

# Constructive Neural-Network Learning Algorithms for Pattern Classification

Rajesh Parekh, *Member, IEEE*, Jihoon Yang, *Member, IEEE*, and Vasant Honavar, *Member, IEEE*

**Abstract**—Constructive learning algorithms offer an attractive approach for the incremental construction of near-minimal neural-network architectures for pattern classification. They help overcome the need for *ad hoc* and often inappropriate choices of network topology in algorithms that search for suitable weights in *a priori* fixed network architectures. Several such algorithms are proposed in the literature and shown to converge to zero classification errors (under certain assumptions) on tasks that involve learning a *binary to binary* mapping (i.e., classification problems involving binary-valued input attributes and two output categories). We present two constructive learning algorithms *M*Pyramid-*real* and *M*Tiling-*real* that extend the *pyramid* and *tiling* algorithms, respectively, for learning *real to M-ary* mappings (i.e., classification problems involving real-valued input attributes and multiple output classes). We prove the convergence of these algorithms and empirically demonstrate their applicability to practical pattern classification problems. Additionally, we show how the incorporation of a local pruning step can eliminate several redundant neurons from *M*Tiling-*real* networks.

**Index Terms**—Artificial neural networks, classification, constructive learning algorithms, multcategory, perceptron, pruning, real-valued pattern.

## I. INTRODUCTION

ARTIFICIAL neural networks have been successfully applied to problems in *pattern classification*, *function approximation*, *optimization*, *pattern matching* and *associative memories* [10], [30]. Multilayer feedforward networks trained using the *backpropagation* learning algorithm [44] are limited to searching for a suitable set of weights in an *a priori* fixed network topology. This mandates the selection of an appropriate network topology for the learning problem on hand. However, there are no known efficient methods for determining the optimal network topology for a given problem. Too small networks are unable to adequately learn the problem well while overly large networks tend to overfit the training data and consequently result in poor generalization performance (see [12] for an analogy to the *curve fitting* problem). In practice, a variety of architectures are tried out and the one that appears best suited to the given problem is picked. Such a *trial-and-error*

approach is not only computationally expensive but also does not guarantee that the selected network architecture will be close to optimal or will generalize well. This suggests the need for algorithms that learn both the network topology and the weights.

### A. Constructive Neural-Network Learning Algorithms

*Constructive* (or *generative*) learning algorithms offer an attractive framework for the incremental construction of near-minimal neural-network architectures. These algorithms start with a small network (usually a single neuron) and dynamically grow the network by adding and training neurons as needed until a satisfactory solution is found [20], [23]. Some key motivations for studying constructive neural-network learning algorithms are the following.

- *Flexibility of Exploring the Space of Neural-Network Topologies:*

Constructive algorithms overcome the limitation of searching for a solution in the weight space of an *a priori* fixed network architecture by extending the search, in a controlled fashion, to the entire space of neural-network topologies. Further, it has been shown that at least in principle, algorithms that are allowed to add neurons and weights represent a class of *universal learners* [3].

- *Potential for Matching the Intrinsic Complexity of the Learning Task:*

It is desirable that a learning algorithm construct networks whose complexity (in terms of relevant criteria such as number of nodes, number of links, and connectivity) is commensurate with the intrinsic complexity of the underlying learning task (implicitly specified by the training data) [26]. Constructive algorithms search for small solutions first and thus offer a potential for discovering a near-minimal network that suitably matches the complexity of the learning task. Smaller networks are also preferred because of their potential for more efficient hardware implementation and greater transparency in extracting the learned knowledge.

- *Estimation of Expected Case Complexity of the Learning Task:*

Most practical learning problems are known to be computationally hard to solve. However, little is known about the *expected* case complexity of problems encountered and successfully solved by living systems primarily because it is difficult to mathematically characterize the properties of such problems. Constructive algorithms, if successful, can provide useful empirical estimates of the expected case complexity of practical learning problems.

Manuscript received May 6, 1997; revised October 29, 1998 and October 28, 1999. This work was supported in part by the National Science Foundation Grants IRI-9409580 and IRI-9643299. The work of V. Honavar was funded in part by grants from the National Science Foundation, the John Deere Foundation, the National Security Agency, and IBM

R. Parekh is with Allstate Research and Planning Center, Menlo Park CA 94025 USA (e-mail: rpare@allstate.com).

J. Yang is with Information Sciences Lab, HRL Laboratories LLC, Malibu CA 90265 USA (e-mail: yang@wins.hrl.com).

V. Honavar is with the Department of Computer Science, Iowa State University, Ames, IA 50011 USA (e-mail: honavar@cs.iastate.edu).

Publisher Item Identifier S 1045-9227(00)02997-0.

- *Tradeoffs Among Performance Measures:*

Different constructive learning algorithms allow trading off certain performance measures (e.g., learning time) for others (e.g., network size and generalization accuracy) [47].

- *Incorporation of Prior Knowledge:*

Constructive algorithms provide a natural framework for incorporating problem-specific knowledge into initial network configurations and for modifying this knowledge using additional training examples [14], [33], [34].

- *Lifelong Learning:*

Recent research in *lifelong learning* [48] has proposed training networks that learn to solve multiple related problems by applying the knowledge acquired from the simpler problems to learn the more difficult ones. Constructive learning algorithms lend themselves well to the lifelong learning framework. A network that has domain knowledge from the simpler task(s) built into its architecture (either by explicitly setting the values of the connection weights or by training them) can form a building block for a system that constructively learns more difficult tasks.

## B. Network Pruning

*Network pruning* offers another approach for dynamically determining an appropriate network topology. Pruning techniques (see [40] for an excellent survey) begin by training a larger than necessary network and then eliminate weights and neurons that are deemed redundant. Constructive algorithms offer several significant advantages over pruning-based algorithms including, the ease of specification of the initial network topology, better economy in terms of training time and number of training examples, and potential for converging to a smaller network with superior generalization [27]. In this paper we will focus primarily on constructive learning algorithms. In Section IV we show how a local pruning step can be integrated into the network construction process to obtain more compact networks.

## C. Constructive Algorithms for Pattern Classification

Neural-network learning can be specified as a *function approximation* problem where the goal is to learn an unknown function  $\mathbf{f}: \mathcal{R}^N \rightarrow \mathcal{R}$  (or a good approximation of it) from a set of input-output pairs  $S = \{(\mathbf{x}^N, y) | \mathbf{x}^N \in \mathcal{R}^N, y \in \mathcal{R}\}$ . A variety of constructive neural-network learning algorithms have been proposed for solving the general function approximation problem (see [27] for a survey). These algorithms typically use a *greedy strategy* wherein each new neuron added to the network is trained to minimize the residual error as much as possible. Often the unknown target function ( $\mathbf{f}$ ) is inherently complex and cannot be closely approximated by a network comprising of a single hidden layer of neurons implementing simple transfer functions (e.g., *sigmoid*). To overcome this difficulty, some constructive algorithms use different transfer functions (e.g., the *Gaussian* [21]) while others such as the *projection pursuit regression* [18] use a summation of several *nonlinear* transfer functions. Alternatively, algorithms such as the *cascade correlation* family construct multilayer networks wherein the structural interconnections among the hidden neurons allow the net-

work to approximate complex functions using relatively simple neuron transfer functions like the sigmoid [13], [39], [49].

*Pattern classification* is a special case of function approximation where the function's output  $y$  is restricted to one of  $M$  ( $M \geq 2$ ) discrete values (or classes) i.e., it involves a *real to M-ary* function mapping. A neural network for solving classification problems typically has  $N$  input neurons and  $M$  output neurons. The  $k$ th output neuron ( $1 \leq k \leq M$ ) is trained to output one (while all the other output neurons are trained to output zero) for patterns belonging to the  $k$ th class.<sup>1</sup> Clearly, the class of constructive algorithms that implement the more general *real to real* mapping can be adapted to pattern classification (see [54] for an example). However, a special class of constructive learning algorithms can be designed to closely match the unique demands of pattern classification. Since it is sufficient for each output neuron to be binary valued (i.e., output zero or one), individual neurons can implement the simple *threshold* or *hard-limiting* activation function (with outputs zero and one) instead of a continuous activation function like the sigmoid. Threshold neurons offer the following advantages over their continuous counterparts: First, they are potentially easier to implement in hardware. Second, the *perceptron learning rule* [42] is a simple iterative procedure for training threshold neurons. The learning rules for sigmoid neurons and the like are more complicated and thus computationally more expensive. Third, threshold functions can be clearly described in terms of simple "if-then-else" rules. This makes it easier to incorporate domain expertise (which is usually available in the form of if-then-else rules) into a network of threshold neurons [14]. Similar argument suggests that the task of extracting learned knowledge from a network of threshold neurons would be considerably simpler. In this paper, we will focus on constructive learning of networks of threshold neurons for pattern classification.

1) *Constructive Learning Using Iterative Weight Update:* A number of algorithms that incrementally construct networks of threshold neurons for learning the *binary to binary* mapping have been proposed in the literature (for example, the *tower*, *pyramid* [19], *tiling* [31], *upstart* [15], *oil-spot* [29], and *sequential* [28] algorithms). These algorithms differ in terms of their choices regarding: restrictions on input representation (e.g., binary or bipolar valued inputs); when to add a neuron; where to add a neuron; connectivity of the added neuron; weight initialization for the added neuron; how to train the added neuron (or a subnetwork affected by the addition); and so on. They can be shown to converge to networks which yield zero classification errors on any noncontradictory training set involving two output classes (see [47] for a unifying framework that explains the convergence properties of these constructive algorithms). A geometrical analysis of the decision boundaries of some of these algorithms is presented in [7]. Practical pattern classification often requires assigning patterns to  $M$  (where  $M > 2$ ) categories. Although in principle, the  $M$  category classification task can be decomposed into  $M$  2-category classification tasks, this approach does not take into account the interrelationships between the  $M$  output categories. Further, the

<sup>1</sup>A single output neuron suffices in the case of problems that involve two category classification.

algorithms mentioned above have been designed to operate on binary (or bipolar) valued attributes only. Real-valued attributes are almost invariably encountered in practical classification tasks. One work around for this problem is to discretize the real-valued attributes prior to training [11]. Discretization can result in a loss of information and also greatly increase the number of input attributes. Thus, it is of interest to design algorithms that can directly accept real-valued attributes. We present constructive neural-network learning algorithms that are capable of handling multiple output categories and real-valued pattern attributes.

2) *Exploiting Geometric Properties for Constructive Learning*: The class of constructive learning algorithms we focus on in this paper trains individual neurons using an iterative weight update strategy (such as the *perceptron* rule). Another class of constructive learning algorithms that use a *one-shot* learning strategy deserves mention. These algorithms exploit the geometric properties of the training patterns to directly (i.e., in one-shot) determine appropriate weights for the neurons added to the network. The *grow and learn* (GAL) algorithm [1] and the *DistAl* algorithm [52] construct a single hidden layer network that implements a kind of *nearest neighbor classification* scheme. Each hidden neuron is an *exemplar* representing a group of patterns that belong to the same class and are close to each other in terms of some suitably chosen distance metric. The *minimizing resources* method [43], the *multisurface* method [4], and the *Voronoi diagram* approach [5] are based on the idea of *partitioning* the input space by constructing linear hyperplanes. Hidden layer neurons are trained to partition the input space into homogeneous regions where each region contains patterns belonging to a single output class. The output layer neurons combine regions that represent the same output class. The geometric approach to constructive learning can be applied successfully in solving small to medium scale problems. However, the global search strategy employed by these algorithms can pose a limitation when learning from very large training sets [47]. Further, the reliance on a suitably chosen distance metric (in the case of *GAL* and *DistAl*) makes it imperative for the user to try out a variety of distance metrics for each learning problem.

In this paper, we present extensions of the *pyramid* and the *tiling* algorithms to handle multiple output classes and real-valued pattern attributes.<sup>2</sup> We prove the convergence of these algorithms and demonstrate their applicability on some practical problems. The remainder of this paper is organized as follows: Section II gives an overview of some elementary concepts and describes the notation used throughout this paper. Sections III and IV describe the *M Pyramid-real* and *MTiling-real* constructive learning algorithms respectively and prove their convergence. Section V illustrates the practical applicability of these algorithms and Section VI concludes with a discussion and some directions for future research.

