

DistAl: An Inter-pattern Distance-based Constructive Learning Algorithm

Jihoon Yang, Rajesh Parekh and Vasant Honavar *

Artificial Intelligence Research Group

Department of Computer Science

226 Atanasoff Hall

Iowa State University

Ames, IA 50011. U.S.A.

{yang|parekh|honavar}@cs.iastate.edu

Abstract

Multi-layer networks of threshold logic units offer an attractive framework for the design of pattern classification systems. A new constructive neural network learning algorithm (**DistAl**) based on inter-pattern distance is introduced. **DistAl** constructs a single hidden layer of hyperspherical threshold neurons. Each neuron is designed to exclude a cluster of training patterns belonging to the same class. The weights and thresholds of the hidden neurons are determined directly by comparing the inter-pattern distances of the training patterns. This offers a significant advantage over other constructive learning algorithms that use an iterative (and often time consuming) weight modification strategy to train individual neurons. The individual clusters (represented by the hidden neurons) are combined by a single output layer of threshold neurons. The speed of **DistAl** makes it a good candidate for datamining and knowledge acquisition from very large datasets. The paper presents results of experiments using several artificial and real-world datasets. The results demonstrate that **DistAl** compares favorably with other neural network learning algorithms for pattern classification.

Keywords: neural networks, constructive learning algorithms, pattern classification

*This research was partially supported by the National Science Foundation (through grant IRI-9409580) and the John Deere Foundation.

1 Introduction

Trainable pattern classifiers find a broad range of applications in data mining and knowledge discovery [1, 2], intelligent agents [3, 4], diagnosis[5], computer vision [6], and automated knowledge acquisition [2, 7, 8, 9] from data. Multi-layer networks of threshold logic units (TLU) [10, 11, 12, 13, 14, 15] offer an attractive framework for the design of trainable pattern classification systems for a number of reasons including: potential for parallelism and fault and noise tolerance; significant representational and computational efficiency over disjunctive normal form (DNF) expressions and decision trees [11]; and simpler digital hardware implementations than their continuous counterparts such as sigmoid neurons used in networks trained with error backpropagation algorithm [16, 17].

A TLU implements an $(N - 1)$ -dimensional hyperplane which partitions N -dimensional Euclidean pattern space into two regions. A single TLU neural network is sufficient to classify patterns in two classes if they are *linearly separable*. A number of learning algorithms that are guaranteed to find a TLU weight setting that correctly classifies a linearly separable pattern set have been proposed in the literature [11, 18, 19, 20, 21, 22, 23, 24]. However, when the given set of patterns is not linearly separable, a multi-layer network of TLUs is needed to learn a complex decision boundary that is necessary to correctly classify the training examples.

Broadly speaking, there are two approaches to the design of multi-layer neural networks for pattern classification:

- *A-priori fixed topology* networks: the number of layers, the number of hidden neurons in each hidden layer, and the connections between each neuron are defined a-priori for each classification task. This is done on the basis of problem-specific knowledge (if available), or in ad hoc fashion (requiring a process of trial and error). Learning in such networks usually amounts to (typically error gradient guided) search for a suitable setting of numerical parameters, weights in a weight space defined by the choice of the

network topology.

- *Adaptive topology* networks: the topology of the target network is determined dynamically by introducing new neurons, layers, and connections in a controlled fashion using generative or constructive learning algorithms. In some cases, pruning mechanisms that discard redundant neurons and connections are used in conjunction with the network construction mechanisms [25, 26].

Constructive algorithms offer the following advantages over the conventional backpropagation style learning approaches [12, 15, 27, 28]:

- They obviate the need for an *ad-hoc, a-priori* choice of the network topology. Instead, they determine the network topology dynamically to give high chance of producing *optimal* (or minimal size) network.
- They are guaranteed to converge to zero classification errors on all finite and non-contradictory datasets.
- They use elementary threshold logic units (TLU) that are trained using the *perceptron* style weight update rules.
- They do not involve extensive parameter fine tuning.
- They provide a natural framework for exploiting problem-specific knowledge into the initial network configuration or heuristic knowledge (e.g., about the general topological constraints on the network) into the network construction algorithm [29].

Several constructive algorithms that incrementally construct networks of threshold neurons for 2-category pattern classification tasks have been proposed in the literature. These include the *tower* [30, 31], *pyramid* [31], *tiling* [32], *upstart* [33], *perceptron cascade* [34], and *sequential* [35]. Recently, provably correct extensions of these algorithms to handle multiple output classes and real-valued pattern attributes were proposed (see [12, 13, 14]). With

the exception of the sequential learning algorithm, these constructive learning algorithms are based on the idea of transforming the hard task of determining the necessary network topology and weights to two subtasks:

- Incremental addition of one or more threshold neurons to the network when the existing network topology fails to achieve the desired classification accuracy on the training set.
- Training the added threshold neuron(s) using some variant of the perceptron training algorithm (e.g., the pocket algorithm [11]) to improve the classification accuracy of the network.

In the case of the sequential learning algorithm, hidden neurons are added and trained by an appropriate weight training rule to exclude patterns belonging to the same class from the rest of the pattern set. The time-consuming, iterative nature of the perceptron training algorithm (though considerably faster than the corresponding error guided backpropagation training) often makes the use of such algorithms impractical for very large datasets (e.g., in largescale datamining and knowledge acquisition tasks), especially in applications where reasonably accurate classifiers have to be learned in almost real time. Similarly, hybrid learning systems that use neural network learning as the inner loop of a more complex optimization process (e.g., feature subset selection using a genetic algorithm where evaluation of fitness of a solution requires training a neural network based on a subset of input features represented by the solution and evaluating its classification accuracy [36, 37, 38]) call for a fast neural network training algorithm.

Instance-based learning (IBL) [39, 40, 41, 42] is an approach to learning in which the learning algorithm typically stores some or all of the training examples as prototypes. Each prototype is stored as an ordered pair (\mathbf{X}, c) where \mathbf{X} is a *pattern* represented in some chosen instance language (typically, in the form of a vector of attribute values), and c is the *class* to which \mathbf{X} belongs. Such a system, when used to classify a new pattern \mathbf{Y} , uses some *distance function* (e.g., Euclidean distance in the case of real-valued patterns) that computes

the distance of \mathbf{Y} from each stored prototype and predicts the classification of \mathbf{Y} using the known classification of the nearest prototype (or prototypes). Such algorithms, also referred to as *nearest neighbor* techniques have been investigated by researchers in pattern recognition [43, 44, 45], case-based reasoning [46, 47, 48], artificial neural networks [49], cognitive psychology [50, 51], and text classification [52]. Such distance-based techniques are also related to *radial basis function* networks [28, 53, 54, 55].

We present a new constructive neural network learning algorithm (DistAl), which can be viewed as a variant of the instance-based, nearest-neighbor, and radial-basis function-based approaches to pattern classification. DistAl replaces the iterative weight update of neurons that is typically used in constructive learning algorithms by a comparison of pair-wise distances among the training patterns. Since the inter-pattern distances are computed only once during the execution of the algorithm our approach achieves a significant speed advantage over other constructive learning algorithms.

The rest of the paper is organized as follows: Section 2 describes DistAl. Section 3 presents the results of various experiments designed to evaluate the performance of neural networks trained using DistAl on some benchmark classification problems. It also presents the results of experiments using DistAl in conjunction with a genetic algorithm-based approach to feature subset selection on several benchmark problems as well as a document classification task. Section 4 concludes with a summary and discussion of some directions for future research.

2 DistAl: A New Constructive Learning Algorithm

DistAl differs from other constructive learning algorithms mentioned above in two respects:

- It uses *spherical* threshold units – a variant of the TLU – as hidden neurons. The regions that are defined (or separated) by TLUs are unbounded. This motivates us to use spherical threshold units that cover locally bounded regions [8]. A spherical threshold neuron i has associated with it a weight vector \mathbf{W}_i , two thresholds – $\theta_{i,low}$ and

$\theta_{i,high}$, and a suitably defined distance metric d . It computes the distance $d(\mathbf{W}_i, \mathbf{X}^p)$ between a given input pattern \mathbf{X}^p and \mathbf{W}_i . The corresponding output $o_i^p = 1$ if $\theta_{i,low} \leq d(\mathbf{W}_i, \mathbf{X}^p) \leq \theta_{i,high}$ and 0 otherwise. The spherical neuron thus identifies a cluster of patterns that lie in the region between two concentric hyperspherical regions. \mathbf{W}_i represents the common center and $\theta_{i,low}$ and $\theta_{i,high}$ respectively represent the boundaries of the two regions.

- **DistAl** does not use an iterative algorithm for finding the weights and the thresholds. Instead, it computes the inter-pattern distance once between each pair of patterns in the training set and determines the weight values for hidden neurons by a greedy strategy (that attempts to correctly classify as many patterns as possible with the introduction of each new hidden neuron). The weights and thresholds are then set without the computationally expensive iterative process (see Section 2.2 for details).

The use of one-time inter-pattern distance calculation instead of (usually) iterative, expensive and time-consuming perceptron training procedure makes the proposed algorithm significantly faster than most other constructive learning algorithms. In fact, the time and space complexities of **DistAl** can be shown to be polynomial in the size of the training set (see Section 2.6 for details). This makes **DistAl** particularly well-suited for largescale datamining tasks.

2.1 Distance Metrics

Each hidden neuron introduced by **DistAl** essentially represents clusters of patterns that fall in the region bounded by two concentric hyperspherical regions in the pattern space. The weight vector of the neuron defines the center of the hyperspherical regions and the thresholds determine the boundaries of the regions (relative to the choice of the distance metric used). The choice of an appropriate distance metric for the hidden layer neurons is critical to achieving a good performance. Different distance metrics represent different

notions of *distance* in the pattern space. They also impose different *inductive biases* [7, 8] on the learning algorithm. Consequently, many researchers have investigated the use of alternative distance functions for instance-based learning [6, 44, 52, 56, 57]. The number and distribution of the clusters that result from specific choices of distance functions is a function of the distribution of the patterns as well as the clustering strategy used. Since it is difficult to identify the best distance metric in the absence of knowledge about the distribution of patterns in the pattern space, we chose to explore a number of different distance metrics proposed in the literature.

