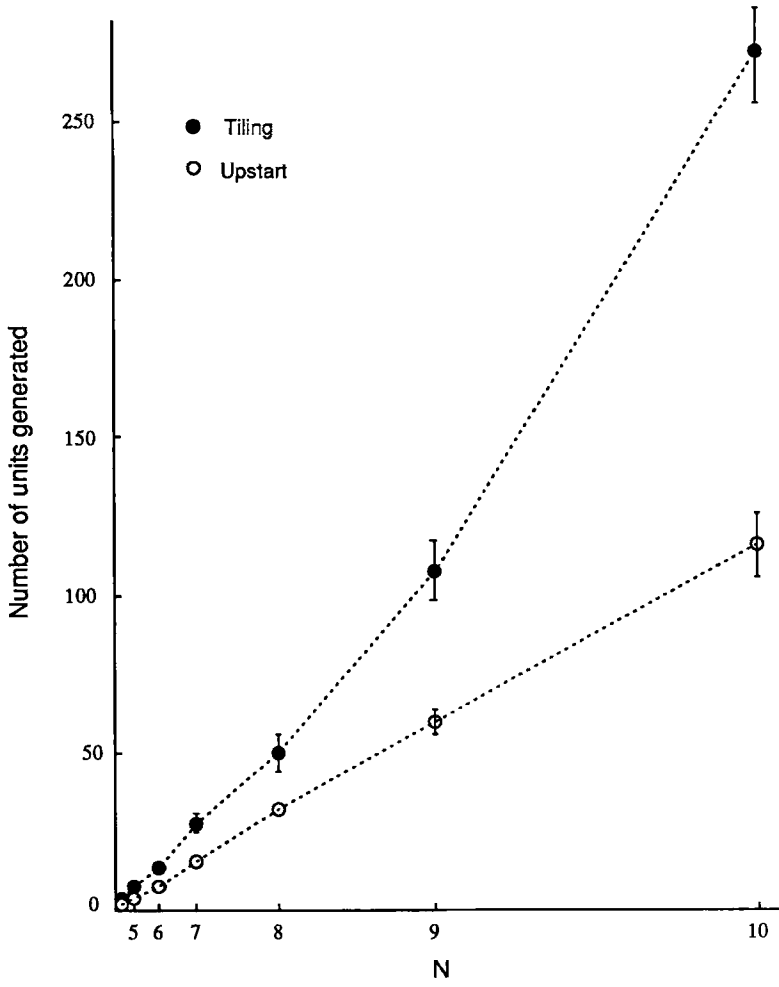


Figure 3: Number of units built vs. the number of patterns (2^N) for the random mapping problem. The slope of the upstart line is approximately $1/9$. Each point is an average of 25 runs, each on a different training set.

a given distribution after having seen only a subset of them. Several workers (Denker et al. 1987; Mezard and Nadal 1989) have looked at generalization using the "2-or-more clumps" predicate. The problem is this: given an input pattern, respond ON if the number of clumps is 2 or greater and OFF otherwise, where a "clump" is a group of adjacent³ 1's bounded on either side by 0's. As with parity, there is a solution consisting of a single hidden layer of N units that would solve the problem exactly. Following Mezard and Nadal, the patterns were generated by a Monte Carlo method (Binder 1979) such that the mean number of clumps is 1.5. I used $N = 25$ input units, with a training set of up to 600 patterns. The set used to test the resulting net's ability to generalize was a further 600 patterns. The results, with comparisons to the tiling algorithm, are summarized in Figure 4.

6 Discussion

The architecture generated by this procedure is unconventional in that it has a hierarchical tree structure. However, in the case where we choose not to eliminate any training patterns there is an equivalent structure with the same units arranged as a single hidden layer. Consider two daughters (say X, Y) and their parent (Z). With primes denoting "corrected" values, the corrected value o'_Z is always equal to $o_Z - o'_X + o'_Y$. This formulation is possible because the combinations that would disobey this never occur. For example, Y would never be correctly ON if Z was ON. Since this is true of *every* unit in the tree, the final output is simply a sum of the "raw" responses. For example,

$$\text{output} = o'_Z = o_A - o_B + o_C + \dots - o_X + o_Y + o_Z$$

Imagine the tree units disconnected from one another and placed in a single layer. A new output unit connected to this "hidden" layer can easily calculate the appropriate sum by, for example, having weights of +1 from each unit that adds to the sum and -1 from each unit that subtracts, with a bias of zero. In effect we can convert a binary tree into a single hidden layer architecture that implements the same mapping, at the expense of adding one unit and being unable to exploit pattern elimination. The algorithm for constructing a single hidden layer architecture is as follows: construct units as before, omitting step 2 (where they are connected into a feedforward tree). Then connect all the units so constructed to a new output unit. The weights can be learned by perceptron learning, or can be inferred from the tree structure: there is a sign reversal for every "X" daughter.

6.1 Extensions. The upstart method can be extended in a number of ways. First, we are not restricted to binary branching trees. Having

³Circular boundary conditions are used: input 1 is "adjacent" to input N .