<sup>2</sup>The framework presented here is more general and can potentially be applied to the entire class of constructive algorithms for pattern classification. The interested reader is referred to [35] for an application of this framework to the *tower*, *upstart*, *perceptron cascade*, and *sequential learning* algorithms.

## II. PRELIMINARIES

### A. Threshold Logic Units

A  $N$ -input *threshold logic unit* (TLU, also known as a *perceptron*) is an elementary processing unit that computes the threshold (hard-limiting) function of the weighted sum of its inputs. The output ( $O^p$ ) of a TLU with weights  $\mathbf{W} = \langle W_0, W_1, W_2, \dots, W_N \rangle$  (where the weight  $W_0$  is referred to as the *threshold* or *bias*) in response to a pattern  $\mathbf{X}^p = \langle X_1^p, X_2^p, \dots, X_N^p \rangle$  is  $O^p = 1$  if  $W_0 + \sum_{i=1}^N W_i \cdot X_i \geq 0$  and  $O^p = -1$  otherwise.<sup>3</sup> For notational convenience, we prefix each pattern  $\mathbf{X}^p$  with a 1 (i.e.,  $\mathbf{X}^p = \langle 1, X_1^p, X_2^p, \dots, X_N^p \rangle$ ) and denote the TLU output  $O^p$  as the hard-limiting function of  $\mathbf{W} \cdot \mathbf{X}^p$ .

1) *Perceptron Learning Rule*: An  $N$ -input TLU implements an  $(N - 1)$ -dimensional hyperplane that partitions the  $N$ -dimensional Euclidean pattern space defined by the coordinates  $X_1, \dots, X_N$  into two regions (or classes). A TLU can thus function as a 2-category classifier. Consider a set of *examples*  $S = S^+ \cup S^-$  where  $S^+ = \{(\mathbf{X}^p, C^p) | C^p = 1\}$  and  $S^- = \{(\mathbf{X}^p, C^p) | C^p = -1\}$  ( $C^p$  is the desired output for the input pattern  $\mathbf{X}^p$  and  $p$  ranges from one to  $|S|$ ). A TLU can be trained using the *perceptron learning rule* [42] ( $\mathbf{W} \leftarrow \mathbf{W} + \eta(C^p - O^p)\mathbf{X}^p$  where  $\eta > 0$  is the learning rate) to attempt to find a weight vector  $\hat{\mathbf{W}}$  such that  $\forall \mathbf{X}^p \in S^+, \hat{\mathbf{W}} \cdot \mathbf{X}^p \geq 0$  and  $\forall \mathbf{X}^q \in S^-, \hat{\mathbf{W}} \cdot \mathbf{X}^q < 0$ . If such a weight vector ( $\hat{\mathbf{W}}$ ) exists then  $S$  is said to be *linearly separable*.

2) *Stable Variants of the Perceptron Rule*: If the set  $S$  is not linearly separable then the perceptron algorithm behaves poorly in the sense that the classification accuracy on the training set can fluctuate widely from one training epoch to the next. Several modifications to the perceptron algorithm (e.g., the *pocket algorithm* with *ratchet modification* [19], the *thermal perceptron algorithm* [16], the *loss minimization algorithm* [24], and the *barycentric correction procedure* [38]) are proposed to find a reasonably good weight vector that correctly classifies a large fraction of the training set  $S$  when  $S$  is not linearly separable and to converge to zero classification errors when  $S$  is linearly separable. Siu *et al.* have established the necessary and sufficient conditions for a training set  $S$  to be nonlinearly separable [46]. They have also shown that the problem of identifying a largest linearly separable subset of  $S$  is NP-complete. Thus, we rely on a suitable heuristic algorithm (such as the *pocket algorithm* with *ratchet modification*) to correctly classify as large a subset of the training patterns as possible in the limited training time allowed. We denote such an algorithm by  $\mathcal{A}$ . In our experiments with constructive learning algorithms we use the *thermal perceptron algorithm* to train individual TLU's. The weight update equation of the *thermal perceptron algorithm* is

$$\mathbf{W} \leftarrow \mathbf{W} + \eta \frac{1}{T} (C^p - O^p) \mathbf{X}^p e^{-|n^p|/T}$$

where  $n^p$  is the net input ( $\mathbf{W} \cdot \mathbf{X}^p$ ) and  $T$  is the temperature.  $T$  is set to an initial value  $T_0$  at the start of learning and gradually

<sup>3</sup>Although *bipolar* TLU's whose outputs are one and  $-1$  are functionally equivalent to *binary* TLU's whose outputs are one and zero, empirical evidence suggests that networks constructed using bipolar TLU's are often smaller than those constructed using binary neurons for the same task [20].

annealed to zero as the training progresses. The damping factor ( $e^{-ln^p/T}$ ) prevents any large weight changes toward the end of the training thereby avoiding any irreversible deterioration in the TLU's classification accuracy.

### B. Multiclass Discrimination

Classification problems involving  $M$  ( $M > 2$ ) output classes typically require a layer of  $M$  TLU's. These TLU's can be trained by *independent* training or as a *winner-take-all* (WTA) group. The former strategy trains the TLU's independently and in parallel. However, this does not take into account the fact that class assignments are crisp (i.e., a pattern assigned to class  $i$  cannot possibly belong to any other class as well) and thus potentially results in scenarios where more than one TLU has an output of one. The WTA training strategy gears the weight changes so that the  $i$ th TLU has the highest net input among the group of  $M$  TLU's in response to a pattern belonging to class  $i$ . The winner (i.e., the neuron with the highest net input) is assigned an output of one while all other neurons are assigned outputs of  $-1$ . In the event of a tie for the highest net input all neurons are assigned outputs of  $-1$ . If a pattern is misclassified then the weights of the TLU's whose output in response to the pattern does not match the desired output are updated using the perceptron rule (or one of its variants). WTA training offers an advantage over independent training in that pattern classes that are only pairwise separable from each other can be correctly classified using WTA training while in independent training only pattern classes that are independently separable from all the other classes can be correctly classified [20].

### C. Preprocessing

Most constructive learning algorithms are designed for binary (or bipolar) valued inputs. An extension of the *upstart* algorithm [45] and the *perceptron cascade* algorithm [6] proposed a preprocessing technique to handle patterns with real-valued attributes wherein the patterns are projected on to a parabolic surface by appending to each pattern ( $\mathbf{X}^p = \langle X_1^p, \dots, X_N^p \rangle$ ) an additional attribute  $X_{N+1}^p = \sum_{i=1}^N (X_i^p)^2$ . With this projection it is possible to train a TLU to exclude any one pattern from all others such that the TLU outputs one for the pattern to be excluded and  $-1$  for all the others. We use this projection idea to demonstrate the convergence of the *MPyramid-real* algorithm on real-valued pattern attributes (see Section III-A).

### D. Notation

The following notation is used in the description of the algorithms and their convergence proofs:

$N$	number of inputs;
$M$	number of outputs;
$I$	input layer index;
$1, 2, \dots, L$	indexes for other layers (hidden and output);
$ l $	number of neurons in layer $l$ ;
$l_1, l_2, \dots, l_{ l }$	indexing of neurons in layer $l$ ;
$\mathbf{W}_{l_i}$	weight vector of neuron $i$ in layer $l$ ;
$\langle W_{l_i,0}, W_{l_i,1}, \dots, W_{l_i, l } \rangle$ ,	
$W_{l_i,k} \in \mathcal{R}, k = 0 \dots  l $	

$W_{l_i, l_j}^p$

$S = \{\mathbf{X}^1, \mathbf{X}^2, \dots\}$

$\mathbf{X}^p = \langle X_1^p, \dots, X_N^p \rangle$

where  $X_i^p \in$

$\mathcal{R}$  for all  $i$  and  $1 \leq p \leq |S|$

$\hat{\mathbf{X}}^p = \langle 1, X_1^p, \dots, X_N^p, X_{N+1}^p \rangle$ ,

$n_{l_j}^p =$

$\langle C_1^p, C_2^p, \dots, C_M^p \rangle$ ,

$C_i^p = 1$  if  $\mathbf{X}^p \in$

class  $i$  and  $C_i^p =$

$-1$  otherwise

$\mathbf{O}_l^p = \langle O_{l_1}^p, O_{l_2}^p, \dots, O_{l_{|l|}}^p \rangle$

$O_{l_i}^p \in \{-1, 1\}$  for all  $i$

$e_l$

connection weight between neuron  $i$  in layer  $l^1$  and neuron  $j$  in layer  $l^2$ ;

pattern set;

pattern  $p$ ;

augmented pattern  $p$ ;

projected pattern  $p$ ;

net input of neuron  $l_j$  in response to pattern  $\mathbf{X}^p$ ;

target output for pattern  $\mathbf{X}^p$ ;

layer  $l$ 's output in response to the pattern  $\mathbf{X}^p$ ;

number of misclassifications at layer  $l$ .

Define a function  $sgn: \mathcal{R} \rightarrow \{-1, 1\}$  as  $sgn(x) = -1$  if  $x < 0$  and  $sgn(x) = 1$  if  $x \geq 0$ . Note that bipolar TLU's implement the *sgn* function of their net input. The input layer neurons are designed to allow the patterns to be input to the network and thus simply copy their input to their output.

## III. THE MPYRAMID-REAL ALGORITHM

The *pyramid* algorithm [19] constructs a layered network of TLU's by successively placing each new TLU above the existing ones. The first neuron receives inputs from the  $N$  input neurons. Each succeeding neuron receives inputs from the  $N$  input neurons and from each of the neurons below itself. Thus, the second neuron receives a total of  $N + 1$  inputs, the third neuron receives a total of  $N + 2$  inputs and so on. Each newly added neuron takes over the role of the output neuron. The network growth continues until the desired classification accuracy is achieved.

The extension of the pyramid algorithm to handle real-valued attributes involves modifying each input pattern by augmenting an extra attribute ( $X_{N+1}^p$ ) as described in Section II-C. The network thus has  $N + 1$  input neurons. To handle multiple output categories the algorithm uses  $M$  neurons (instead of one) in each layer of the network. The newly added layer of  $M$  neurons becomes the network output layer. Each of the  $M$  neurons in the new layer are connected to the  $N + 1$  input neurons and to each of the  $M$  neurons in each layer below the current one. This algorithm is described in Fig. 1 and an example network is shown in Fig. 2.

### A. Convergence Proof

*Theorem 1:* There exists a set of weights for neurons in the newly added layer  $L$  of the network such that the number of misclassifications is reduced by at least one (i.e.,  $\forall L > 1, e_L < e_{L-1}$ ).

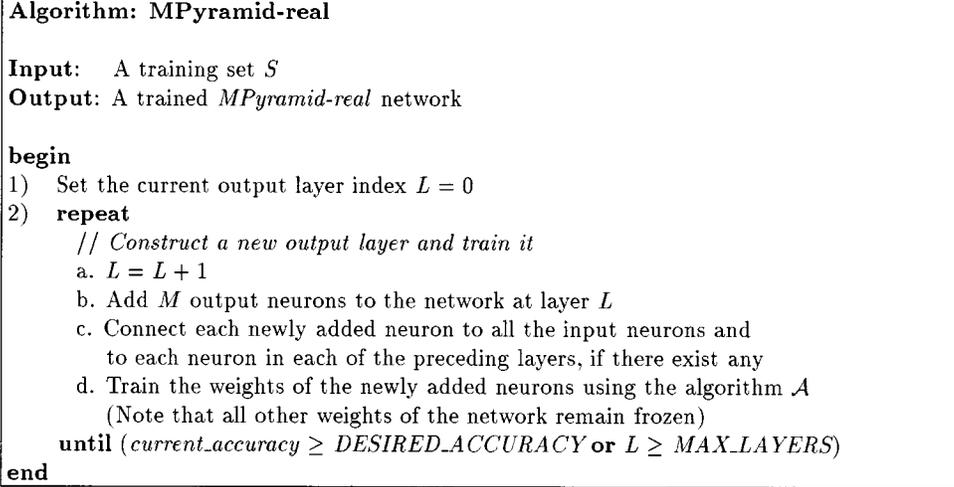


Fig. 1. MPyramid-real algorithm.