The distance between two patterns is often skewed by attributes that have high values. *Normalization* of individual attributes overcomes this problem in the distance computation. Normalization can be achieved by dividing each pattern attribute by the *range* of possible values for that attribute, or by 4 times the standard deviation for that attribute [57].

Normalization also allows attributes with nominal and/or missing values to be considered in distance computation. The distance for attributes with nominal values (say with attribute values x and y) is computed as follows [57]:

- *Overlap*: $d_{ol}(x, y) = 0$ if $x = y$; 1 otherwise.

- *Value difference*:

$$d_{vd}(x, y) = \sum_{c=1}^C \left| \frac{N_{a,x,c}}{N_{a,x}} - \frac{N_{a,y,c}}{N_{a,y}} \right|^q$$

where

- $N_{a,x}(N_{a,y})$: number of patterns in the training set that have value $x(y)$ for attribute a
- $N_{a,x,c}(N_{a,y,c})$: number of patterns in the training set that have value $x(y)$ for attribute a and output class c
- C : number of output classes
- q : a constant (Euclidean: 2, Manhattan: 1)

If there is a missing value in either of the patterns, the distance for that component (of the entire pattern vector) is taken to be 1.

Let $\mathbf{X}^p = [X_1^p, \dots, X_n^p]$ and $\mathbf{X}^q = [X_1^q, \dots, X_n^q]$ be two pattern vectors. Let max_i , min_i and σ_i be the maximum, minimum, and the standard deviation of values of the i th attribute of patterns in a dataset, respectively. Then the distance between \mathbf{X}^p and \mathbf{X}^q , for different choices of the distance metric d is defined as follows:

1. Range, value-difference based Euclidean:

$$\sqrt{\frac{1}{n} \sum_{i=1}^n \left[\left(\frac{X_i^p - X_i^q}{max_i - min_i} \right)^2 \text{ or } d_{vd}(X_i^p, X_i^q) \right]^2}$$

2. Range, value-difference based Manhattan:

$$\frac{1}{n} \sum_{i=1}^n \left[\frac{|X_i^p - X_i^q|}{max_i - min_i} \text{ or } d_{vd}(X_i^p, X_i^q) \right]$$

3. Range, value-difference based Maximum Value:

$$\max_i \left[\frac{|X_i^p - X_i^q|}{max_i - min_i} \text{ or } d_{vd}(X_i^p, X_i^q) \right]$$

Similarly, $4 * \sigma_i$ can be used instead of $max_i - min_i$ for standard deviation based metrics, and $d_{ol}(X_i^p, X_i^q)$ can be used instead of $d_{vd}(X_i^p, X_i^q)$ for overlap based metrics in above formulas.

4. Dice coefficient:

$$1 - \frac{2 \sum_{i=1}^n X_i^p X_i^q}{\sum_{i=1}^n (X_i^p)^2 + \sum_{i=1}^n (X_i^q)^2}$$

5. Cosine coefficient:

$$1 - \frac{\sum_{i=1}^n X_i^p X_i^q}{\sqrt{\sum_{i=1}^n (X_i^p)^2 \cdot \sum_{i=1}^n (X_i^q)^2}}$$

6. Jaccard coefficient:

$$1 - \frac{\sum_{i=1}^n X_i^p X_i^q}{\sum_{i=1}^n (X_i^p)^2 + \sum_{i=1}^n (X_i^q)^2 - \sum_{i=1}^n X_i^p X_i^q}$$

7. Canberra:

$$\sum_{i=1}^n \frac{|X_i^p - X_i^q|}{|X_i^p + X_i^q|}$$

Attribute based clustering:

Occasionally, the values of a single attribute between two bounds (say a_{lo} and a_{hi}) might exclusively identify patterns belonging to a particular output class. Thus, a hidden neuron that remembers the name of the attribute a and the two thresholds (a_{lo} and a_{hi}) can be used to form a cluster of patterns belonging to the same class. We use the attribute based comparison to obtain homogeneous clusters in conjunction with the inter-pattern distance based clustering.

2.2 Network Construction

Let $S = \{\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^N\}$ represents the N training patterns. **DistAl** calculates the pairwise inter-pattern distances for the training set (using the chosen distance metric d) and stores them in the distance matrix \mathcal{D} . Each row of \mathcal{D} is sorted in ascending order. Thus, row k of \mathcal{D} corresponds to the training pattern \mathbf{X}^k and the elements $\mathcal{D}[k, i]$ correspond to the distance of \mathbf{X}^k to the other training patterns. $\mathcal{D}[k, 0]$ is the distance to the closest pattern and $\mathcal{D}[k, N]$ is the distance to the farthest pattern from \mathbf{X}^k . Simultaneously, the attribute values of the training patterns are stored in \mathcal{D}' . \mathcal{D}' is essentially the entire training set with $\mathcal{D}'[k, i]$ representing the i th attribute value of the k th training pattern. Each column (attribute) of \mathcal{D}' is sorted in ascending order.

The key idea behind **DistAl** is to generate a single layer of hidden neurons each of which separates a subset of patterns in a training set using \mathcal{D} (or \mathcal{D}'). Then, they are fully connected to M output TLUs (1 for each output class) in an output layer. The representation of the patterns at the hidden layer is linearly separable [35]. Thus, an iterative perceptron learning rule can be used to train the output weights. However, the output weights can be directly set as follows: The weights between output and hidden neurons are chosen such that each hidden neuron overwhelms the effect of the hidden neurons generated later. If there are a

total of h hidden neurons (numbered $1, 2, \dots, h$ from left to right) then the weight between the output neuron j and the hidden neuron i is set to 2^{h-i} if the hidden neuron i excludes patterns belonging to class j and zero otherwise.

Let \mathbf{W}_l^h be the weights between the l th hidden neuron and inputs. Let \mathbf{W}_m^o be the weights between the output neuron for class m and hidden neurons, and W_{ml}^o be the weight between the output neuron for class m and the l th hidden neuron, respectively. The following pseudo-code summarizes the process of network construction:

Initialize the number of hidden neurons: $h = 0$;

while $S \neq \phi$

do

1. Double all existing weights (if any) between hidden and output neurons:

$$\mathbf{W}_m^o = \mathbf{W}_m^o * 2 \quad \forall m$$

2. Increment the number of hidden neurons: $h = h + 1$

3. Inter-pattern distance based:

Identify a row k of \mathcal{D} that excludes the largest subset of patterns in S that belong to the same class m as follows:

- (a) **For** each row $r = 1, \dots, N$ **do**

- i. Let i_r and j_r be column indices (corresponding to row r) for the matrix \mathcal{D} such that the patterns corresponding to the elements $\mathcal{D}[r, i_r], \mathcal{D}[r, i_r + 1], \dots, \mathcal{D}[r, j_r]$ all belong to the same class and also belong to S .

- ii. Let $c_r = j_r - i_r + 1$ (the number of patterns excluded).

- (b) Select k to be the one for which the corresponding c_k is the largest: $k = \arg \max_r c_r$

- (c) Let S_k be the corresponding set of patterns that are excluded by pattern \mathbf{X}^k , $d_{low}^k = \mathcal{D}[k, i_k]$ (distance to the closest pattern of the cluster) and $d_{high}^k = \mathcal{D}[k, j_k]$

(distance to the farthest pattern of the cluster).

4. Attribute based:

Analogously, using \mathcal{D}' identify an attribute a that excludes the largest number of patterns in S that belong to the same output class m (i.e., identify a for which c_a is the largest among all attributes.); Let S_a be the corresponding set of patterns from S that are excluded by attribute a , d_{low}^a and d_{high}^a be the minimum and maximum values respectively for attribute a among the patterns in set S_a .

5. **if** [Inter-pattern distance based] **then**

(a) Define a spherical threshold neuron with $\mathbf{W}^h = \mathbf{X}^k, \theta_{low} = d_{low}^k, \theta_{high} = d_{high}^k$.

(b) $S = S - S_k$

else

(a) Define a neuron corresponding to attribute a with $\theta_{low} = d_{low}^a, \theta_{high} = d_{high}^a$.

(b) $S = S - S_a$

6. Connect the new hidden neuron to output neurons: $W_{mh}^o = 1; W_{nh}^o = 0 \quad \forall n \neq m$

end while

2.3 Use of Network in Classification

The outputs in the output layer are computed by the *winner-take-all (WTA)* strategy. The output neuron m that has the highest net input produces 1 and all the other neurons produce 0's. The WTA strategy and the weight setting explained in Section 2.2 guarantee 100% training accuracy for any finite non-contradictory set of training patterns. (See Section 2.5 for detailed convergence proof).

The generalization accuracy of a test set is computed by the same way. Each test pattern is fed into the network and the outputs are computed by the WTA strategy. If there is one

or more hidden neurons that produce 1 (i.e., there exist one or more hidden neurons that include the test pattern within their thresholds), the outputs are computed by the WTA strategy in the output layer. Otherwise (i.e., all hidden neurons produce 0's and all output neurons produce 0's as well), the distance between the test pattern and the thresholds of each hidden neuron is computed. The hidden neuron that has the minimum distance is chosen to produce 1. Then the outputs are computed again in the output layer to compare with the desired classification.