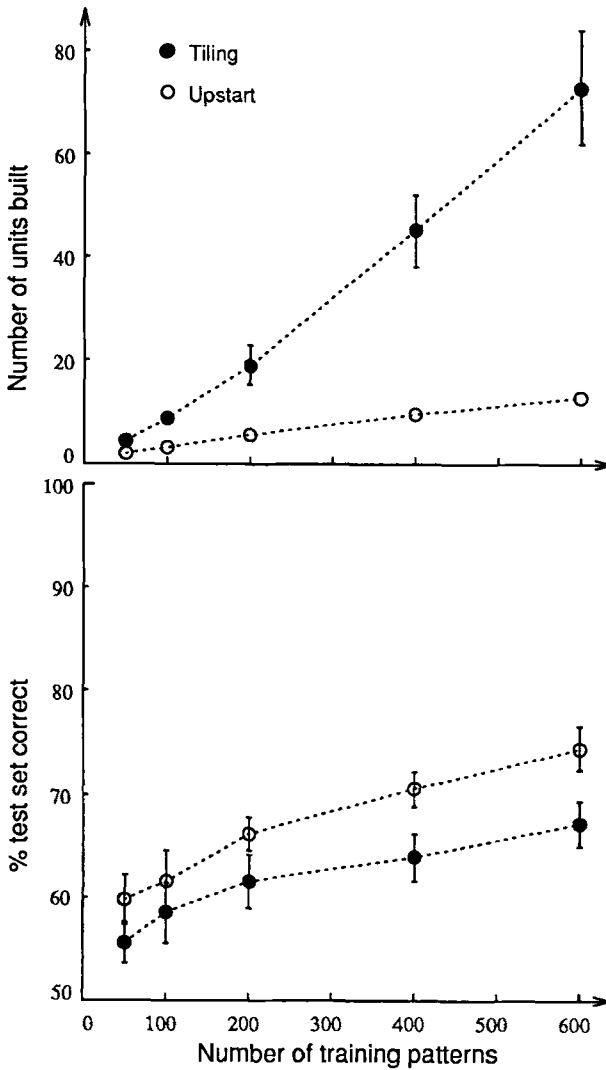


Figure 4: Performance of the method on the "2-or-more clumps" problem. The lower graph shows the percentage generalization as the size of the training set is increased. Plotted above this and with the same abscissa is the size of the corresponding network. $T_0 = 4.0$. There are 25 runs per point, each on a different set. Where not shown, error bars are smaller than the points.

trained a daughter unit and connected it to the output, a new daughter can now be trained using targets derived from the *partially corrected* output, instead of the daughter, and so on until no more mistakes are made. This would build a single hidden layer. Hybrid methods, building trees of variable breadth, will also work.

Second, these algorithms can be extended to problems involving multiple output units. A good method should build considerably fewer units than would be obtained by treating each output separately (especially if the output targets are correlated); in other words, maximum mutual use should be made of hidden units. Consider the following algorithm, where steps 1 and 2 are repeated until every output unit makes no mistakes:

Start. There are no hidden units and no connections, so the output units are always OFF.

Step 1. Choose an output unit (say, the one which makes the most errors). Build the appropriate hidden unit to correct some of the mistakes being made by this output unit, as described above. Connect this new unit to *all* the output units.

Step 2. Train the weights from each unit in this enlarged hidden layer to each of the output units. Reevaluate the numbers of errors made by each output unit.

Hence a single hidden layer is constructed.

6.2 Conclusions. The design of general purpose supervised learning algorithms for neural networks involves two important considerations: the network should succeed in correctly classifying the patterns on which it is trained, and nontrivial solutions should involve as few computational elements as possible, avoiding redundant computation. The upstart algorithm can build a network to implement correctly any boolean mapping. Because each “daughter” cell makes as few errors as possible, it corrects as many “parent” errors as possible, which results in small networks. These networks are smaller than those produced by the tiling algorithm. In general the minimum number of units required for any given problem cannot be calculated. However, for a few special cases such as parity and the clumps problem it is known that N or fewer units are needed, and the upstart algorithm achieves this. In the potentially worst case where the targets are randomly assigned, m patterns can be correctly classified by approximately $m/9$ units. The basic idea can be implemented in different architectures and is extendable to the case of multiple outputs.

Acknowledgments

I am grateful to David Willshaw who helped greatly with the manuscript, and also David Wallace for useful comments. I particularly thank Peter

Dayan for suggesting that the tree could be “squashed” into a single hidden layer, as well as for helpful discussions.

References

- Binder, K. 1979. Monte Carlo methods in statistical physics. *Topics Curr. Phys.* 7.
- Denker, J., Schwartz, D., Wittner, B., Solla, S., Howard, R., Jackel, L., and Hopfield, J. 1987. Large automatic learning, rule extraction and generalization. *Complex Syst.* 1, 877–922.
- Frean, M. R. 1990. A “Thermal” perceptron for efficient linear discrimination. Unpublished.
- Gallant, S. I. 1986a. Three constructive algorithms for network learning. *Proc. 8th Ann. Conf. Cog. Sci. Soc.*, Amherst, MA, Aug. 15–17, pp. 652–660.
- Gallant, S. I. 1986b. Optimal linear discriminants. *IEEE Proc. 8th Conf. Pattern Recognition, Paris*, Oct. 28–31, pp. 849–852.
- Mezard, M., and Nadal, J. 1989. Learning in feedforward layered networks: The tiling algorithm. *J. Phys. A* 22, 12, 2191–2203.
- Minsky, M., and Papert, S. 1969. *Perceptrons*. MIT Press, Cambridge.
- Nadal, J. 1989. Study of a growth algorithm for neural networks. *Int. J. Neural Syst.* 1, 55–59.
- Rosenblatt, F. 1962. *Principles of Neurodynamics*. Spartan Books, New York.
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. 1986. Learning internal representations by error propagation. In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition. Volume I Foundations*, D. E. Rumelhart, J. L. McClelland, and the PDP Research Group, eds. MIT Press, Cambridge.

Received 12 December 1989; accepted 19 March 1990.