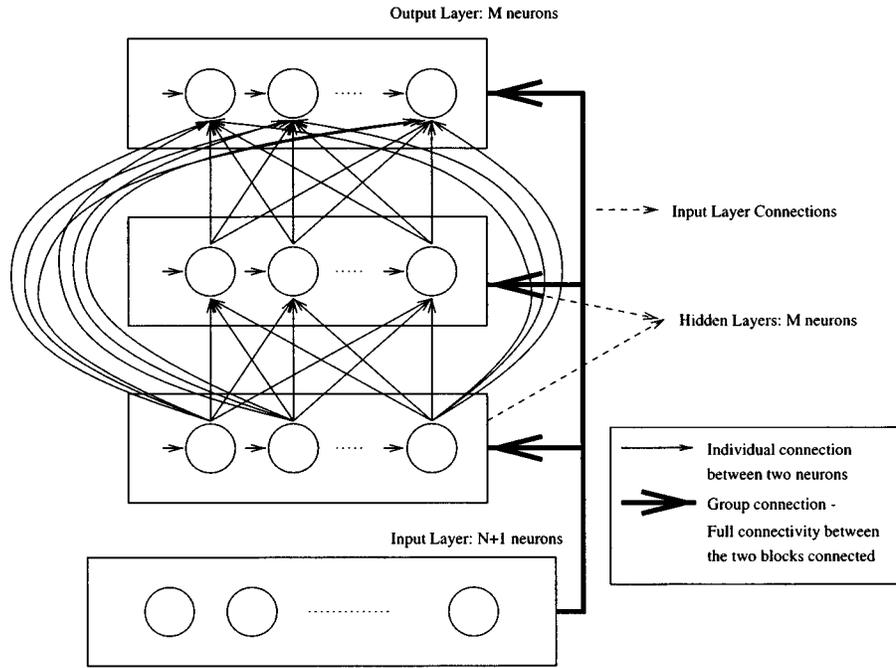


Fig. 2. MPyramid-real network.

*Proof:* Define  $\kappa = \max_{p,q} \sum_{i=1}^N (X_i^p - X_i^q)^2$ . For each pattern  $\hat{X}^p$ , define  $\epsilon_p = (1/2) \min_{q \neq p} \sum_{i=1}^N (X_i^p - X_i^q)^2$ . It is clear that  $0 < \epsilon_p < \kappa$  for all patterns  $\hat{X}^p$ . Assume that a pattern  $\hat{X}^p$  is not correctly classified at layer  $L-1$  (i.e.,  $C^p \neq O_{L-1}^p$ ). Further, let the output vector  $O_{L-1}^p$  for the misclassified pattern  $\hat{X}^p$  be such that  $O_{L-1,\beta}^p = 1$  and  $O_{L-1,k}^p = -1, \forall k = 1 \dots M, k \neq \beta$ ; whereas the target output  $C^p$  is such that  $C_\gamma^p = 1$  and  $C_l^p = -1, \forall l = 1 \dots M, l \neq \gamma$ , and  $\gamma \neq \beta$ .

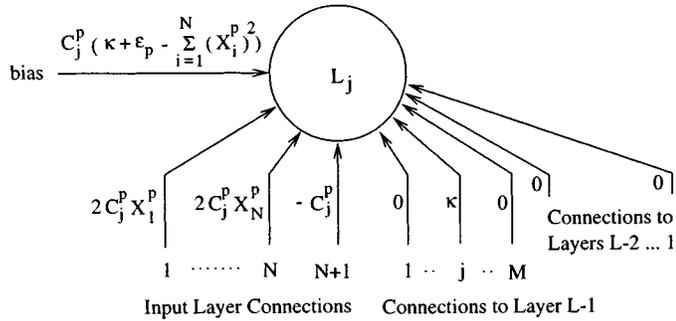
The network adds a new layer  $L$  of  $M$  neurons. Let the weights of these new neurons  $L_j (j = 1 \dots M)$  be set as follows (see Fig. 3):

$$W_{L_j,0} = C_j^p \left( \kappa + \epsilon_p - \sum_{i=1}^N (X_i^p)^2 \right)$$

$$\begin{aligned} W_{L_j, I_i} &= 2C_j^p X_i^p \text{ for } i = 1 \dots N \\ W_{L_j, I_{N+1}} &= -C_j^p \\ W_{L_j, L-k_i} &= 0 \text{ for } k = 2 \dots L-1, \text{ and } i = 1 \dots M \\ W_{L_j, L-1_j} &= \kappa \\ W_{L_j, L-1_i} &= 0 \text{ for } i = 1 \dots M, i \neq j. \end{aligned} \quad (1)$$

For the pattern  $\hat{X}^p$  the net input  $n_{L_j}^p$  of neuron  $L_j$  is

$$\begin{aligned} n_{L_j}^p &= W_{L_j,0} + \sum_{i=1}^{N+1} W_{L_j, I_i} X_i^p + \sum_{k=1}^{j-1} \sum_{i=1}^M W_{L_j, L-k_i} O_{L-k_i}^p \\ &= W_{L_j,0} + \sum_{i=1}^{N+1} W_{L_j, I_i} X_i^p + \sum_{i=1}^M W_{L_j, L-1_i} O_{L-1_i}^p \end{aligned}$$


 Fig. 3. Weight setting for the output neuron  $L_j$  of the *MPyramid-real* network.

$$\begin{aligned}
 & \text{since } W_{L_j, L-k_i} = 0 \text{ for } k = 2, \dots, L-1 \\
 & \text{and } i = 1, \dots, M \text{ [see (1)]} \\
 & = C_j^p (\kappa + \epsilon_p - \sum_{i=1}^N (X_i^p)^2) + 2C_j^p \sum_{i=1}^N (X_i^p)^2 \\
 & \quad - C_j^p \sum_{i=1}^N (X_i^p)^2 + \kappa O_{L-1,j}^p \\
 & = C_j^p (\kappa + \epsilon_p) + \kappa O_{L-1,j}^p. \tag{2}
 \end{aligned}$$

Thus, the net inputs of the output neurons  $L_\gamma$ ,  $L_\beta$ , and  $L_j$  where  $j = 1 \dots M$ ;  $j \neq \gamma$ ,  $j \neq \beta$  are

$$\begin{aligned}
 n_{L_\gamma}^p &= C_\gamma^p (\kappa + \epsilon_p) + \kappa O_{L-1,\gamma}^p \\
 &= \epsilon_p \\
 n_{L_\beta}^p &= C_\beta^p (\kappa + \epsilon_p) + \kappa O_{L-1,\beta}^p \\
 &= -\epsilon_p \\
 n_{L_j}^p &= C_j^p (\kappa + \epsilon_p) + \kappa O_{L-1,j}^p \\
 &= -2\kappa - \epsilon_p.
 \end{aligned}$$

Since  $\epsilon_p > 0$ , for all  $p$ , the net input of neuron  $L_\gamma$  is higher than that of every other neuron in the layer  $L$ . Thus,  $O_{L_\gamma}^p = 1$  and  $O_{L_j}^p = -1$ ,  $\forall j \neq \gamma$  which means that pattern  $\hat{\mathbf{X}}^p$  is correctly classified at layer  $L$ . Even if as a result of a tie for the highest net input, the output of each neuron in layer  $L-1$  in response to  $\hat{\mathbf{X}}^p$  is  $O_{L-1,j}^p = -1$  the weights of the new neurons in layer  $L$  would result in a correct classification of  $\hat{\mathbf{X}}^p$ .

Consider the pattern  $\hat{\mathbf{X}}^q \neq \hat{\mathbf{X}}^p$  that is correctly classified at layer  $L-1$  (i.e.,  $O_{L-1}^q = \mathbf{C}^q$ )

$$\begin{aligned}
 n_{L_j}^q &= W_{L_j,0} + \sum_{i=1}^{N+1} W_{L_j, I_i} X_i^q + \sum_{k=1}^{L-1} \sum_{i=1}^M W_{L_j, L-k_i} O_{L-k_i}^q \\
 &= W_{L_j,0} + \sum_{i=1}^{N+1} W_{L_j, I_i} X_i^q + \sum_{i=1}^M W_{L_j, L-1_i} O_{L-1_i}^q \\
 & \quad \text{since } W_{L_j, L-k_i} = 0 \text{ for } k = 2, \dots, L-1 \\
 & \quad \text{and } i = 1, \dots, M \text{ [see (1)]} \\
 & = C_j^q \left( \kappa + \epsilon_p - \sum_{i=1}^N (X_i^p)^2 \right) + 2C_j^q \sum_{i=1}^N (X_i^p)(X_i^q) \\
 & \quad - C_j^q \sum_{i=1}^N (X_i^q)^2 + \kappa O_{L-1,j}^q
 \end{aligned}$$

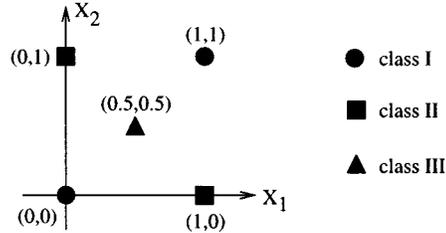


Fig. 4. Example dataset to illustrate the convergence proofs.

 TABLE I  
 DESCRIPTION OF THE EXAMPLE  
 DATASET

Number	Input			Output			$\epsilon_p$
	$X_1$	$X_2$	$X_3 = \sum_{i=1}^2 x_i^2$	$c_1$	$c_2$	$c_3$	
1	0	0	0	-1	-1	1	0.25
2	0	1	1	-1	1	-1	0.25
3	1	0	1	-1	1	-1	0.25
4	1	1	2	-1	-1	1	0.25
5	0.5	0.5	0.5	1	-1	-1	0.25

$$\begin{aligned}
 & = C_j^p (\kappa + \epsilon_p) + \kappa O_{L-1,j}^q \\
 & \quad - C_j^p \sum_{i=1}^N [(X_i^p)^2 - 2(X_i^p)(X_i^q) + (X_i^q)^2] \\
 & = C_j^p (\kappa + \epsilon_p) + \kappa O_{L-1,j}^q - C_j^p \left[ \sum_{i=1}^N (X_i^p - X_i^q)^2 \right] \\
 & = C_j^p (\kappa + \epsilon_p - \epsilon') + \kappa O_{L-1,j}^q \\
 & \quad \text{where } \epsilon' = \sum_{i=1}^N (X_i^p - X_i^q)^2; \text{ note } \epsilon' > \epsilon_p \\
 & = \kappa' C_j^p + \kappa O_{L-1,j}^q \text{ where } \kappa + \epsilon_p - \epsilon' = \kappa', \text{ i.e., } \kappa' < \kappa. \tag{3}
 \end{aligned}$$

The neuron  $L_\gamma$  such that  $O_{L-1,\gamma}^q = 1$  has the highest net input among all output neurons irrespective of the value assumed by  $C_\gamma^p$ . Thus,  $O_{L_\gamma}^q = O_{L-1,\gamma}^q = C_\gamma^q$  i.e., the classification of previously correctly classified patterns remains unchanged. We have shown the existence of weights that will reduce the number of misclassifications whenever a new layer is added to the network. We rely on the TLU weight training algorithm  $\mathcal{A}$  to find such weights. Since the training set is finite in size, eventual convergence to zero errors is guaranteed.  $\square$

### B. Example

The following example illustrates the concepts described in the above proof. Consider a simple dataset shown in Fig. 4. The patterns belong to three output classes and are clearly not linearly separable. Table I summarizes the dataset.

By definition,  $\kappa = \max_{p,q} \sum_{i=1}^N (X_i^p - X_i^q)^2$ . For the example dataset,  $\kappa = 2$ . The first layer of the network (let us designate it by  $L^1$ ) is trained using the algorithm  $\mathcal{A}$ . One possible set of weights for the neurons is depicted in Fig. 5. The response of each of the neurons to the patterns in the dataset is summarized in Table II.