2.4 Example

Although DistAI works on tasks with multi-category real-valued patterns, we will illustrate its operation using the simple XOR problem. We will assume the use of Manhattan distance metric. There are four training patterns ($S = \{\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}^3, \mathbf{X}^4\}$):

input	class
\mathbf{X}^1 : 0 0	A
\mathbf{X}^2 : 0 1	B
\mathbf{X}^3 : 1 0	B
\mathbf{X}^4 : 1 1	A

This yields the following distance matrix after sorted:

$$D = \begin{pmatrix} 0 & 1 & 1 & 2 \\ 0 & 1 & 1 & 2 \\ 0 & 1 & 1 & 2 \\ 0 & 1 & 1 & 2 \end{pmatrix}$$

The first row of the matrix is the distance of $\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}^3$ and \mathbf{X}^4 from pattern \mathbf{X}^1 . The second row of the matrix is the distance of $\mathbf{X}^2, \mathbf{X}^1, \mathbf{X}^4$ and \mathbf{X}^3 from \mathbf{X}^2 . The third row of the matrix is the distance of $\mathbf{X}^3, \mathbf{X}^1, \mathbf{X}^4$ and \mathbf{X}^2 from \mathbf{X}^3 . The last row of the matrix is the distance of $\mathbf{X}^4, \mathbf{X}^2, \mathbf{X}^3$ and \mathbf{X}^1 from \mathbf{X}^4 .

\mathbf{X}^1 excludes the maximum number of patterns from a single class (i.e., $S_k = \{\mathbf{X}^2, \mathbf{X}^3\}$, class = B). A hidden neuron is introduced for this cluster with $\mathbf{W}_1^h = [0 \ 0]$, $\theta_{low} = \theta_{high} = 1$, $W_{B1}^o = 1$, $W_{A1}^o = 0$. \mathbf{X}^2 and \mathbf{X}^3 are now eliminated from further consideration (i.e., $S = S - S_k = \{\mathbf{X}^1, \mathbf{X}^4\}$) The remaining patterns ($S_k = \{\mathbf{X}^1, \mathbf{X}^4\}$, class = A) can be excluded by any pattern (say, \mathbf{X}^1 again) with another hidden neuron with $\mathbf{W}_2^h = [0 \ 0]$, $\theta_{low} = 0$, $\theta_{high} = 2$, $W_{A2}^o = 1$, $W_{B2}^o = 0$, $W_{A1}^o = W_{A1}^o * 2 = 0$, $W_{B1}^o = W_{B1}^o * 2 = 2$. Now the algorithm stops since the entire training set is correctly classified (i.e., $S = S - S_k = \phi$). Figure 1 shows the network construction process.

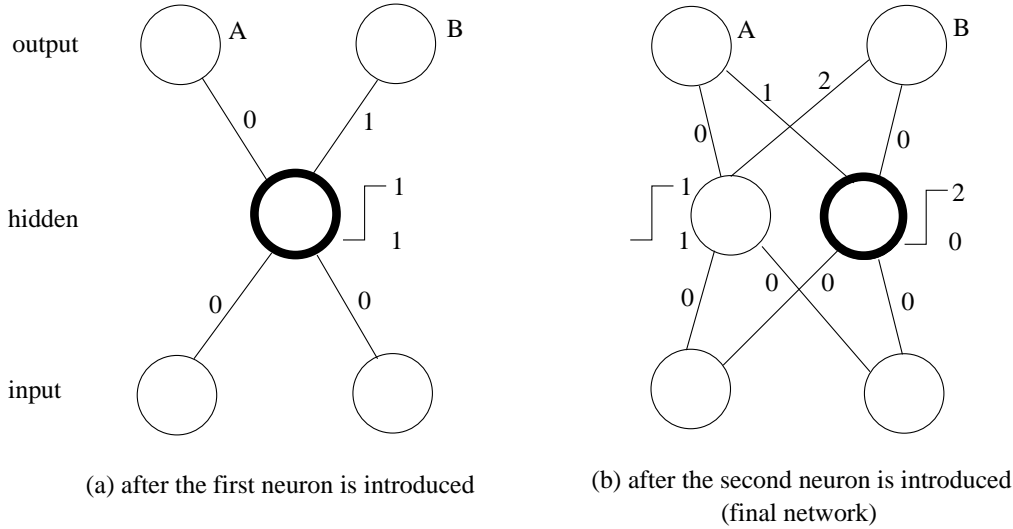


Figure 1: Process of Network Construction for the Example in DistAl

2.5 Convergence Proof

Theorem:

Given a finite non-contradictory set of training examples E , DistAl is guaranteed to converge to zero classification error after adding a finite number (h) of hidden neurons, where $h \leq |E|$. (In practice, $h \ll |E|$).

(Proof)

Let Z_i be the set of patterns that are excluded by i th hidden neuron. Each hidden neuron

finds the largest subset of patterns to be excluded. DistAl keeps introducing a hidden neuron until S becomes an empty set (i.e., $S = S - Z_i$). Since $S = \{\mathbf{X}^1, \dots, \mathbf{X}^N\}$ is the training set with the cardinality of N , $h = |Z_1, Z_2, \dots, Z_h| \leq N$ where Z_h is the last subset of patterns to be eliminated. It is clear that at least one pattern (\mathbf{X}^p) can be excluded by a new hidden neuron i with $\mathbf{W}_i^h = \mathbf{X}^p$ and 0 thresholds.¹ Since there are a finite number of patterns in the training set, and since each added hidden neuron is guaranteed to correctly classify a non-empty subset of the training set which is then eliminated from further consideration, no more than $|E|$ hidden neurons are needed.

The internal representation of the hidden layer for a pattern \mathbf{X}^p (which is a member of the i th cluster) has the form

$$\mathbf{H}^p = (0, 0, \dots, 0, 1, *, \dots, *) \quad (1)$$

(it has 0's in the first $i - 1$ hidden neurons, 1 in the i th hidden neuron and either 0 or 1 in the remaining hidden neurons) for a network with h hidden neurons. The weights from hidden to output neurons are set directly as explained in Section 2.2:

$$W_{ji}^o = \begin{cases} 2^{h-i} & \text{if } j \text{ is the right class of hidden neuron } i \\ 0 & \text{otherwise} \end{cases}$$

Consider a pattern \mathbf{X}^p which belongs to the subset Z_i of patterns excluded by the i th hidden neuron that represents the pattern \mathbf{X}^k . Let c_j be the classification of \mathbf{X}^k . Then $W_{ji}^o > W_{li}^o \ \forall j \neq l$. Also, the internal representation (1) guarantees the net input of output neuron j to be larger than that of any other output neuron. Consequently, \mathbf{X}^p is correctly classified in the output layer by the WTA strategy. As an example, assume $\mathbf{H}^p = (1, 1, 1)$ for a pattern \mathbf{X}^p belonging to class A, and the hidden neurons represent clusters for class A, B and B, respectively. Then, when \mathbf{X}^p is fed into input neurons, the net input to the

¹Note that this is not always true for maximum value distance metric and attribute-based approach. That is because there can be many patterns of different classifications that have the same maximum values/attributes values. Therefore, the convergence proof given here and the complexity analysis in Section 2.6 apply to distance-based approaches (excluding Maximum value metric), but not attribute-based approach.

output neuron for class A will be $2^{3-1} = 4$ and that to the output neuron for class B will be $2^{3-2} + 2^{3-3} = 3$. Thus, \mathbf{X}^p will be correctly classified as class A.

Therefore, DistAl is guaranteed to converge to zero classification error after adding a finite number of hidden neurons for a finite non-contradictory set of training examples. \square

2.6 Complexity Analysis

This section presents the complexity analysis for DistAl. The complexity analysis assumes that network construction is based on a single distance metric.

Let N_{pat} be the number of training patterns and N_{att} be the number of attributes in a dataset, respectively. Let N_{out} be the number of output neurons. Assume $N_{pat} > N_{att}$ and $N_{pat} \gg \max[N_{out}, h]$.

2.6.1 Time Complexity

Computing and sorting the distance matrix \mathcal{D} takes $\mathcal{O}(\max[N_{pat}^2 \cdot N_{att}, N_{pat}^2 \cdot \log N_{pat}])$.² Now, consider the pseudo-code given in Section 2.2. Step 1 takes $\mathcal{O}(N_{out} \cdot h)$. Step 2 takes $\mathcal{O}(1)$. Step 3 takes $\mathcal{O}(N_{pat}^2)$ because we need to go through the entire matrix \mathcal{D} to determine S_k .³ Step 5 takes $\mathcal{O}(N_{pat})$ to update S . Step 6 takes $\mathcal{O}(N_{out})$. Thus, the **while** loop takes $\mathcal{O}(N_{pat}^3)$ in the worst case. Therefore, the overall worst-case time complexity is $\mathcal{O}(N_{pat}^3)$. In practice, DistAl runs significantly faster than the worst-case time complexity because it eliminates a subset of elements from the original training set instead of a single pattern.

2.6.2 Space Complexity

The space requirement for the input patterns and their targets is $\mathcal{O}(N_{pat} \cdot [N_{att} + N_{out}])$. The weights require $\mathcal{O}(N_{out} \cdot h + h \cdot N_{in})$. The distance matrix requires $\mathcal{O}(N_{pat}^2)$. Thus, the total space complexity is $\mathcal{O}(N_{pat}^2)$.

²Computation of \mathcal{D}' in attribute-based approach takes only $\mathcal{O}(N_{att} \cdot N_{pat} \log N_{pat})$ because distance computation is not necessary.

³Step 4 is not considered here because it is used only with the attribute-based metric. The time required for step 4 is comparable to the time required for step 3.