The pattern  $\hat{\mathbf{X}}^4 = \langle 1, 1, 2 \rangle$  is misclassified at layer  $L^1$ . Let  $\hat{\mathbf{X}}^4$  represent the pattern  $\hat{\mathbf{X}}^p$  in the proof. The algorithm adds a new layer of neurons ( $L^2$ ) to the network. Since for  $\hat{\mathbf{X}}^4$  the neuron  $L_2^1$  outputs one whereas the neuron  $L_3^1$  should have output one,  $L_3^2$ ,  $L_2^2$ , and  $L_1^2$  correspond to the neurons  $L_\gamma$ ,  $L_\beta$ , and  $L_j$ , respectively in the proof. Equation (1) specifies one particular set of weights for the newly added neurons (as shown in Fig. 5). The response of the neurons in  $L^2$  to each pattern is given in Table III. The net inputs of the neurons  $L_\gamma$ ,  $L_\beta$ , and  $L_j$  in response to pattern  $\hat{\mathbf{X}}^4$  are  $\epsilon_p$ ,  $-\epsilon_p$  and  $-2\kappa - \epsilon_p$ , respectively, as derived in the proof. Further, the classification of all previously correctly classified patterns such as  $\hat{\mathbf{X}}^1$ ,  $\hat{\mathbf{X}}^2$ ,  $\hat{\mathbf{X}}^3$ , and  $\hat{\mathbf{X}}^5$  that represent the pattern  $\hat{\mathbf{X}}^q$  in the proof remains unaltered.

#### IV. THE *MTILING-REAL* ALGORITHM

The *tiling* algorithm [31] constructs a strictly layered network of threshold neurons. The bottom-most layer receives inputs from each of the  $N$  input neurons. The neurons in each subsequent layer receive inputs from those in the layer immediately below itself. Each layer maintains a *master neuron* and a set (possibly empty) of ancillary neurons that are added and trained to ensure a *faithful representation* of the training patterns. The *faithfulness* criterion states that no two training examples belonging to different classes should produce identical output at any given layer. Faithfulness is clearly a necessary condition for convergence in strictly layered networks [31].

The proposed extension to multiple output classes involves constructing layers with  $M$  master neurons (one for each of the output classes).<sup>4</sup> Unlike the *MPyramid-real* algorithm, it is not necessary to preprocess the dataset using projection. However, it should be noted that such preprocessing will not hamper the convergence properties of the algorithm. Groups of one or more ancillary neurons are trained at a time in an attempt to make the current layer faithful. The algorithm is described in Fig. 6 and an example network is shown in Fig. 7.

##### A. Convergence Proof

Each hidden layer contains  $M$  master neurons and  $K$  ( $K \geq 0$ ) ancillary neurons that are trained to achieve a faithful representation of the patterns. Let  $\bar{S}$  be a subset of the training set  $S$  such that for each pattern  $\mathbf{X}^p$  belonging to  $\bar{S}$  the outputs  $O_1^p, O_2^p, \dots, O_{M+K}^p$  are exactly the same. We designate this output vector  $\langle O_1^p, O_2^p, \dots, O_{M+K}^p \rangle$  as a prototype  $\mathbf{\Pi}^p = \langle \pi_1^p, \pi_2^p, \dots, \pi_{M+K}^p \rangle$ ,  $\pi_i^p = \pm 1$  for all  $i = 1 \dots (M+K)$ . If all the patterns of  $\bar{S}$  belong to exactly one class (i.e., they have the same desired output) then the prototype  $\mathbf{\Pi}^p$  is a faithful representation of the patterns in  $\bar{S}$ . Otherwise,  $\mathbf{\Pi}^p$  is an unfaithful representation of  $\bar{S}$ . Further, if  $\langle \pi_1^p, \pi_2^p, \dots, \pi_M^p \rangle = \langle C_1^p, C_2^p, \dots, C_M^p \rangle$  (i.e., the observed output for the patterns is the same as the desired output) then the patterns in  $\bar{S}$  are said to be correctly classified.

The algorithm's convergence is proved in two parts: first we show that it is possible to obtain a faithful representation of the training set (with real-valued attributes) at the first layer ( $L^1$ ).

<sup>4</sup>An earlier version of this algorithm appeared in [51].

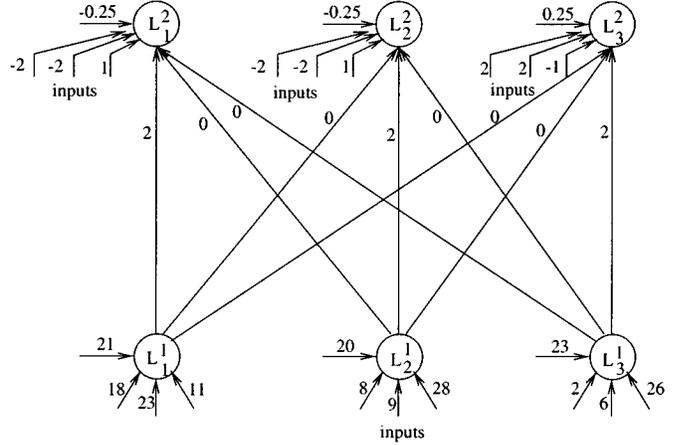


Fig. 5. *MPyramid-real* network for the example dataset.

TABLE II  
RESPONSE OF THE LAYER  $L^1$  NEURONS

Number	Net Input			Output		
	$n_{L_1^1}$	$n_{L_2^1}$	$n_{L_3^1}$	$o_{L_1^1}$	$o_{L_2^1}$	$o_{L_3^1}$
1	21	20	23	-1	-1	1
2	55	57	55	-1	1	-1
3	50	56	51	-1	1	-1
4	84	93	83	-1	1	-1
5	47	42.5	40	1	-1	-1

TABLE III  
RESPONSE OF THE LAYER  $L^2$  NEURONS

Number	Net Input			Output		
	$n_{L_1^2}$	$n_{L_2^2}$	$n_{L_3^2}$	$o_{L_1^2}$	$o_{L_2^2}$	$o_{L_3^2}$
1	-2.25	-2.25	2.25	-1	-1	1
2	-3.25	0.75	-0.75	-1	1	-1
3	-3.25	0.75	-0.75	-1	1	-1
4	-4.25	-0.25	0.25	-1	-1	1
5	0.25	-3.75	-0.25	1	-1	-1

Then we prove that with each additional layer the number of classification errors is reduced by at least one.

*Theorem 2:* For any finite noncontradictory dataset it is possible to train a layer of threshold neurons such that the output of these neurons is a faithful representation of the dataset.

*Proof:* Assume that a layer ( $L^1$ ) with  $M$  master neurons is trained on the dataset ( $S$ ). Consider a subset  $\bar{S}$  of  $S$  such that the master neurons give the same output for each pattern in  $\bar{S}$ . Further assume that  $\bar{S}$  is not faithfully represented by the master neurons. We demonstrate that it is possible to add ancillary neurons (with appropriately set weights) to the layer  $L^1$  in order to obtain a faithful representation of  $\bar{S}$ .

Consider a pattern  $\mathbf{X}^p$  belonging to the *convex hull*<sup>5</sup> of  $\bar{S}$ . If  $\mathbf{X}^p$  is such that for some attribute  $i$  ( $i = 1, \dots, N$ ),  $|X_i^p| > |X_i^q|$  for all  $\mathbf{X}^q \in \bar{S}$  and  $\mathbf{X}^q \neq \mathbf{X}^p$  then an ancillary neuron

<sup>5</sup>The convex hull for a set of points  $Q$  is the smallest convex polygon  $P$  such that each point in  $Q$  lies either on the boundary of  $P$  or in its interior. The interested reader is referred to [8] for a detailed description of convex hulls and related topics in computational geometry.

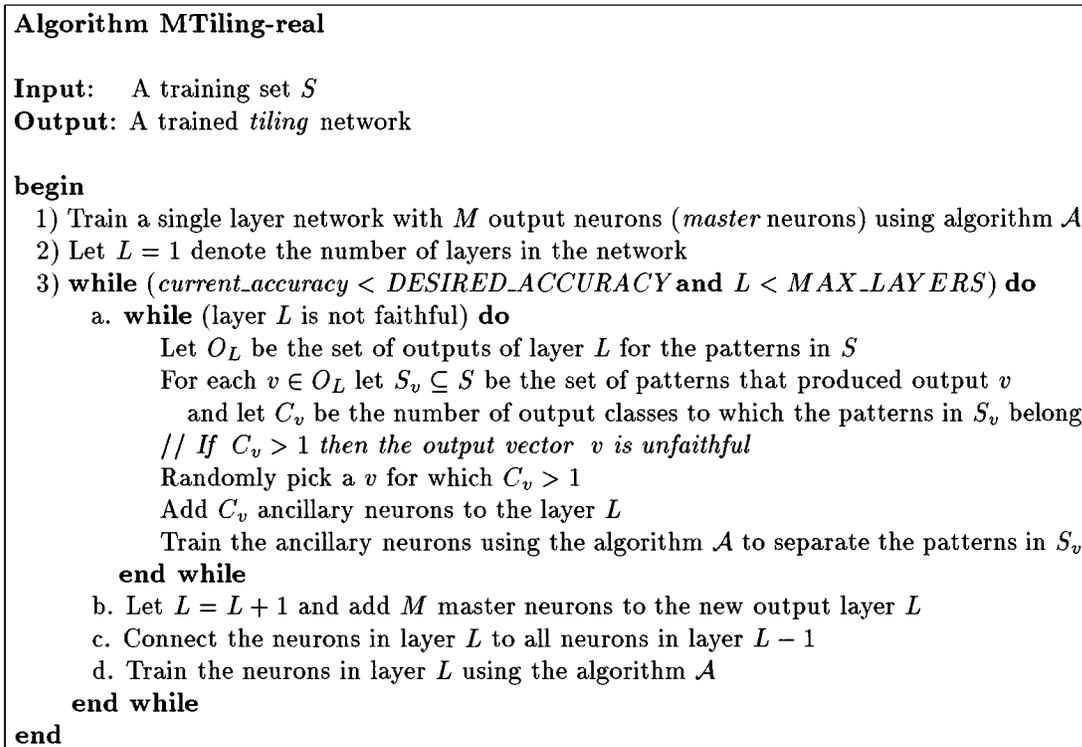


Fig. 6. *MTiling-real* algorithm.

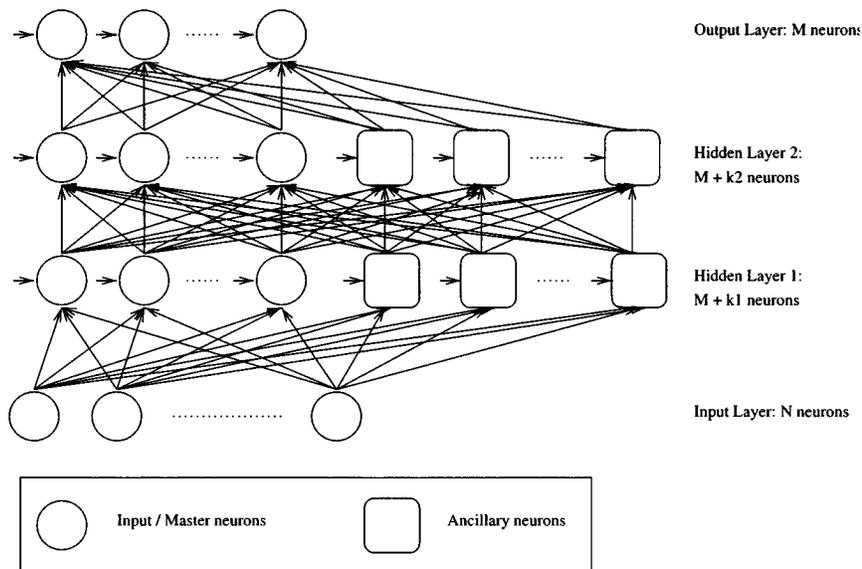


Fig. 7. *MTiling-real* network.

with weights  $\mathbf{W} = \{-(X_i^p)^2, 0, \dots, 0, X_i^p, 0, \dots, 0\}$  (i.e., all weights except  $W_0$  and  $W_i$  set to zero) will output one for  $\mathbf{X}^p$  and  $-1$  for all other patterns. If however,  $\bar{S}$  is such that there is a tie for the highest value of each attribute  $X_i$  among the patterns then there must exist a pattern  $\mathbf{X}^p$  on the convex hull of  $\bar{S}$  that dominates all others in the sense that for each attribute  $X_i$ ,  $X_i^p \geq X_i^q$  for all  $\mathbf{X}^q$  in  $\bar{S}$  (note that  $\mathbf{X}^q \neq \mathbf{X}^p$ ). Clearly,  $\mathbf{X}^p \cdot \mathbf{X}^p > \mathbf{X}^p \cdot \mathbf{X}^q$ . An ancillary neuron with weights  $\mathbf{W} = \{-\sum_{i=1}^N (X_i^p)^2, X_1^p, \dots, X_N^p\}$  will output one for  $\mathbf{X}^p$  and  $-1$  for all other patterns in  $\bar{S}$ .