2.7 Improving the Performance of DistAl Using Feature-Subset Selection

In pattern classification tasks, the choice of features (or attributes) used to represent patterns affect:

- *learning time*: The attributes used to describe the patterns implicitly determine the search space that needs to be explored by the learning algorithm. The larger the search space, the more time the learning algorithm needs for learning a sufficiently accurate classification function [7, 58].
- *number of examples needed*: All other things being equal, the larger the number of attributes used to describe the patterns, the larger is the number of examples need to learn a classification function to a desired accuracy [7, 58].
- *cost of classification*: In many real-world pattern classification tasks (e.g., medical diagnosis), some of the attributes may be observable symptoms and others might require diagnostic tests. Different diagnostic tests might have different costs as well as risks associated with them.

This presents us with a *feature subset selection problem* in automated design of pattern classifiers. The feature subset selection problem refers the task of identifying and selecting a useful subset of attributes to be used to represent patterns from a larger set of attributes. Satisfactory solution of this problem is particularly critical if instance-based, nearest-neighbor, or similarity-based learning algorithms like DistAl are used to build the classifier. This is due to the fact that such classifiers rely on the use of inter-pattern distances which are intricately linked to the choice of features used to represent the patterns. Presence of irrelevant or misleading features (e.g., social security numbers in a medical diagnosis task) can skew the distance calculation and hence adversely affect the generalization performance of the resulting classifier.

A detailed discussion of feature subset selection is beyond the scope of this paper. The interested reader is referred to [37, 38] for discussion of a variety of alternative approaches to feature subset selection. Since exhaustive search over all possible subsets of features is computationally infeasible, most approaches make restrictive assumptions (e.g., monotonicity – which states that the addition of features does not worsen classification accuracy) or use a variety of heuristics. Genetic algorithms [59, 60, 61] offer a particularly promising approach to feature subset selection for a number of reasons [36, 37, 38]:

- They do not have to rely on the often unrealistic monotonicity assumption.
- They are particularly effective tools for exploring large search spaces for near-optimal solutions [59, 60, 61].

The use of a genetic algorithm in any search or optimization problem requires:

- choice of a representation for encoding candidate solutions to be manipulated by the genetic algorithm
- definition of a fitness function that is used to evaluate the candidate solutions
- definition of a selection-scheme (e.g., fitness-proportionate selection)
- definition of suitable genetic operators that are used to transform candidate solutions (and thereby explore the search space)
- setting of user-controlled parameters (e.g., probability of applying a particular genetic operator, size of the population, etc.)

In our use of genetic algorithm for feature subset selection for `DistAl`, each candidate solution represented a subset of features used to encode patterns as input to `DistAl`. The fitness of the candidate solution was computed as the generalization accuracy (computed using a 10-fold cross-validation) of a classifier constructed using `DistAl`. Standard mutation and crossover operators were used on a fixed length binary vector representation of candidate

solutions (with a 1 indicating a selected feature). Experiments were run using the rank-based selection strategy with the following parameter settings: Population size is 50; Number of generation is 300; The probability of crossover is 0.5; The probability of mutation is 0.01; The probability of selection of the highest ranked individual is 0.6. (See [37, 38] for detailed explanations on the experiments).

3 Experimental Evaluation of DistAl

This section presents results of experiments using DistAl on several benchmark problems both with and without feature subset selection and compares them with the results presented by Wilson and Martinez in a recent paper [57]. It also presents the performance of DistAl on a real-world document classification task.

3.1 Datasets

Two artificial datasets (parity and two spirals) and a wide range of real-world datasets from the machine learning data repository at the University of California at Irvine [62] were chosen to test the performance of DistAl. DistAl is also used for classifying paper abstracts and news articles. The paper abstracts were chosen from three different sources: IEEE Expert magazine, Journal of Artificial Intelligence Research and Neural Computation. The news articles were obtained from Reuters dataset. Each document is represented in the form of a vector of numeric weights for each of the words (terms) in the vocabulary. The weights correspond to the term frequency and inverse document frequency (TFIDF) [63, 64] values for the corresponding words. The training sets for paper abstracts were generated based on the classification of the corresponding documents into two classes (interesting and not interesting) by two different individuals, resulting in two different data sets (**Abstract1** and **Abstract2**). The classifications for news articles were given based on their topics (6, 4 and 8 classes) following [65], resulting in three different datasets (**Reuters1, Reuters2** and

Reuters3), respectively. Table 1 summarizes the characteristics of the datasets selected for our experiments.

3.2 Experimental Results

DistAl is deterministic in the sense that its behavior is always identical for a given training set. Most other constructive learning algorithms are non-deterministic because their behavior is not always identical in different runs with the same training set and even with the same learning parameters due to the randomness in selecting initial weights, pattern presentations, and so on. Therefore, just one run of DistAl per dataset is sufficient to study the performance.

3.2.1 Parity Datasets

The seven, eight and nine-bit parity datasets (**P7**, **P8**, **P9**) were used to evaluate the performance of DistAl in terms of the network size. The Manhattan distance metric was used to train the entire set of patterns. Table 2 presents the size of the network generated by several algorithms. It shows that DistAl is capable of generating compact networks comparable to other algorithms for non-trivial tasks like the parity problem. Note that DistAl is also very fast. Since DistAl does not require iterative perceptron training procedure and keeps eliminating a subset of patterns that are not considered further in the learning process, it converges significantly fast. ⁴

3.2.2 Various Datasets from UCI Repository

DistAl was run once for each distance metric to compare the performance in terms of the generalization accuracy and the network size. A simple pruning technique was implemented to produce compact networks: When a new hidden neuron is introduced, the generalization

⁴It is not feasible to make a fair, thorough comparison of speeds of different algorithms. DistAl converged fairly quickly for almost all datasets. (See Section 2.6 for detailed analysis of time complexity). GA-MLP [66] is based on a genetic algorithm and thus it usually takes significant amount of time to get a quality solution. Cascade correlation [67] uses *Quickprop* [69]. *Quickprop* uses an iterative gradient descent method based on a second order heuristic.

accuracy of the network is computed. The current best generalization accuracy is stored in a *pocket* along with the network size. After the training is completed (i.e., 100% training accuracy is obtained) or no further training is possible (i.e., the limit of allowable hidden neurons (currently set to 100) is reached or no more patterns can be eliminated in Maximum value metric or attribute-based approach), the network with the best generalization accuracy in the pocket is restored by pruning the unnecessary hidden neurons.

A 10-fold cross-validation was performed for each dataset and its performance was shown in Table 3. The entries in the tables correspond to means and standard deviations and are shown in the form *mean* \pm *standard deviation*. An ‘*’ indicates that the distance computation was not possible (e.g., the denominator might be zero in Canberra metric) and a ‘-’ indicates that the distance metric was not applicable (e.g., Dice coefficient metric can not be used for nominal or missing values). As we can see from Table 3, no single distance metric outperformed other metrics on all datasets. That is because the performance depends on the distribution of the data. A distance metric might be appropriate for certain kinds of datasets while it might not for others.

It is impossible to do a thorough and fair comparison between various learning algorithms since each algorithm has its own *optimal* parameter settings which is usually unknown and not feasible to obtain within a reasonable amount of time. Also, the training and test sets that had been generated and used are not identical in general under the assumption that the experiments have been done a finite number of times. (An infinite number of experiments with random partitions of training and test sets from the same distributions of data can increase the confidence level). Following comparisons should be interpreted in light of those considerations. The best results of DistAl are compared with the best results in [57]. The results in [57] are chosen since they are recent and also obtained by a nearest-neighbor algorithm with a 10-fold cross-validation. Table 4 summarizes the comparison.

As we can see from Table 4, DistAl gave comparable results on most datasets (except **Soylarge** and **Vowel**). In case of **Vowel** dataset, the nearest neighbor algorithm [57] reports

a even higher accuracy than DistAl ⁵

The network size of three algorithms (*perceptron cascade* [34], *cascade correlation* [67], *upstart* [33]) for the two spirals problem is shown in [34]: 17.8 (*perceptron cascade*), 15.2 (*cascade correlation*), 91.4 (*upstart*). DistAl generated more compact networks with 7.7 hidden neurons.

Table 5 shows that the combination of DistAl and feature subset selection yield fairly good results. The results indicate that the networks constructed using GA-selected subset of features compare quite favorably with networks that use all of the features. In particular, feature subset selection resulted in significant improvement in generalization. For detailed explanation of implementation, related work and comparisons with other approaches see [37, 38].

3.2.3 Document Datasets

The same experimental setup was used as in Section 3.2.2. Table 6 shows that DistAl gives fairly good results for document classification as well. It gave reasonably high (over 80%) generalization accuracy for all datasets. Also, the GA-selected subset of features produced improved generalization accuracy with slightly larger network size. For detailed explanation of implementation, related work and comparisons with other approaches see [64].

4 Summary and Discussion

A fast inter-pattern distance-based constructive learning algorithm, DistAl, is introduced and its performance on a number of datasets is demonstrated. DistAl is different from other constructive learning algorithms in two aspects. First, it does not require an iterative perceptron style weight update rules for determining the connections between neurons. Instead, it computes the distance (using one of the pre-defined distance metrics) between each pattern pair and uses it to set the weights (and the thresholds) between hidden neurons and inputs.

⁵The best results reported in the literature [62] is 56% for **Vowel** dataset.

The weights between the hidden and output neurons are set using a one-shot (as opposed to iterative) learning algorithm. Thus, **DistAl** is relatively fast compared in comparison with most neural network training algorithms that rely on an iterative update of weights and consequently require multiple passes through the training set. Furthermore, **DistAl** is guaranteed to converge to 100% classification accuracy on any non-contradictory training set for most of the distance metrics used in this paper. Second, it generates a single hidden layer composed of *hyperspherical* threshold neurons instead of threshold logic units. Thus, the induced network can potentially discover natural clusters that exist in the data.

Despite its simplicity, experiments reported in this paper show that **DistAl** yields good performance on almost all real-world datasets that were considered. It also produced good performance on difficult artificial tasks such as parity and the two spirals data which have been used by numerous researchers for evaluation of supervised learning algorithms.

DistAl, because of its reliance on inter-pattern distances, is sensitive to the presence of irrelevant or misleading attributes in the pattern representation. Consequently, its classification accuracy can be further improved by incorporating a suitable feature subset selection algorithm. This is borne out by the experiments using **DistAl** in conjunction with a genetic algorithm for feature subset selection [37, 38].

A potential disadvantage of **DistAl** is its need for maintaining the inter-pattern distance matrix during learning. The memory needed to store this matrix grows quadratically with the size of the training set. This problem can be mitigated by freeing the memory for those patterns that are excluded by a new hidden neuron as learning progresses. It would be interesting to explore variants of **DistAl** that can avoid the need for maintaining the entire inter-pattern distance matrix during learning.

Because of its speed, **DistAl** is particularly well-suited to many real-world applications involving large amount of data and/or requesting real-time response such as largescale datamining and knowledge acquisition tasks and hybrid learning systems that use neural network learning as the inner loop of a more complex knowledge discovery process. An interesting

direction for future research is the design of versions of DistAl that can be used to incremental learning and assimilation of classification knowledge from multiple, distributed, dynamic data sources. Some preliminary results based on experiments using DistAl to design mobile agents for text classification and retrieval from distributed document collections are reported in [64].

Constructive algorithms in general provide an natural framework for exploration of cumulative (life long) learning [70] and for knowledge-based theory refinement [29, 71]. An interesting direction for future research would be to explore the use of DistAl for this task using real-world datasets e.g., the genome data used in [29].

References

- [1] U. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy. *Advances in Knowledge Discovery and Data Mining*. MIT Press, Cambridge, MA, 1996.
- [2] V. Honavar. Machine learning: Principles and applications. In J. Webster, editor, *Encyclopedia of Electrical and Electronics Engineering*. Wiley, New York, 1998. To appear.
- [3] J. Bradshaw. *Software Agents*. MIT Press, Cambridge, MA, 1997.
- [4] V. Honavar. Intelligent agents. In J. Williams and K. Sochats, editors, *Encyclopedia of Information Technology*. Marcel Dekker, New York, 1998. To appear.
- [5] K. Balakrishnan and V. Honavar. Intelligent diagnosis systems. *International Journal of Intelligent Systems*, 1998. In press.
- [6] R. Duda and P. Hart. *Pattern Classification and Scene Analysis*. Wiley, New York, 1973.
- [7] T. Mitchell. *Machine Learning*. McGraw Hill, New York, 1997.

- [8] P. Langley. *Elements of Machine Learning*. Morgan Kaufmann, Palo Alto, CA, 1995.
- [9] V. Honavar. Toward learning systems that integrate multiple strategies and representations. In V. Honavar and L. Uhr, editors, *Artificial Intelligence and Neural Networks: Steps Toward Principled Integration*, pages 615–644. Academic Press: New York., 1994.
- [10] C-H. Chen, R. Parekh, J. Yang, K. Balakrishnan, and V. Honavar. Analysis of decision boundaries generated by constructive neural network learning algorithms. In *Proceedings of WCNN'95, July 17-21, Washington D.C.*, volume 1, pages 628–635, 1995.
- [11] S. Gallant. *Neural Network Learning and Expert Systems*. MIT Press, Cambridge, MA, 1993.
- [12] R. Parekh, J. Yang, and V. Honavar. Constructive neural network learning algorithms for multi-category real-valued pattern classification. Technical Report ISU-CS-TR97-06, Department of Computer Science, Iowa State University, 1997. (Submitted for review to the IEEE Transactions on Neural Networks).
- [13] R. Parekh, J. Yang, and V. Honavar. **MUpstart** - a constructive neural network learning algorithm for multi-category pattern classification. In *Proceedings of the IEEE/INNS International Conference on Neural Networks, ICNN'97*, pages 1924–1929, 1997.
- [14] J. Yang, R. Parekh, and V. Honavar. **MTiling** - a constructive neural network learning algorithm for multi-category pattern classification. In *Proceedings of the World Congress on Neural Networks '96*, pages 182–187, San Diego, 1996.
- [15] V. Honavar. Structural learning. In J. Webster, editor, *Encyclopedia of Electrical and Electronics Engineering*. Wiley, New York, 1998. To appear.
- [16] D. Rumelhart, G. Hinton, and R. Williams. Learning internal representations by error propagation. In *Parallel Distributed Processing: Explorations into the Microstructure of Cognition*, volume 1 (Foundations). MIT Press, Cambridge, Massachusetts, 1986.

- [17] P. Werbos. *Beyond Regression: New Tools for Prediction and Analysis in Behavioral Sciences*. PhD thesis, Harvard University, 1974.
- [18] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65:386–408, 1958.
- [19] N. Nilsson. *The Mathematical Foundations of Learning Machines*. McGraw-Hill, New York, 1965.
- [20] W. Krauth and M. Mézard. Learning algorithms with optimal stability in neural networks. *J. Phys. A: Math. Gen.*, 20:L745–L752, 1987.
- [21] J. Anlauf and M. Biehl. Properties of an adaptive perceptron algorithm. In *Parallel Processing in Neural Systems and Computers*, pages 153–156. 1990.
- [22] M. Frean. *Small Nets and Short Paths: Optimizing Neural Computation*. PhD thesis, Center for Cognitive Science, Edinburgh University, UK, 1990.
- [23] H. Poulard. Barycentric correction procedure: A fast method of learning threshold units. In *Proceedings of WCNN'95, July 17-21, Washington D.C.*, volume 1, pages 710–713, 1995.
- [24] B. Raffin and M. Gordon. Learning and generalization with minimerror, a temperature-dependent learning algorithm. *Neural Computation*, 7:1206–1224, 1995.
- [25] R. Reed. Pruning algorithms — a survey. *IEEE Transactions on Neural Networks*, 4(5):740–747, 1993.
- [26] R. Parekh, J. Yang, and V. Honavar. Pruning strategies for constructive neural network learning algorithms. In *Proceedings of the IEEE/INNS International Conference on Neural Networks, ICNN'97*, pages 1960–1965, 1997.
- [27] V. Honavar. *Generative Learning Structures and Processes for Generalized Connectionist Networks*. PhD thesis, University of Wisconsin, Madison, 1990.

- [28] V. Honavar and Uhr. L. Generative learning structures for generalized connectionist networks. *Information Sciences*, 70(1-2):75–108, 1993.
- [29] R. Parekh and V. Honavar. Constructive theory refinement in knowledge based neural networks. In *Proceedings of the International Joint Conference on Neural Networks*, Anchorage, Alaska, 1998. To appear.
- [30] J. Nadal. Study of a growth algorithm for a feedforward network. *International Journal of Neural Systems*, 1(1):55–59, 1989.
- [31] S. Gallant. Perceptron based learning algorithms. *IEEE Transactions on Neural Networks*, 1(2):179–191, June 1990.
- [32] M. Mézard and J. Nadal. Learning feed-forward networks: The tiling algorithm. *J. Phys. A: Math. Gen.*, 22:2191–2203, 1989.
- [33] M. Freat. The upstart algorithm: A method for constructing and training feedforward neural networks. *Neural Computation*, 2:198–209, 1990.
- [34] N. Burgess. A constructive algorithm that converges for real-valued input patterns. *International Journal of Neural Systems*, 5(1):59–66, 1994.
- [35] M. Marchand, M. Golea, and P. Rujan. A convergence theorem for sequential learning in two-layer perceptrons. *Europhysics Letters*, 11(6):487–492, 1990.
- [36] J. Yang and V. Honavar. Feature subset selection using a genetic algorithm. In *Proceedings of the Genetic Programming Conference, GP'97*, pages 380–385, Stanford University, CA, 1997.
- [37] J. Yang and V. Honavar. Feature subset selection using a genetic algorithm. *IEEE Expert (Special Issue on Feature Transformation and Subset Selection)*, 1998. To appear.

- [38] J. Yang and V. Honavar. Feature subset selection using a genetic algorithm. In *Feature Extraction, Construction and Selection - A Data Mining Perspective*. Kluwer: New York, 1998. To appear.
- [39] D. Aha. Incremental constructive induction: An instance-based approach. In *Proceedings of the Eighth International Workshop on Machine Learning*, pages 117–121, Evanston, IL, 1991. Morgan Kaufmann.
- [40] D. Aha, D. Kibler, and M. Albert. Instance-based learning algorithms. *Machine Learning*, 6:37–66, 1991.
- [41] P. Turney. Theoretical analyses of cross-validation error and voting in instance-based learning. *Journal of Experimental and Theoretical Artificial Intelligence*, pages 331–360, 1994.
- [42] P. Domingos. Rule induction and instance-based learning: A unified approach. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI-95)*, 1995.
- [43] T. Cover and P. Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13:21–27, 1967.
- [44] E. Diday. Recent progress in distance and similarity measures in pattern recognition. In *Proceedings of the Second International Joint Conference on Pattern Recognition*, pages 534–539, 1974.
- [45] B. Dasarathy. *Nearest Neighbor (NN) Norms: NN Pattern Classification Techniques*. IEEE Computer Society Press, Los Alamitos, CA, 1991.
- [46] C. Stanfill and D. Waltz. Toward memory-based reasoning. *Communications of the ACM*, 29(12):1213–1228, 1986.