After adding an ancillary neuron, the output of the layer  $L^1$  is faithful in response to  $\mathbf{X}^p$ . Note that this output is distinct from the outputs for all the other patterns in the entire training set  $S$ . In effect, the pattern  $\mathbf{X}^p$  has been *excluded* from the remaining patterns in the training set. Similarly, using additional TLU's (up to  $|\bar{S}|$  TLU's in all) it can be shown that the outputs of the neurons in the layer provide a faithful representation of  $\bar{S}$ .  $\square$

In practice, by training a groups of one or more ancillary neurons at a time it is possible to attain a faithful representation

of the input pattern set at the first hidden layer using far fewer TLU's as compared to the total number of training patterns.

**Theorem 3:** There exists a set of weights for the master neurons of a newly added layer  $L$  in the network such that the number of misclassifications is reduced by at least one (i.e.,  $\forall L > 1, e_L < e_{L-1}$ ).

*Proof:* Consider a set  $S_1 \subseteq S$  of patterns that belong to the same output class but are not correctly classified by the master neurons in layer  $L - 1$ . Let the prototype  $\mathbf{\Pi}^p = \langle \pi_1^p, \pi_2^p, \dots, \pi_{M+K}^p \rangle$  represent the output of layer  $L - 1$  in response to the patterns in  $S_1$ . Since the patterns in  $S_1$  are not correctly classified at layer  $L - 1$ ,  $\langle \pi_1^p, \pi_2^p, \dots, \pi_M^p \rangle \neq \langle C_1^p, C_2^p, \dots, C_M^p \rangle$  (the desired output for the patterns in  $S_1$ ). For the incorrectly classified prototype  $\mathbf{\Pi}^p$  assume that  $\pi_\beta^p = 1, 1 \leq \beta \leq M$  and  $\forall j = 1 \dots M, j \neq \beta, \pi_j^p = -1$ . Clearly,  $C_\beta^p = -1$  and  $\exists \gamma 1 \leq \gamma \leq M, \gamma \neq \beta$  such that  $C_\gamma^p = 1$ .

The algorithm adds a new layer ( $L$ ) of  $M$  master neurons. The following weights for the master neuron  $L_j$  ( $j = 1 \dots M$ ) shown in Fig. 8 results in the correct classification of patterns in  $S_1$ . It also ensures that the output of any other set of patterns  $S_2 \subseteq S$  (with corresponding prototype  $\mathbf{\Pi}^q = \langle \pi_1^q, \pi_2^q, \dots, \pi_{M+K}^q \rangle$  and  $S_1 \cap S_2 = \emptyset$ ) for which the master neurons of layer  $L - 1$  produce correct outputs (i.e.,  $\langle \pi_1^q, \pi_2^q, \dots, \pi_M^q \rangle = \langle C_1^q, C_2^q, \dots, C_M^q \rangle$ ) remains unchanged

$$\begin{aligned} W_{L_j,0} &= 2C_j^p \\ W_{L_j,L-1_k} &= C_j^p \pi_k^p \text{ for } k = 1 \dots |L-1|, k \neq j \\ W_{L_j,L-1_j} &= |L-1|. \end{aligned} \quad (4)$$

From (4), the net input of neuron  $L_j$  in response to the prototype  $\mathbf{\Pi}^p$  is

$$\begin{aligned} n_{L_j}^p &= W_{L_j,0} + \sum_{k=1}^{|L-1|} W_{L_j,L-1_k} \pi_k^p \\ &= 2C_j^p + |L-1|\pi_j^p + \sum_{k=1, k \neq j}^{|L-1|} C_j^p \pi_k^p \pi_k^p \\ &= 2C_j^p + |L-1|\pi_j^p + (|L-1| - 1)C_j^p \\ &= |L-1|\pi_j^p + (|L-1| + 1)C_j^p. \end{aligned} \quad (5)$$

Thus,

$$\begin{aligned} n_{L_\gamma}^p &= |L-1|(-1) + (|L-1| + 1)(1) \\ &= 1 \\ n_{L_\beta}^p &= |L-1|(1) + (|L-1| + 1)(-1) \\ &= -1 \text{ where } 1 \leq \beta \leq M, \beta \neq \gamma \\ n_{L_k}^p &= |L-1|(-1) + (|L-1| + 1)(-1) \\ &\text{ for } k = 1 \dots M, k \neq \gamma, k \neq \beta \\ &= -2|L-1| - 1. \end{aligned}$$

The master neuron  $L_\gamma$  has the highest net input among all master neurons in layer  $L$  which means that  $O_{L_\gamma}^p = 1$  and  $O_{L_j}^p = -1, \forall j = 1 \dots M, j \neq \gamma$  and  $C^p = O_{L_\gamma}^p$ . Thus, the patterns in  $S_1$  are now correctly classified. Even if as a result of a tie in the value of the highest net input among the master neurons of

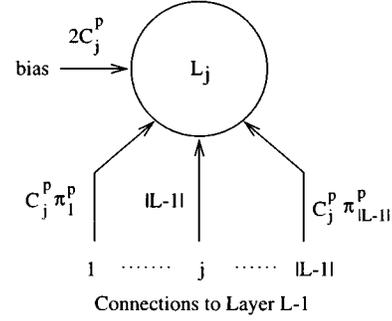


Fig. 8. Weight setting for the output neuron  $L_j$  of the *MTiling-real* network.

layer  $L - 1$ ,  $\mathbf{\Pi}^p$  is such that  $\pi_i^p = -1$  for all  $i = 1 \dots M$ , we see from (5) that the net input of neuron  $L_\gamma$  is still the largest among the net inputs of the  $M$  master neurons in layer  $L$ . Thus, the patterns in  $S_1$  would be correctly classified.

Now consider the prototype  $\mathbf{\Pi}^q$  of the set of patterns  $S_2$  that are correctly classified by the network at layer  $L - 1$  (as described earlier). The net input of the master neurons at layer  $L$  in response to the prototype  $\mathbf{\Pi}^q$  is

$$\begin{aligned} n_{L_j}^q &= W_{L_j,0} + \sum_{k=1}^{|L-1|} W_{L_j,L-1_k} \pi_k^q \\ &= 2C_j^p + |L-1|\pi_j^q + \sum_{k=1, k \neq j}^{|L-1|} W_{L_j,L-1_k} \pi_k^q \\ &= 2C_j^p + |L-1|\pi_j^q + \sum_{k=1, k \neq j}^{|L-1|} C_j^p \pi_k^p \pi_k^q \\ &= 2C_j^p + |L-1|\pi_j^q + C_j^p \sum_{k=1}^{|L-1|} \pi_k^p \pi_k^q - C_j^p \pi_j^p \pi_j^q. \end{aligned} \quad (6)$$

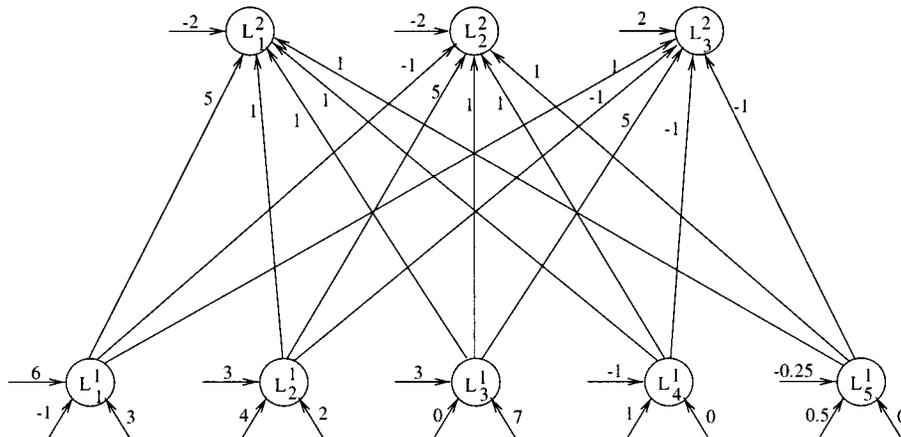
Since  $\mathbf{\Pi}^q \neq \mathbf{\Pi}^p$ ,  $-|L-1| \leq \sum_{k=1}^{|L-1|} C_j^p \pi_k^p \pi_k^q \leq |L-1| - 2$ . Consider a neuron  $L - 1_\alpha$  ( $1 \leq \alpha \leq M$ ) such that  $\pi_\alpha^q = 1$ . From (6)

$$\begin{aligned} n_{L_\alpha}^q &= 2C_\alpha^p + |L-1|\pi_\alpha^q + C_\alpha^p \sum_{k=1}^{|L-1|} \pi_k^p \pi_k^q - C_\alpha^p \pi_\alpha^p \pi_\alpha^q \\ &\geq 2C_\alpha^p + |L-1|(1) - C_\alpha^p |L-1| - C_\alpha^p \pi_\alpha^p (1). \end{aligned}$$

If  $C_\alpha^p = 1$  then  $\pi_\alpha^p = -1$  (since the prototype  $\mathbf{\Pi}^p$  was not correctly classified). On the other hand, if  $C_\alpha^p = -1$  then  $\pi_\alpha^p$  could be either one or  $-1$ . In either case, after substituting these values in the above equation we see that  $n_{L_\alpha}^q \geq 3$ . Further for any other neuron  $L - 1_j$  where  $j = 1 \dots M, j \neq \alpha, \pi_j^q = -1$

$$\begin{aligned} n_{L_j}^q &= 2C_j^p + |L-1|\pi_j^q + C_j^p \sum_{k=1}^{|L-1|} \pi_k^p \pi_k^q - C_j^p \pi_j^p \pi_j^q \\ &\leq 2C_j^p + |L-1|(-1) + C_j^p (|L-1| - 2) - C_j^p \pi_j^p (-1). \end{aligned}$$

If  $C_j^p = 1$  then  $\pi_j^p = -1$  (since the prototype  $\mathbf{\Pi}^p$  was not correctly classified). On the other hand, if  $C_j^p = -1$  then  $\pi_j^p$  could be either one or  $-1$ . In either case, after substituting these values in the above equation we see that  $n_{L_j}^q \leq -1$ . The neuron  $L_\alpha$  has the highest net input among all the master neurons  $L_j$  ( $j =$


 Fig. 9. *MTiling-real* network for the example dataset.

$1 \cdots M$ ). Thus,  $O_{L^q_\alpha}^q = 1$  and  $O_{L^q_j}^q = -1 \forall j = 1 \cdots M, j \neq \alpha$ , which means that  $\mathbf{C}^q = \mathbf{O}^q$ .

We have shown that there exists weights for the master neurons of a newly added layer which will reduce the number of misclassifications by at least one. We rely on the algorithm  $\mathcal{A}$  to find such weights. Since the training set is finite the algorithm would eventually converge to zero classification errors after adding a sufficient number of layers.  $\square$

### B. Example

The following example illustrates the concepts described in the above proof. Consider the simple dataset  $S$  shown in Fig. 4.<sup>6</sup>

The first step in the network construction involves training a layer of  $M = 3$  master neurons. Let us designate this layer by  $L^1$ . One possible set of weights for the master neurons is depicted in Fig. 9. The response of each of master neurons to the patterns in the dataset is summarized in Table IV.

The output of the master neurons is unfaithful with respect to the following two sets of patterns:  $S_1 = \{\mathbf{X}^1, \mathbf{X}^5\}$  and  $S_2 = \{\mathbf{X}^2, \mathbf{X}^4\}$ . Consider the set  $S_2$ .  $\mathbf{X}^4$  has a dominating attribute  $X_1^4 = 1$  ( $X_1^4 > X_1^2 = 0$ ) thus an ancillary neuron with weights  $\langle -(X_1^4)^2, X_1^4, 0 \rangle$  (i.e.,  $\langle -1, 1, 0 \rangle$ ) makes the layer  $L^1$  faithful with respect to  $S_2$  as per theorem 2. Similarly, a second ancillary neuron with weights set to  $\langle -0.25, 0.5, 0 \rangle$  (see Fig. 9) makes  $L^1$  faithful with respect to  $S_1$ . Note that the ancillary neurons are added later (to make the layer faithful) and hence they function as independent TLU's and not as part of the winner-take-all group formed by the master neurons. The output of the patterns in  $S$  at layer  $L^1$  is given in Table V. It can be verified that this output is faithful with respect to all the patterns in  $S$ . Note that patterns  $\mathbf{X}^1$  and  $\mathbf{X}^2$  of the example dataset  $S$  are misclassified at layer  $L^1$ .