- [47] S. Cost and S. Salzberg. A weighted nearest neighbor algorithm for learning with symbolic features. *Machine Learning*, 10(1):57–78, 1993.
- [48] J. Kolodner. *Case-Based Reasoning*. Morgan Kaufmann, San Francisco, 1993.
- [49] G. Carpenter and S. Grossberg. *Pattern Recognition by Self-Organizing Neural Networks*. MIT Press, Cambridge, MA, 1991.
- [50] A. Tversky. Features of similarity. *Psychological Review*, 84:327–352, 1977.
- [51] R. Nosofsky. Attention, similarity, and the identification-categorization relationship. *Journal of Experimental Psychology: General*, 115:39–57, 1986.
- [52] G. Salton and M. McGill. *Introduction to Modern Information Retrieval*. McGraw Hill, New York, 1983.
- [53] D. Broomhead and D. Lowe. Multivariable functional interpolation and adaptive networks. *Complex Systems*, 2:321–355, 1988.
- [54] M. Powell. Radial basis functions for multivariable interpolation: A review. In J. Mason and M. Cox, editors, *Algorithms for Approximation*, pages 143–167. Clarendon Press, Oxford, 1987.
- [55] F. Girosi, M. Jones, and T. Poggio. Regularization theory and neural networks architectures. *Neural Computation*, 7:219–269, 1995.
- [56] B. Batchelor. *Pattern Recognition: Ideas in Practice*. Plenum Press, New York, 1978.
- [57] D. Wilson and T. Martinez. Improved heterogeneous distance functions. *Journal of Artificial Intelligence Research*, 6:1–34, 1997.
- [58] B. Natarajan. *Machine Learning: A Theoretical Approach*. Morgan Kauffman, San Mateo, CA, 1991.

- [59] D. Goldberg. *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley, New York, 1989.
- [60] M. Mitchell. *An Introduction to Genetic algorithms*. MIT Press, Cambridge, MA, 1996.
- [61] Z. Michalewicz. *Genetic Algorithms + Data Structures = Evolution Programs*. Springer-Verlag, New York, 3rd edition, 1996.
- [62] P. Murphy and D. Aha. Repository of machine learning databases. Department of Information and Computer Science, University of California, Irvine, CA, 1994.
- [63] G. Salton. Developments in automatic text retrieval. *Science*, 253:974–979, 1991.
- [64] J. Yang, P. Pai, V. Honavar, and L. Miller. Mobile intelligent agents for document classification and retrieval: A machine learning approach. In *14th European Meeting on Cybernetics and Systems Research. Symposium on Agent Theory to Agent Implementation*, Vienna, Austria, 1998.
- [65] D. Koller and M. Sahami. Hierarchically classifying documents using very few words. In *International Conference on Machine Learning*, pages 170–178, 1997.
- [66] H. Andersen and A. Tsoi. A constructive algorithm for the training of a multilayer perceptron based on the genetic algorithm. *Complex Systems*, 7:249–268, 1993.
- [67] S. Fahlman and C. Lebiere. The cascade correlation learning algorithm. In D. Touretzky, editor, *Neural Information Systems 2*, pages 524–532. Morgan-Kaufman, 1990.
- [68] M. Golea and M. Marchand. A growth algorithm for neural network decision trees. *Europhysics Letters*, 12(3):205–210, 1990.
- [69] S. Fahlman. Faster-learning variations on backpropagation: an empirical study. In D. Touretzky, G. Hinton, and T. Sejnowsky, editors, *Proceedings of the 1988 Connectionist Models Summer School*, pages 38–51. Morgan-Kaufman, 1988.

- [70] S. Thrun. Lifelong learning: A case study. Technical Report CMU-CS-95-208, Carnegie Mellon University, 1995.
- [71] J. W. Shavlik. A framework for combining symbolic and neural learning. In *Artificial Intelligence and Neural Networks: Steps Toward Principled Integration*. Academic Press, Boston, 1994.

Table 1: Datasets used in the experiments. **Size** is the number of patterns in the dataset, **Dimension** is the number of input attributes, **Missing?** is whether there are any missing values, and **Class** is the number of output classes.

Dataset	Size	Dimension	Attribute Type	Missing?	Class
7-bit parity (P7)	128	7	numeric	No	2
8-bit parity (P8)	256	8	numeric	No	2
9-bit parity (P9)	512	9	numeric	No	2
two spirals (2SP)	192	2	numeric	No	2
annealing database (Annealing)	798	38	numeric, nominal	Yes	5
audiology database (Audiology)	200	69	nominal	Yes	24
pittsburgh bridges (Bridges)	105	11	numeric, nominal	Yes	6
breast cancer (Cancer)	699	9	numeric	Yes	2
credit screening (CRX)	690	15	numeric, nominal	Yes	2
flag database (Flag)	194	28	numeric, nominal	No	8
glass identification (Glass)	214	9	numeric	No	6
heart disease (Heart)	270	13	numeric, nominal	No	2
heart disease [Cleveland](HeartCle)	303	13	numeric, nominal	Yes	2
heart disease [Hungarian](HeartHun)	294	13	numeric, nominal	Yes	2
heart disease [Long Beach](HeartLB)	200	13	numeric, nominal	Yes	2
heart disease [Swiss](HeartSwi)	123	13	numeric, nominal	Yes	2
hepatitis domain (Hepatitis)	155	19	numeric, nominal	Yes	2
horse colic (Horse)	300	22	numeric, nominal	Yes	2
ionosphere structure (Ionosphere)	351	34	numeric	No	2
iris plants (Iris)	150	4	numeric	No	3
liver disorders (Liver)	345	6	numeric	No	2
monks problems (Monks-1,2,3)	432	6	nominal	No	2
pima indians diabetes (Pima)	768	8	numeric	No	2
DNA sequences (Promoters)	106	57	nominal	No	2
sonar classification (Sonar)	208	60	numeric	No	2
large soybean (Soylarge)	307	35	nominal	Yes	19
small soybean (Soysmall)	47	35	nominal	No	4
vehicle silhouettes (Vehicle)	846	18	numeric	No	4
house votes (Votes)	435	16	nominal	Yes	2
vowel recognition (Vowel)	528	10	numeric	No	11
wine recognition (Wine)	178	13	numeric	No	3
zoo database (Zoo)	101	16	numeric, nominal	No	7
paper abstracts 1 (Abstract1)	100	790	numeric	No	2
paper abstracts 2 (Abstract2)	100	790	numeric	No	2
news articles 1 (Reuters1)	939	1568	numeric	No	6
news articles 2 (Reuters1)	139	435	numeric	No	4
news articles 3 (Reuters2)	834	1440	numeric	No	8

Table 2: Comparison of the network size generated by different algorithms for the parity datasets. A ‘-’ indicates that the result is not reported in the corresponding reference.

Algorithm	P7	P8	P9
DistAl	5	5	6
GA-MLP [66]	9	15	-
Perceptron cascade [34]	3	4	4
Cascade correlation [67]	4-5	5-6	-
Upstart [33]	6	7	8
Growth algorithm [68]	7	8	9
Sequential [35]	7	8	9
Tiling [32]	7	8	9
Tower [30]	3.5	4	4.5

Table 3: Results of various distance metrics (range, value-difference based Euclidean, Manhattan and Maximum value metrics). The best generalization accuracy among different distance metrics are shown in bold face.

Dataset	Euclidean [r,v]		Manhattan [r,v]		Maximum value [r,v]	
	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>
2SP	79.5 ± 10.1	8.7 ± 1.5	72.1 ± 5.8	8.5 ± 3.7	70.5 ± 8.9	10.1 ± 3.7
Annealing	96.6 ± 2.0	12.1 ± 2.4	93.3 ± 2.8	15.5 ± 3.4	0.0 ± 0.0	0.0 ± 0.0
Audiology	66.0 ± 9.7	24.7 ± 4.8	59.0 ± 8.6	26.7 ± 3.0	1.0 ± 2.0	0.6 ± 1.2
Bridge	55.0 ± 15.6	3.5 ± 2.1	59.0 ± 10.4	3.2 ± 2.6	45.0 ± 16.9	3.4 ± 1.5
Cancer	97.8 ± 1.2	2.9 ± 1.2	97.5 ± 1.7	3.9 ± 1.4	95.1 ± 1.5	6.2 ± 5.1
CRX	87.7 ± 3.3	7.7 ± 6.9	87.5 ± 3.8	7.3 ± 4.3	86.4 ± 3.3	7.7 ± 3.5
Flag	63.7 ± 8.0	5.7 ± 3.2	64.7 ± 11.3	6.0 ± 3.6	57.9 ± 5.8	7.0 ± 2.7
Glass	70.5 ± 8.5	9.8 ± 6.9	66.2 ± 4.5	9.9 ± 6.5	67.6 ± 7.6	10.1 ± 5.6
Heart	83.7 ± 5.3	3.3 ± 1.8	84.8 ± 4.8	5.7 ± 3.6	85.2 ± 3.3	7.6 ± 4.5
HeartCle	85.3 ± 7.2	4.6 ± 3.8	85.3 ± 3.4	6.0 ± 2.9	82.3 ± 4.5	10.9 ± 6.7
HeartHun	84.5 ± 5.8	6.7 ± 2.8	84.8 ± 5.6	6.6 ± 2.9	0.0 ± 0.0	0.0 ± 0.0
HeartLB	78.5 ± 9.2	5.0 ± 3.5	77.5 ± 6.8	4.9 ± 3.0	0.0 ± 0.0	0.0 ± 0.0
HeartSwi	93.3 ± 3.3	2.0 ± 0.0	93.3 ± 5.0	2.2 ± 0.8	0.0 ± 0.0	0.0 ± 0.0
Hepatitis	83.3 ± 4.5	3.0 ± 1.3	83.3 ± 6.2	2.5 ± 0.8	79.3 ± 8.7	2.0 ± 0.0
Horse	86.0 ± 3.6	5.3 ± 4.5	84.7 ± 4.3	5.1 ± 3.2	63.7 ± 5.9	2.0 ± 0.0
Ionosphere	93.1 ± 4.5	6.8 ± 1.4	90.0 ± 5.8	5.8 ± 2.1	91.4 ± 4.2	5.5 ± 1.7
Iris	93.1 ± 4.5	6.8 ± 1.4	90.0 ± 5.8	5.8 ± 2.1	91.4 ± 4.2	5.5 ± 1.7
Liver	67.7 ± 6.8	7.8 ± 4.5	63.5 ± 8.2	6.4 ± 6.8	67.4 ± 5.3	7.1 ± 3.7
Monks-1	90.0 ± 9.3	7.4 ± 4.4	89.1 ± 7.7	7.4 ± 5.1	82.8 ± 9.1	9.6 ± 4.9
Monks-2	79.8 ± 10.4	8.4 ± 4.5	79.5 ± 10.4	13.0 ± 9.5	82.8 ± 9.4	6.4 ± 2.9
Monks-3	99.1 ± 1.5	3.0 ± 0.0	98.6 ± 1.9	3.0 ± 0.6	98.6 ± 1.9	2.1 ± 0.3
Pima	74.3 ± 3.2	9.5 ± 6.7	73.4 ± 4.0	13.2 ± 7.8	73.7 ± 5.2	8.3 ± 5.0
Promoters	87.0 ± 11.0	2.8 ± 0.4	88.0 ± 7.5	2.2 ± 0.4	85.0 ± 8.1	2.8 ± 0.6
Sonar	83.0 ± 7.8	6.4 ± 2.7	81.5 ± 9.5	4.8 ± 2.4	78.5 ± 8.1	7.5 ± 3.8
Soylarge	81.0 ± 5.6	20.2 ± 3.2	74.3 ± 9.3	21.6 ± 5.0	67.7 ± 4.5	16.7 ± 2.4
Soysmall	90.0 ± 16.6	3.4 ± 0.5	92.5 ± 16.0	3.6 ± 0.5	97.5 ± 7.5	3.6 ± 0.5
Vehicle	64.1 ± 6.5	29.5 ± 13.3	61.7 ± 3.2	25.9 ± 18.3	57.0 ± 4.7	49.4 ± 22.2
Votes	96.1 ± 1.5	3.2 ± 1.5	95.4 ± 2.3	3.7 ± 1.2	78.8 ± 8.1	3.6 ± 1.4
Vowel	65.2 ± 6.9	34.6 ± 8.5	65.8 ± 6.4	40.7 ± 8.5	61.7 ± 8.3	39.5 ± 7.7
Wine	92.9 ± 5.8	4.3 ± 0.8	92.9 ± 5.8	4.1 ± 0.7	94.1 ± 6.4	4.7 ± 0.6
Zoo	96.0 ± 4.9	6.1 ± 1.1	96.0 ± 8.0	6.1 ± 0.9	93.9 ± 4.6	6.0 ± 1.2