Consider a set  $S_1 = \{(0, 0)\}$  with corresponding prototype  $\mathbf{\Pi}^P = \langle 1, -1, -1, -1, -1 \rangle$ .  $S_1$  is misclassified at layer  $L^1$ . The algorithm adds a new layer  $L^2$  of master neurons. The weights of these neurons as per (4) are shown in Fig. 9. The net input of the  $L^2$  neurons in response to  $\mathbf{\Pi}^P$  is  $\langle -1, -11, 1 \rangle$  and thus the output is  $\langle -1, -1, 1 \rangle$  which

<sup>6</sup>Note that the *MTiling-real* algorithm does not require patterns to be pre-processed.

 TABLE IV  
RESPONSE OF THE LAYER  $L^1$  MASTER NEURONS

Number	Net Input			Output		
	$n_{L^1_1}$	$n_{L^1_2}$	$n_{L^1_3}$	$o_{L^1_1}$	$o_{L^1_2}$	$o_{L^1_3}$
1	6	3	3	1	-1	-1
2	9	5	10	-1	-1	1
3	5	7	3	-1	1	-1
4	8	9	10	-1	-1	1
5	7	6	6.5	1	-1	-1

shows that the patterns in  $S_1$  are now correctly classified. Next consider a set  $S_2 = \{(1, 0)\}$  with corresponding prototype  $\mathbf{\Pi}^q = \langle -1, 1, -1, 1, 1 \rangle$ .  $S_2$  is correctly classified at layer  $L^1$ . The net input of the  $L^2$  neurons in response to  $\mathbf{\Pi}^q$  is  $\langle -5, 5, -7 \rangle$  and thus the output is  $\langle -1, 1, -1 \rangle$  which shows that the classification of patterns in  $S_2$  remains unchanged. Similarly, it can be verified that the classification of the patterns  $\langle 1, 1 \rangle$  and  $\langle 0.5, 0.5 \rangle$  also remains unchanged. Thus, we have shown that the addition of layer  $L^2$  reduces the number of misclassifications by at least one.

### C. Pruning Redundant Ancillary Neurons

Each layer of a network constructed using the *MTiling-real* learning algorithm comprises of  $M$  master neurons and  $K$  (where  $K \geq 0$ ) ancillary neurons. The latter are trained to make the layer faithful with respect to the set of training patterns. During training, if the current layer is *unfaithful* then groups of one or more ancillary neurons are trained for each *unfaithful class* of patterns (i.e., patterns that have exactly the same output at the current layer but belong to different output classes). Ideally, one would expect that each layer contains a minimal number of ancillary neurons necessary to achieve faithfulness. However, in practice, hidden layers often have redundant ancillary neurons. This can be attributed to the following two reasons: first, owing to the inherent biases of the TLU weight training algorithm  $\mathcal{A}$  (used to train the ancillary neurons) and the fact that  $\mathcal{A}$  is allowed only a limited training time (typically 500–1000 iterations over the training set) more than one group of ancillary neurons might be trained before

TABLE V  
RESPONSE OF THE LAYER  $L^1$  MASTER AND ANCILLARY NEURONS

Number	Net Input Master			Net Input Ancillary		Output				
	$n_{L_1^1}$	$n_{L_2^1}$	$n_{L_3^1}$	$n_{L_4^1}$	$n_{L_5^1}$	$o_{L_1^1}$	$o_{L_2^1}$	$o_{L_3^1}$	$o_{L_4^1}$	$o_{L_5^1}$
1	6	3	3	-1	-0.25	1	-1	-1	-1	-1
2	9	5	10	-1	-0.25	-1	-1	1	-1	-1
3	5	7	3	0	0.25	-1	1	-1	1	1
4	8	9	10	0	0.25	-1	-1	1	1	1
5	7	6	6.5	-0.5	0	1	-1	-1	-1	1

TABLE VI  
EXAMPLE OF PRUNING

Output Class $C^p$	Master Neurons			Ancillary Neurons			
	$O_{M_1}^p$	$O_{M_2}^p$	$O_{M_3}^p$	$O_{K_1}^p$	$O_{K_2}^p$	$O_{K_3}^p$	$O_{K_4}^p$
I	-1	-1	1	-1	-1	-1	1
II	-1	-1	1	-1	1	-1	1
I	-1	-1	1	-1	-1	1	-1
II	1	-1	-1	-1	1	-1	1
III	1	-1	-1	-1	1	-1	-1

a faithful representation is attained for a particular unfaithful class. Second, as a result of the *locality of training*, whereby each group of ancillary neurons is trained using only a subset of the training patterns, it is possible that not all ancillary neurons are absolutely necessary for faithfulness.

We incorporate a local pruning step in the *MTiling-real* algorithm to remove redundant ancillary neurons. This step is invoked immediately after a layer is made faithful. The check for redundant neurons is simple. Each of the ancillary neurons are systematically dropped (one at a time) and the outputs of the remaining neurons are checked for faithfulness. If the current representation (with an ancillary neuron dropped) is faithful then that ancillary neuron is redundant and is pruned. However, if the current representation is not faithful then the ancillary neuron that was dropped is necessary for faithfulness and hence is brought back. The search for redundant ancillary neurons incurs an additional cost. Let  $K$  be the number of ancillary neurons when the layer is first made faithful and  $|S|$  be the number of training patterns. The ancillary neurons are dropped one at a time and the outputs of the remaining neurons (including the  $M$  master neurons) are checked for faithfulness. The faithfulness test takes  $O((M+K) \cdot |S|)$  time and is repeated  $K$  times (once for each ancillary neuron). Thus, the total time complexity of the pruning step is  $O(K \cdot (M+K) \cdot |S|) \approx O(K^2 \cdot |S|)$ . The outputs of the neurons in response to each training pattern are compared for equality during the faithfulness test. These outputs are computed by the *MTiling-real* algorithm when the layer is determined to be faithful and hence do not have to be recomputed. Further, since each neuron only outputs one or  $-1$  and the faithfulness test only requires an equality check, the search for redundant neurons can be performed very efficiently. We conducted an experimental study of pruning in *MTiling-real* networks (see [36] for details) and found that the total time spent in searching for and pruning redundant neurons is less than 10% of the total training time. Further, pruning resulted in a moderate to substan-

tial (at times as high as 50%) reduction in network size. Below we illustrate pruning with a simple example.

Consider a training set  $S$  comprising of five patterns belonging to three output classes. Assume that a layer  $L$  consisting of  $M = 3$  master neurons (say  $M_1$ ,  $M_2$ , and  $M_3$ ) and  $K = 4$  ancillary neurons (say  $K_1$ ,  $K_2$ ,  $K_3$ , and  $K_4$ ) is trained to achieve a faithful representation of  $S$ . The outputs of the individual neurons in response to the training patterns are depicted in Table VI. During the pruning step when the ancillary neuron  $K_1$  is dropped,  $\langle M_1, M_2, M_3, K_2, K_3, K_4 \rangle$  is faithful with respect to  $S$ . Thus,  $K_1$  is redundant and is pruned. Next, when  $K_2$  is dropped  $\langle M_1, M_2, M_3, K_3, K_4 \rangle$  is not faithful. Thus,  $K_2$  is not redundant and is restored. Next, when  $K_3$  is dropped  $\langle M_1, M_2, M_3, K_2, K_4 \rangle$  (note that  $K_2$  has been restored) is faithful and hence  $K_3$  is pruned. Similarly, we can see that  $K_4$  is not redundant and the final representation of  $L$  is  $\langle M_1, M_2, M_3, K_2, K_4 \rangle$ .

## V. CONSTRUCTIVE LEARNING ALGORITHMS IN PRACTICE

The preceding discussion focused on the theoretical proofs of convergence of the *MPyramid-real* and *MTiling-real* constructive learning algorithms. In this section we present results of a few focused experiments designed to address the following issues.

- 1) The convergence proofs presented above rely on the ability of the network construction strategy to connect a new neuron to an existing network so as to guarantee the existence of weights that will enable the added neuron to improve the resulting network's classification accuracy and the TLU weight training algorithm  $\mathcal{A}$ 's ability to find such a weight setting. Finding an optimal weight setting for each added neuron such that the classification error is maximally reduced when the data is nonseparable is a NP-hard problem [46]. Further, in practice the heuristic

TABLE VII  
DATASETS

<i>Dataset</i>	<b>Size</b>	<b>Inputs</b>	<b>Outputs</b>	<b>Attributes</b>
3 concentric circles ( <b>3-circles</b> )	1800	2	3	real
ionosphere structure ( <b>ionosphere</b> )	351	34	2	real, int
pima indians diabetes ( <b>pima</b> )	768	8	2	real,int
image segmentation ( <b>segmentation</b> )	2310	19	7	real, int
vehicle silhouette ( <b>vehicle</b> )	846	18	4	int

TABLE VIII  
CONVERGENCE RESULTS

<b>Dataset</b>	<b>Performance Parameter</b>	<i>perceptron</i>	<i>MPyramid-real</i>	<i>MTiling-real</i>
<b>3-circles</b>	Network Size	$3.0 \pm 0.0$	$6.0 \pm 0.0$	$46.7 \pm 3.7$
	Train Accuracy %	$44.5 \pm 5.5$	$100.0 \pm 0.0$	$100.0 \pm 0.0$
	Test Accuracy %	$41.4 \pm 5.4$	$99.9 \pm 0.2$	$97.0 \pm 1.2$
<b>ionosphere</b>	Network Size	$1.0 \pm 0.0$	$5.0 \pm 1.3$	$5.5 \pm 2.3$
	Train Accuracy %	$97.5 \pm 1.0$	$100.0 \pm 0.0$	$100.0 \pm 0.0$
	Test Accuracy %	$85.4 \pm 6.4$	$90.6 \pm 4.3$	$86.0 \pm 6.2$

algorithms such as the *thermal perceptron algorithm* that are used in constructive learning are only allowed limited training time (say 500 or 1000 epochs). This makes it important to study the convergence of the proposed constructive algorithms in practice.

- 2) The convergence proofs only guarantee the existence of a set of weights for each newly added neuron (or group of neurons) that will reduce the number of misclassifications by at least one. A network that recruits one neuron to simply memorize each training example can trivially attain zero classification errors. A comparison of the actual size of a trained constructive network with the number of patterns in the training set, would at least partially, address this issue of memorization.
- 3) Regardless of the convergence of the constructive learning algorithms to zero classification errors, a question of practical interest is the algorithms' ability to improve generalization on the test set as the network grows in size. One would expect *over-fitting* to set in eventually as neurons continue to get added in an attempt to reduce the classification error, but we wish to examine whether the addition of neurons improves generalization before over-fitting sets in.

#### A. Datasets

A cross-section of datasets having real-valued pattern attributes and patterns belonging to multiple output classes was selected for the experiments. Table VII summarizes the characteristics of the datasets. **Size** denotes the number of patterns in the dataset, **inputs** indicates the total number of input attributes (of the unmodified dataset), **outputs** represents the number of output classes, and **attributes** describes the type of input attributes of the patterns. The real-world datasets **ionosphere**, **pima**, **segmentation**, and **vehicle** are available at the UCI Machine Learning Repository [32] while the **3-circles** dataset was artificially generated. The **3-circles** dataset comprises of 1800

randomly drawn points from the two dimensional Euclidean space. These points are labeled as belonging to classes 1, 2, and 3 if their distance from the origin is less than one, between one and two, and between two and three, respectively. Each of these datasets is nonlinearly separable.