Table 3: Results of various distance metrics (standard deviation, value-difference based Euclidean, Manhattan and Maximum value metrics).

Dataset	Euclidean [s,v]		Manhattan [s,v]		Maximum value [s,v]	
	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>
2SP	83.7 \pm 7.6	7.7 \pm 1.8	69.5 \pm 6.1	7.1 \pm 3.9	72.1 \pm 6.7	8.8 \pm 3.1
Annealing	96.3 \pm 1.4	10.6 \pm 2.8	93.9 \pm 2.3	13.7 \pm 3.4	0.0 \pm 0.0	0.0 \pm 0.0
Audiology	66.0 \pm 9.7	24.7 \pm 4.8	59.0 \pm 8.6	26.7 \pm 3.0	1.0 \pm 2.0	0.6 \pm 1.2
Bridge	56.0 \pm 17.4	4.0 \pm 3.6	59.0 \pm 13.0	3.4 \pm 2.7	52.0 \pm 14.7	4.0 \pm 2.8
Cancer	96.8 \pm 2.0	4.0 \pm 1.6	96.8 \pm 1.9	4.5 \pm 2.6	95.4 \pm 1.7	10.4 \pm 4.4
CRX	87.4 \pm 3.6	7.2 \pm 3.7	87.0 \pm 4.1	7.0 \pm 4.5	86.4 \pm 4.5	6.1 \pm 4.6
Flag	60.5 \pm 8.2	6.4 \pm 4.6	65.8 \pm 9.5	9.1 \pm 6.2	55.3 \pm 10.3	11.1 \pm 9.2
Glass	68.1 \pm 7.7	11.5 \pm 7.7	66.2 \pm 5.8	7.3 \pm 3.6	69.5 \pm 6.8	9.5 \pm 7.7
Heart	82.6 \pm 5.0	3.6 \pm 1.6	85.6 \pm 5.1	4.7 \pm 3.1	81.5 \pm 6.2	7.7 \pm 5.9
HeartCle	81.7 \pm 4.8	3.9 \pm 2.7	83.7 \pm 4.6	4.6 \pm 3.3	83.7 \pm 6.7	5.6 \pm 4.5
HeartHun	84.8 \pm 6.4	7.0 \pm 4.1	83.1 \pm 4.7	5.5 \pm 3.4	0.0 \pm 0.0	0.0 \pm 0.0
HeartLB	76.5 \pm 8.7	3.3 \pm 2.8	77.5 \pm 6.8	3.9 \pm 2.5	0.0 \pm 0.0	0.0 \pm 0.0
HeartSwi	94.2 \pm 3.8	2.2 \pm 0.6	94.2 \pm 3.8	2.3 \pm 0.9	0.0 \pm 0.0	0.0 \pm 0.0
Hepatitis	84.7 \pm 9.5	6.2 \pm 4.0	84.0 \pm 6.8	4.6 \pm 3.1	84.7 \pm 9.5	4.2 \pm 1.8
Horse	83.0 \pm 5.7	4.4 \pm 3.5	83.7 \pm 6.9	7.9 \pm 4.4	63.7 \pm 5.9	2.0 \pm 0.0
Ionosphere	92.9 \pm 5.5	6.9 \pm 2.1	91.4 \pm 5.9	5.8 \pm 1.4	92.6 \pm 4.1	5.3 \pm 1.9
Iris	92.9 \pm 5.5	6.9 \pm 2.1	91.3 \pm 6.0	3.3 \pm 1.0	94.0 \pm 6.3	4.0 \pm 1.2
Liver	66.5 \pm 5.1	9.9 \pm 5.7	64.4 \pm 4.6	10.6 \pm 8.3	66.2 \pm 6.2	11.3 \pm 7.6
Monks-1	90.0 \pm 9.3	7.4 \pm 4.4	89.1 \pm 7.7	7.4 \pm 5.1	82.8 \pm 9.1	9.6 \pm 4.9
Monks-2	79.8 \pm 10.4	8.4 \pm 4.5	79.5 \pm 10.4	13.0 \pm 9.5	82.8 \pm 9.4	6.4 \pm 2.9
Monks-3	99.1 \pm 1.5	3.0 \pm 0.0	98.6 \pm 1.9	3.0 \pm 0.6	98.6 \pm 1.9	2.1 \pm 0.3
Pima	74.2 \pm 3.9	10.0 \pm 3.8	76.3 \pm 5.1	8.1 \pm 4.9	74.7 \pm 4.0	13.1 \pm 10.8
Promoters	87.0 \pm 11.0	2.8 \pm 0.4	88.0 \pm 7.5	2.2 \pm 0.4	85.0 \pm 8.1	2.8 \pm 0.6
Sonar	73.5 \pm 7.4	4.6 \pm 3.1	78.5 \pm 8.1	4.8 \pm 2.1	73.5 \pm 7.4	5.6 \pm 3.6
Soylarge	81.0 \pm 5.6	20.2 \pm 3.2	74.3 \pm 9.3	21.6 \pm 5.0	67.7 \pm 4.5	16.7 \pm 2.4
Soysmall	90.0 \pm 16.6	3.4 \pm 0.5	92.5 \pm 16.0	3.6 \pm 0.5	97.5 \pm 7.5	3.6 \pm 0.5
Vehicle	65.1 \pm 4.0	24.4 \pm 9.6	65.4 \pm 3.5	23.7 \pm 5.0	62.1 \pm 4.7	52.9 \pm 18.6
Votes	96.1 \pm 1.5	3.2 \pm 1.5	95.4 \pm 2.3	3.7 \pm 1.2	78.8 \pm 8.1	3.6 \pm 1.4
Vowel	66.7 \pm 7.5	31.2 \pm 10.1	65.0 \pm 7.7	36.3 \pm 8.5	57.9 \pm 8.8	39.2 \pm 13.5
Wine	95.9 \pm 4.6	4.9 \pm 0.3	92.4 \pm 7.5	4.4 \pm 0.7	92.9 \pm 5.8	4.4 \pm 0.9
Zoo	96.0 \pm 4.9	6.1 \pm 1.1	96.0 \pm 8.0	6.1 \pm 0.9	93.0 \pm 4.6	6.0 \pm 1.2

Table 3: Results of various distance metrics (range, overlap based Euclidean, Manhattan and Maximum value metrics).