#### B. Experimental Results

We used the ten-fold cross validation method in our experiments. Each dataset was divided into ten equal sized folds and ten independent runs of each algorithm were conducted for each dataset. For the  $i$ th run, the  $i$ th fold was designated as the test set and the patterns in the remaining nine folds were used for training. At the end of training the network's generalization was measured on the test set. Individual TLU's were trained using the *thermal perceptron algorithm*. The weights of each neuron were initialized at random to a value in the interval  $[-1 \dots 1]$ , the learning rate  $\eta$  was held constant at 1.0, and each neuron was trained for 500 epochs where each epoch involves presenting a set of  $|S|$  randomly drawn patterns from the training set  $S$ . The initial temperature  $T_0$  was set to 1.0 and was dynamically updated at the end of each epoch to match the average net input of the neuron(s) during the entire epoch [6].

Table VIII summarizes the results of experiments designed to test the convergence properties of the constructive learning algorithms. It lists the mean and standard deviation of the network size (the number of hidden and output neurons), the training accuracy, and the test accuracy of the *MPyramid-real* and *MTiling-real* algorithms on the **3-circles** and **ionosphere** datasets. For comparison we include the results of training a single layer network (labeled by *perceptron*) using the *thermal perceptron algorithm*. The training accuracy of the *perceptron* algorithm on both datasets is less than 100% (which confirms the nonlinear separability of the datasets). These results show that not only do the constructive algorithms converge to zero classification errors on the training set but they also generalize fairly well on the unseen test data. Further, a comparison of

TABLE IX  
GENERALIZATION RESULTS

Dataset	Performance Parameter	<i>perceptron</i>	<i>MPyramid-real</i>	<i>MTiling-real</i>
<b>pima</b>	Network Size	1.0 ± 0.0	6.5 ± 4.8	5.7 ± 14.9
	Train Accuracy %	79.3 ± 0.9	79.4 ± 1.9	81.0 ± 4.3
	Test Accuracy %	77.5 ± 3.4	76.8 ± 3.5	77.1 ± 3.5
<b>segmentation</b>	Network Size	7.0 ± 0.0	119.8 ± 34.2	47.1 ± 23.2
	Train Accuracy %	96.0 ± 0.2	94.2 ± 0.9	99.1 ± 1.5
	Test Accuracy %	94.8 ± 2.0	93.0 ± 2.7	99.1 ± 1.7
<b>vehicle</b>	Network Size	4.0 ± 0.0	35.2 ± 28.7	19.4 ± 23.5
	Train Accuracy %	85.5 ± 0.7	87.8 ± 3.3	87.5 ± 4.4
	Test Accuracy %	79.7 ± 5.4	78.2 ± 4.9	77.5 ± 6.2

the network sizes with the total number of patterns in each dataset (see Table VII) conclusively shows that the constructive learning algorithms are not simply memorizing the training patterns by recruiting one neuron per pattern.

Our experiments indicated that the convergence of the algorithms (particularly the *MPyramid-real*) was quite slow on the other three datasets. The slow-down was quite pronounced toward the end of the learning cycle where several new layers were added with minimal increase in classification accuracy. Further, at this stage we observed that the generalization accuracy (measured on an independent test set) of the network was deteriorating with the addition of each new layer. This suggests that in an attempt to correctly classify all patterns the algorithms were over-fitting the training data. In practice we are mostly interested in the network's generalization capability. Most backpropagation type learning algorithms use a separate *hold-out* set (distinct from the *test set*) to stop training when over-fitting sets in. In our experiments to measure the generalization performance of the constructive algorithms we used a similar hold-out sample as follows: The ten-fold cross validation was still used but this time during the  $i$ th run, the  $i$ th fold was designated as the test set, the  $i + 1$ th fold as the independent hold-out set, and the remaining eight folds formed the training set. During the network construction process, the accuracy of the network on the hold-out sample was recorded after each new layer was added and trained. At the end of the training (i.e., when the network converged to 100% classification accuracy or the when the network size reached the maximum number of layers—25 in our case) we pruned the network up to the layer that had the highest accuracy on the hold-out sample. For example, if the trained network had five layers and the accuracy on the hold-out was recorded as 78, 82, 86, 83, and 81% at each of the five layers, respectively, then the layers above layer three were pruned from the network. Note that as a result of the pruning the network's accuracy on the training set will no longer be 100%. At this point we measure the accuracy of the network on the test dataset. It is important to keep in mind that since the test data set is independent of the hold-out set and is not used at all during training the results are not biased or overly optimistic.

Table IX lists the mean and standard deviation of the network size, training accuracy, and test accuracy of the *MPyramid-real* and *MTiling-real* algorithms for the **pima**, **segmentation**, and **vehicle** datasets. We see that the pruned networks generated

by the *MTiling-real* algorithm are smaller than those generated by the *MPyramid-real* algorithm. This is due to the different network construction schemes adopted by the two algorithms. The *MPyramid-real* algorithm uses the entire training set for training each new layer. Thus, on harder training sets it tends to add several layers of neurons without substantial benefits. On the other hand, the *MTiling-real* algorithm breaks up the dataset into smaller subsets (the *unfaithful classes*). Training of the ancillary neurons on these smaller datasets is considerably simpler. Further, given a faithful representation of the patterns at each layer, the master neurons of the succeeding layer are able to significantly reduce the number of misclassifications. The *MTiling-real* algorithm's focus on smaller subsets for training ancillary neurons might actually prove to be disadvantageous on certain datasets (see for example the **3-circles** in Table VIII) because it might expend considerable effort in making the current layer faithful.

As can be seen from Table IX the test accuracy of the *MPyramid-real* and *MTiling-real* algorithms is almost the same as or even slightly worse than that of the single layer network (except in the case of the **segmentation** dataset where *MTiling-real* performs better). This suggests that in the case of the **pima** and **vehicle** datasets the constructive learning algorithms do not add much value. It is possible that these datasets contain irrelevant or noisy attributes that unduly complicate the learning task. Experiments have shown that using a genetic algorithm based feature selection scheme significantly improves the generalization performance of the *DistAl* constructive learning algorithm [50]. In other experiments it has been shown that the choice of the algorithm for training the individual TLU's during constructive learning can significantly impact the convergence and generalization properties of the constructive learning algorithms [35]. It was shown that when the *thermal perceptron algorithm* was replaced by other algorithms such as *barycentric correction procedure* or *pocket algorithm* with *ratchet modification* as the algorithm for training individual TLU's, the performance of the constructive learning algorithms on certain datasets was superior both in terms of convergence properties and generalization ability. It is definitely of interest to further explore the impact of feature subset selection and the choice of different TLU weight training algorithms on the performance of the constructive algorithms. Unfortunately, these issues are beyond the scope of this paper.

The issue of network training times is critical for very large training sets. The *perceptron* algorithm that trains just a single layer of TLU's is clearly faster than the *MPyramid-real* and *MTiling-real* algorithms. From our experiments we have observed that the constructive learning algorithms take between 1.5 and five times as long as the *perceptron* algorithm to train on the datasets considered in this paper. The constructive algorithms typically achieve a reasonably good accuracy relatively quickly. A significant amount of time is expended in adding and training units that only marginally improve the training accuracy. As mentioned earlier, this over-training can potentially worsen the network's generalization performance. In our experiments we allowed the networks to train until either convergence to zero training errors was achieved or 25 layers were of TLU's were added to the network. In order to overcome the problem of over-fitting, the networks were then pruned back based on their performance on an independent hold-out set. In practice, a substantial reduction in training time can be achieved if training is actually stopped as soon as it is observed that the network's performance on the hold-out set is not improving significantly. This form of early-stopping is commonly used in training back-propagation networks.

## VI. CONCLUSIONS

Constructive algorithms offer an attractive approach for the automated design of neural networks. In particular, they eliminate the need for the *ad hoc*, and often inappropriate, *a priori* choice of network architecture, they potentially provide a means of constructing networks whose size (complexity) is commensurate with the complexity of the pattern classification task on hand, and they offer natural ways to incorporate prior knowledge to guide learning and to use constructive learning algorithms in the *lifelong* learning framework. We have focused on a family of algorithms that incrementally construct feedforward networks of threshold neurons.<sup>7</sup> Although a number of such algorithms have been proposed in the literature, most of them are limited to 2-category pattern classification tasks with binary/bipolar valued input attributes. We have presented two constructive learning algorithms *MPyramid-real* and *MTiling-real* that extend the *pyramid* and the *tiling* algorithms, respectively, to handle multicategory classification of patterns that have real-valued attributes. For each of these algorithms we have provided rigorous proofs of convergence to zero classification errors on finite, noncontradictory training sets. This proof technique is sufficiently general (see [35] for an application of this technique to several other constructive learning algorithms).

The convergence of the two algorithms was established by showing that each modification of the network topology guarantees the existence of weights that would reduce the classification error and assuming that there exists a weight modification algorithm  $\mathcal{A}$  that would find such weights. We do not have a rigorous proof that any of the graceful variants of perceptron learning algorithms can in practice, satisfy the requirements imposed on  $\mathcal{A}$ , let alone find an *optimal* (in a suitable

sense of the term—e.g., so as to yield minimal networks) set of weights. The design of TLU training algorithms that (with a high probability) satisfy the requirements imposed on  $\mathcal{A}$  and are at least approximately optimal remains an open research problem. These characteristics of the TLU training algorithm often result in the generation of redundant units during network construction. We have proposed a local pruning strategy that can be used to eliminate redundant neurons (in the *MTiling-real* networks). Experiments with nonlinearly separable datasets demonstrate the practical usefulness of the proposed algorithms. On simpler datasets both the *MPyramid-real* and *MTiling-real* algorithms do converge to fairly compact networks with zero classification errors and good generalizability. However, on more difficult tasks convergence is slow. Further, the network might end up memorizing the hard to classify examples thereby resulting in poor generalization. To address this issue we have used an independent *hold-out* set during training to determine the appropriate final network topology. This technique enhances the capability of constructive learning algorithms to generate compact networks with improved generalization. Although it is hard to determine *a priori* which of the two constructive learning algorithms would be suitable for a particular problem, we recommend using the *MTiling-real* algorithm first (during the preliminary analysis) as it tends to have better convergence properties than the *MPyramid-real* algorithm in practice.

Some directions for future research include the following.

- *Evaluating the Performance of Constructive Learning Algorithms:*  
A systematic experimental and theoretical comparisons of constructive algorithms with other neural network as well as other machine learning algorithms for pattern classification is of interest. Further, a characterization of the inductive and representational biases of the different algorithms will guide users in selecting algorithms for specific problems based on easily measurable properties of the datasets.
- *Hybrid Constructive Learning Algorithms:*  
In related work it was shown that the choice of the specific TLU weight training algorithm can have a significant impact on the performance of constructive learning algorithms [37]. A study of hybrid network training schemes that dynamically select an appropriate network construction strategy, an appropriate TLU weight training algorithm, an appropriate output computation strategy and such to obtain locally optimal performance at each step of the classification task is worth pursuing.
- *Combining Constructive Learning with Feature Selection:*  
The generalization performance of learning algorithms can be often be improved with the help of suitable feature selection techniques. Several feature subset selection algorithms have been proposed in the pattern recognition literature [41]. The effectiveness of genetic algorithms for feature subset selection in conjunction with the *DistAl* algorithm has been demonstrated in [50].
- *Using Boosting and Error-Correcting Output Codes for Improved Generalization:*

Recent advances in machine learning have resulted in the development of techniques such as *boosting* [17] and

<sup>7</sup>Constructive algorithms have also been proposed for the incremental construction of recurrent neural networks (RNN's) that learn *finite state automata* from labeled examples. The interested reader is referred to [22] and [25] for a discussion on constructive learning of RNN.

*error-correcting output codes* [2] for improving the generalization capability of learning algorithms. An application of these techniques in the constructive learning framework is clearly of interest.

• *Knowledge Extraction from Trained Constructive Neural Networks:*

Constructive neural-network learning algorithms have been successfully used for theory refinement. The available domain specific knowledge is incorporated into the initial network topology and is refined based on additional labeled examples using constructive learning [14], [33], [34], [53]. The question now is whether we can use some of the existing strategies (see, for example, [9]) or design suitable new methods for extracting the learned knowledge from a trained constructive network.

#### REFERENCES

- [1] E. Alpaydin, "GAL: Networks that grow when they learn and shrink when they forget," Int. Comput. Sci. Inst., Berkeley, CA, Tech. Rep. TR91-032, 1991.
- [2] G. Bakiri and T. Dietterich, "Solving multiclass learning problems via error-correcting output codes," *J. Artificial Intell. Res.*, vol. 2, pp. 263–286, 1995.
- [3] E. Baum, "A proposal for more powerful learning algorithms," *Neural Comput.*, vol. 1, no. 2, pp. 201–207, 1989.
- [4] K. Bennett and O. Mangasarian, "Neural-network training via linear programming," Dept. Comput. Sci., Univ. Wisconsin, Madison, Tech. Rep. 948, 1990.
- [5] N. Bose and A. Garga, "Neural-network design using Voronoi diagrams," *IEEE Trans. Neural Networks*, vol. 4, pp. 778–787, 1993.
- [6] N. Burgess, "A constructive algorithm that converges for real-valued input patterns," *Int. J. Neural Syst.*, vol. 5, no. 1, pp. 59–66, 1994.
- [7] C.-H. Chen, R. Parekh, J. Yang, K. Balakrishnan, and V. Honavar, "Analysis of decision boundaries generated by constructive neural-network learning algorithms," in *Proc. WCNN'95*, vol. 1, Washington, D.C., Jul. 17–21, 1995, pp. 628–635.
- [8] T. Cormen, C. Leiserson, and R. Rivest, *Introduction to Algorithms*. Cambridge, MA: MIT Press, 1991.
- [9] M. Craven, "Extracting comprehensible models from trained neural networks," Ph.D. dissertation, Dept. Comput. Sci., Univ. Wisconsin, Madison, WI, 1996.
- [10] J. Dayhoff, *Neural-Network Architectures: An Introduction*. New York: Van Nostrand Reinhold, 1990.
- [11] J. Dougherty, R. Kohavi, and M. Sahami, "Supervised and unsupervised discretization of continuous features," in *Proc. 12th Int. Conf. Machine Learning*, San Francisco, CA, 1995, pp. 194–202.
- [12] R. Duda and P. Hart, *Pattern Classification and Scene Analysis*. New York: Wiley, 1973.
- [13] S. Fahlman and C. Lebiere, "The cascade correlation learning algorithm," in *Neural Inform. Syst.* 2, D. Touretzky, Ed. San Mateo, CA: Morgan-Kaufman, 1990, pp. 524–532.
- [14] J. Fletcher and Z. Obradović, "Combining prior symbolic knowledge and constructive neural-network learning," *Connection Sci.*, vol. 5, no. 3/4, pp. 365–375, 1993.
- [15] M. Frean, "The upstart algorithm: A method for constructing and training feedforward neural networks," *Neural Comput.*, vol. 4, pp. 198–209, 1990.
- [16] —, "A thermal perceptron learning rule," *Neural Comput.*, vol. 4, pp. 946–957, 1992.
- [17] Y. Freund and R. Schapire, "A decision-theoretic generalization of on-line learning algorithms and an application to boosting," *J. Comput. Syst. Sci.*, vol. 55, no. 1, pp. 119–139, 1997.
- [18] J. Friedman and W. Stuetzle, "Projection pursuit regression," *J. Amer. Statist. Assoc.*, vol. 76, no. 376, pp. 817–823, 1981.
- [19] S. Gallant, "Perceptron based learning algorithms," *IEEE Trans. Neural Networks*, vol. 1, pp. 179–191, 1990.
- [20] —, *Neural-Network Learning and Expert Systems*. Cambridge, MA: MIT Press, 1993.
- [21] S. Geva and J. Sitte, "A constructive method for multivariate function approximation by multilayer perceptrons," *IEEE Trans. Neural Networks*, vol. 3, pp. 621–624, 1992.
- [22] C. Giles, D. Chen, G.-Z. Sun, H.-H. Chen, Y.-C. Lee, and M. Goudreau, "Constructive learning of recurrent neural networks: Limitations of recurrent cascade correlation and a simple solution," *IEEE Trans. Neural Networks*, vol. 6, pp. 829–836, 1997.
- [23] V. Honavar and V. L. Uhr, "Generative learning structures for generalized connectionist networks," *Inform. Sci.*, vol. 70, no. 1/2, pp. 75–108, 1993.
- [24] T. Hrycej, *Modular Learning in Neural Networks*. New York: Wiley, 1992.
- [25] S. Kremer, "Comments on constructive learning of recurrent neural networks: Limitations of recurrent cascade correlation and a simple solution," *IEEE Trans. Neural Networks*, vol. 7, pp. 1047–1049, 1996.
- [26] T.-Y. Kwok and D.-Y. Yeung, "Objective functions for training new hidden units in constructive neural networks," *IEEE Trans. Neural Networks*, vol. 8, pp. 1131–1148, 1997.
- [27] —, "Constructive algorithms for structure learning in feedforward neural networks for regression problems," *IEEE Trans. Neural Networks*, to be published.
- [28] M. Marchand, M. Golea, and P. Rujan, "A convergence theorem for sequential learning in two-layer perceptrons," *Europhys. Lett.*, vol. 11, no. 6, pp. 487–492, 1990.
- [29] F. Mascioli and G. Martinelli, "A constructive algorithm for binary neural networks: The oil-spot algorithm," *IEEE Trans. Neural Networks*, vol. 6, pp. 794–797, 1995.
- [30] K. Mehrotra, C. Mohan, and S. Ranka, *Elements of Artificial Neural Networks*. Cambridge, MA: MIT Press, 1997.
- [31] M. Mézard and J. Nadal, "Learning feedforward networks: The tiling algorithm," *J. Phys. A: Math. Gen.*, vol. 22, pp. 2191–2203, 1989.
- [32] P. Murphy and D. Aha, "Repository of machine learning databases," Dept. Inform. Comput. Sci., Univ. California, Irvine, CA, <http://www.ics.uci.edu/AI/ML/MLDBRepository.html>, 1994.
- [33] D. W. Opitz and J. W. Shavlik, "Dynamically adding symbolically meaningful nodes to knowledge-based neural networks," *Knowledge-Based Syst.*, vol. 8, no. 6, pp. 301–311, 1995.
- [34] R. Parekh and V. Honavar, "Constructive theory refinement in knowledge-based neural networks," in *Proc. Int. Joint Conf. Neural Networks (IJCNN'98)*, Anchorage, AK, 1998, pp. 2318–2323.
- [35] R. Parekh, J. Yang, and V. Honavar, "Constructive neural-network learning algorithms for multicategory real-valued pattern classification," Dept. Comput. Sci., Iowa State Univ., Tech. Rep. ISU-CS-TR97-06, 1997.
- [36] —, "Pruning strategies for constructive neural-network learning algorithms," in *Proc. IEEE/INNS Int. Conf. Neural Networks, ICNN'97*, 1997, pp. 1960–1965.
- [37] —, "An empirical comparison of the performance of single-layer algorithms for training threshold logic units," *Neural, Parallel, Sci. Comput.*, 2000, to be published.
- [38] H. Poulard, "Barycentric correction procedure: A fast method of learning threshold units," in *Proc. WCNN'95*, vol. 1, Washington, D.C., July 17–21, 1995, pp. 710–713.
- [39] L. Prechelt, "Investigating the cascor family of learning algorithms," *Neural Networks*, vol. 10, no. 5, pp. 885–896, 1997.
- [40] R. Reed, "Pruning algorithms—A survey," *IEEE Trans. Neural Networks*, vol. 4, pp. 740–747, 1993.
- [41] B. Ripley, *Pattern Recognition and Neural Networks*. New York: Cambridge Univ. Press, 1996.
- [42] F. Rosenblatt, "The perceptron: A probabilistic model for information storage and organization in the brain," *Psych. Rev.*, vol. 65, pp. 386–408, 1958.
- [43] P. Rujan and M. Marchand, "Learning by minimizing resources in neural networks," *Complex Syst.*, vol. 3, pp. 229–241, 1989.
- [44] D. Rumelhart, G. Hinton, and R. Williams, "Learning internal representations by error propagation," in *Parallel Distributed Processing: Explorations into the Microstructure of Cognition*. Cambridge, MA: MIT Press, 1986, vol. 1.
- [45] J. Saffery and C. Thornton, "Using stereographic projection as a preprocessing technique for upstart," in *Proc. Int. Joint Conf. Neural Networks*, vol. II, July 1991, pp. 441–446.
- [46] K.-Y. Siu, V. Roychowdhury, and T. Kailath, *Discrete Neural Computation—A Theoretical Foundation*. Englewood Cliffs, NJ: Prentice-Hall, 1995.
- [47] F. Śmieja, "Neural-network constructive algorithms: Trading generalization for learning efficiency?," *Circuits, Syst., Signal Processing*, vol. 12, no. 2, pp. 331–374, 1993.
- [48] S. Thrun, "Lifelong learning: A case study," Carnegie Mellon Univ., Tech. Rep. CMU-CS-95-208, 1995.
- [49] J. Yang and V. Honavar, "Experiments with the cascade-correlation algorithm," *Microcomput. Applicat.*, vol. 17, no. 2, pp. 40–46, 1998.

- [50] —, “Feature subset selection using a genetic algorithm,” *IEEE Intell. Syst. (Special Issue on Feature Transformation and Subset Selection)*, vol. 13, no. 2, pp. 44–49, 1998.
- [51] J. Yang, R. Parekh, and V. Honavar, “MTiling—A constructive neural-network learning algorithm for multi-category pattern classification,” in *Proc. World Congr. Neural Networks '96*, San Diego, CA, 1996, pp. 182–187.
- [52] —, “DistAl: An inter-pattern distance-based constructive learning algorithm,” *Intell. Data Anal.*, vol. 3, pp. 55–73, 1999.
- [53] J. Yang, R. Parekh, V. Honavar, and D. Dobbs, “Data-driven theory refinement using kBDistal,” in *Proc. 3rd Symp. Intell. Data Anal. (IDA'99)*, Amsterdam, The Netherlands, 1999, pp. 331–342.
- [54] D. Yeung, “Constructive neural networks as estimators of bayesian discriminant functions,” *Pattern Recognition*, vol. 26, no. 1, pp. 189–204, 1993.



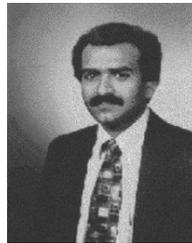
**Rajesh Parekh** (S'89–M'98) received the B.E. degree in computer technology from VJTI, Bombay, India, in 1991, and the M.S. and Ph.D. degrees in computer science with specialization in artificial intelligence from Iowa State University, Ames, in 1993 and 1997, respectively.

He is currently with the Data Mining Group at the Allstate Research and Planning Center, Menlo Park, CA. His research interests include artificial intelligence, applied machine learning, intelligent autonomous agents, knowledge discovery and data mining, neural networks, constructive learning algorithms, computational learning theory, grammatical inference, and distributed artificial intelligence.



**Jihoon Yang** (S'98–M'99) received the B.S. degree in computer science from Sogang University, Seoul, Korea, in 1987, and the M.S. and Ph.D. degrees in computer science with specialization in artificial intelligence from Iowa State University, Ames, in 1989 and 1999, respectively.

He is currently with the Networking and Information Exploitation Department at HRL Laboratories, LLC, Malibu, CA. His research interests include information retrieval, knowledge discovery and data mining, intelligent agents and multiagent systems, machine learning, neural networks, pattern recognition, evolutionary computing, bioinformatics, and distributed artificial intelligence.



**Vasant Honavar** (M'99) received the B.E. degree in electronics engineering from Bangalore University, India, the M.S. degree in electrical and computer engineering from Drexel University, Philadelphia, PA, and the M.S. and Ph.D. degrees in computer science from the University of Wisconsin, Madison.

He founded and directs the Artificial Intelligence Research Laboratory in the Department of Computer Science at Iowa State University, where he is currently an Associate Professor. His research and teaching interests include artificial intelligence, artificial neural networks, machine learning, adaptive systems, bioinformatics and computational biology, evolutionary computing, grammatical inference, intelligent agents and multiagent systems, neural and cognitive modeling, distributed artificial intelligence, data mining and knowledge discovery, evolutionary robotics, parallel and distributed artificial intelligence, knowledge based systems, distributed knowledge networks, and applied artificial intelligence. He has published more than 80 research articles in journals, books, and conferences and has coedited three books.

Dr. Honavar is a coeditor-in-chief of the *Journal of Cognitive Systems Research*. He is a member of ACM, AAAI, and the New York Academy of Sciences.