Dataset	Euclidean [r,o]		Manhattan [r,o]		Maximum value [r,o]	
	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>
2SP	79.5 ± 10.1	8.7 ± 1.5	72.1 ± 5.8	8.5 ± 3.7	70.5 ± 8.9	10.1 ± 3.7
Annealing	94.7 ± 1.8	14.6 ± 3.6	93.2 ± 2.5	15.6 ± 5.0	0.0 ± 0.0	0.0 ± 0.0
Audiology	66.0 ± 10.0	27.3 ± 7.4	66.0 ± 10.0	27.3 ± 7.4	1.0 ± 2.0	0.6 ± 1.2
Bridge	60.0 ± 10.0	4.8 ± 3.0	60.0 ± 10.0	6.3 ± 5.0	36.0 ± 18.0	10.0 ± 17.1
Cancer	97.8 ± 1.2	2.9 ± 1.2	97.5 ± 1.7	3.9 ± 1.4	95.1 ± 1.5	6.2 ± 5.1
CRX	83.8 ± 5.3	9.7 ± 5.1	83.9 ± 3.8	9.1 ± 4.4	61.9 ± 7.9	56.4 ± 20.7
Flag	47.4 ± 7.1	6.9 ± 5.4	50.5 ± 10.8	10.3 ± 6.4	24.7 ± 9.4	2.1 ± 0.3
Glass	70.5 ± 8.5	9.8 ± 6.9	66.2 ± 4.5	9.9 ± 6.5	67.6 ± 7.6	10.1 ± 5.6
Heart	86.7 ± 7.6	5.7 ± 4.4	86.3 ± 5.8	4.1 ± 3.0	73.3 ± 4.9	26.5 ± 25.0
HeartCle	83.0 ± 5.5	4.9 ± 2.7	85.3 ± 2.7	3.4 ± 1.1	71.7 ± 10.3	23.6 ± 15.3
HeartHun	85.9 ± 6.3	5.0 ± 2.9	84.8 ± 3.8	4.5 ± 3.0	0.0 ± 0.0	0.0 ± 0.0
HeartLB	77.0 ± 9.8	3.4 ± 2.5	80.0 ± 7.4	5.1 ± 2.6	0.0 ± 0.0	0.0 ± 0.0
HeartSwi	94.2 ± 3.8	2.3 ± 0.9	94.2 ± 3.8	2.3 ± 0.9	0.0 ± 0.0	0.0 ± 0.0
Hepatitis	83.3 ± 4.5	3.0 ± 1.3	83.3 ± 6.2	2.5 ± 0.8	79.3 ± 8.7	2.0 ± 0.0
Horse	84.0 ± 6.3	5.6 ± 2.2	85.7 ± 7.9	4.1 ± 2.5	63.7 ± 5.9	2.0 ± 0.0
Ionosphere	93.1 ± 4.5	6.8 ± 1.4	90.0 ± 5.8	5.8 ± 2.1	91.4 ± 4.2	5.5 ± 1.7
Iris	96.0 ± 4.4	3.4 ± 0.7	96.0 ± 3.3	3.4 ± 0.7	91.3 ± 6.7	3.3 ± 0.5
Liver	67.7 ± 6.8	7.8 ± 4.5	63.5 ± 8.2	6.4 ± 6.8	67.4 ± 5.3	7.1 ± 3.7
Monks-1	90.9 ± 7.1	26.9 ± 7.5	90.9 ± 7.1	26.9 ± 7.5	49.3 ± 7.1	2.0 ± 0.0
Monks-2	100 ± 0.0	2.7 ± 2.1	100 ± 0.0	2.7 ± 2.1	33.0 ± 4.3	2.0 ± 0.0
Monks-3	91.6 ± 4.4	16.2 ± 4.4	91.6 ± 4.4	16.2 ± 4.4	49.3 ± 6.6	2.0 ± 0.0
Pima	74.3 ± 3.2	9.5 ± 6.7	73.4 ± 4.0	13.2 ± 7.8	73.7 ± 5.2	8.3 ± 5.0
Promoters	83.0 ± 6.4	3.4 ± 1.4	83.0 ± 6.4	3.4 ± 1.4	56.0 ± 6.6	20.0 ± 36.0
Sonar	83.0 ± 7.8	6.4 ± 2.7	81.5 ± 9.5	4.8 ± 2.4	76.0 ± 9.2	7.5 ± 3.8
Soylarge	75.0 ± 5.2	26.3 ± 4.7	75.0 ± 5.2	26.3 ± 4.7	12.3 ± 6.8	2.0 ± 0.0
Soysmall	97.5 ± 7.5	3.9 ± 0.3	97.5 ± 7.5	0.9 ± 0.3	30.0 ± 21.8	13.3 ± 18.2
Vehicle	64.1 ± 6.5	29.5 ± 13.3	61.7 ± 3.2	25.9 ± 18.3	57.0 ± 4.7	49.4 ± 22.2
Votes	95.6 ± 2.6	6.1 ± 2.3	95.6 ± 2.6	6.1 ± 2.3	47.0 ± 8.1	42.5 ± 29.2
Vowel	65.2 ± 6.9	34.6 ± 8.5	65.8 ± 6.4	40.7 ± 8.5	61.7 ± 8.3	39.5 ± 7.7
Wine	92.9 ± 5.8	4.3 ± 0.8	92.9 ± 5.8	4.1 ± 0.7	94.1 ± 6.4	4.7 ± 0.6
Zoo	92.0 ± 7.5	6.2 ± 0.9	92.0 ± 7.5	6.2 ± 0.9	75.0 ± 12.9	33.4 ± 17.4

Table 3: Results of various distance metrics (standard deviation, overlap based Euclidean, Manhattan and Maximum value metrics).

Dataset	Euclidean [s,o]		Manhattan [s,o]		Maximum value [s,o]	
	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>
2SP	83.7 ± 7.6	7.7 ± 1.8	69.5 ± 6.1	7.1 ± 3.9	72.1 ± 6.7	8.8 ± 3.1
Annealing	95.2 ± 1.9	14.5 ± 2.4	94.8 ± 3.0	17.5 ± 2.5	0.0 ± 0.0	0.0 ± 0.0
Audiology	66.0 ± 10.0	27.3 ± 7.4	66.0 ± 10.0	27.3 ± 7.4	1.0 ± 2.0	0.6 ± 1.2
Bridge	63.0 ± 7.8	5.2 ± 3.3	60.0 ± 7.8	4.2 ± 2.7	38.0 ± 14.7	15.3 ± 21.9
Cancer	96.8 ± 2.0	4.0 ± 1.6	96.8 ± 1.9	4.5 ± 2.6	95.4 ± 1.7	10.4 ± 4.4
CRX	85.2 ± 5.6	10.5 ± 5.4	84.9 ± 6.3	9.1 ± 6.3	58.7 ± 6.5	44.0 ± 24.0
Flag	46.8 ± 7.2	7.2 ± 4.4	51.1 ± 8.8	9.1 ± 8.2	31.6 ± 10.8	3.6 ± 1.4
Glass	68.1 ± 7.7	11.5 ± 7.7	66.2 ± 5.8	7.3 ± 3.6	69.5 ± 6.8	9.5 ± 7.7
Heart	85.9 ± 6.4	5.5 ± 3.3	85.6 ± 4.5	5.2 ± 3.6	71.1 ± 5.2	26.4 ± 17.0
HeartCle	82.0 ± 4.5	3.9 ± 2.4	83.3 ± 7.0	5.4 ± 4.2	67.0 ± 7.1	26.7 ± 17.5
HeartHun	82.1 ± 4.8	5.4 ± 4.4	85.5 ± 4.6	5.2 ± 3.7	0.0 ± 0.0	0.0 ± 0.0
HeartLB	77.0 ± 7.5	4.3 ± 3.4	79.0 ± 6.2	5.9 ± 4.4	0.0 ± 0.0	0.0 ± 0.0
HeartSwi	94.2 ± 3.8	2.2 ± 0.6	94.2 ± 3.8	2.2 ± 0.6	0.0 ± 0.0	0.0 ± 0.0
Hepatitis	84.7 ± 9.5	6.2 ± 4.0	84.0 ± 6.8	4.6 ± 3.1	84.7 ± 9.5	4.2 ± 1.8
Horse	80.0 ± 5.2	10.1 ± 6.2	84.3 ± 4.2	4.5 ± 2.0	63.7 ± 5.9	2.0 ± 0.0
Ionosphere	92.9 ± 5.5	6.9 ± 2.1	91.4 ± 5.9	5.8 ± 1.4	92.6 ± 4.1	5.3 ± 1.9
Iris	94.0 ± 3.6	3.8 ± 1.2	91.3 ± 6.0	3.3 ± 1.0	94.0 ± 6.3	4.0 ± 1.2
Liver	66.5 ± 5.1	9.9 ± 5.7	64.4 ± 4.6	10.6 ± 8.3	66.2 ± 6.2	11.3 ± 7.6
Monks-1	90.9 ± 7.1	26.9 ± 7.5	90.9 ± 7.1	26.9 ± 7.5	49.3 ± 7.1	2.0 ± 0.0
Monks-2	100 ± 0.0	2.7 ± 2.1	100 ± 0.0	2.7 ± 2.1	33.1 ± 0.5	2.0 ± 0.0
Monks-3	91.6 ± 4.4	16.2 ± 4.4	91.6 ± 4.4	16.2 ± 4.4	49.3 ± 6.6	2.0 ± 0.0
Pima	74.2 ± 3.9	10.0 ± 3.8	76.3 ± 5.1	8.1 ± 4.9	74.7 ± 4.0	13.1 ± 10.8
Promoters	83.0 ± 6.4	3.4 ± 1.4	83.0 ± 6.4	3.4 ± 1.4	56.0 ± 6.6	20.0 ± 36.0
Sonar	82.0 ± 6.8	4.6 ± 3.1	78.5 ± 8.9	4.8 ± 2.1	73.5 ± 7.4	5.6 ± 3.6
Soylarge	75.0 ± 5.2	26.3 ± 4.7	75.0 ± 5.2	26.3 ± 4.7	12.3 ± 6.9	2.0 ± 0.0
Soysmall	97.5 ± 7.5	3.9 ± 0.3	97.5 ± 7.5	3.9 ± 0.3	30.0 ± 21.8	13.3 ± 18.2
Vehicle	65.1 ± 4.0	24.4 ± 9.6	65.4 ± 3.5	23.7 ± 5.0	62.1 ± 4.7	52.9 ± 18.6
Votes	95.6 ± 2.6	6.1 ± 2.3	95.6 ± 2.6	6.1 ± 2.3	47.0 ± 8.1	42.5 ± 29.2
Vowel	66.7 ± 7.5	31.2 ± 10.1	65.0 ± 7.7	36.3 ± 8.5	57.9 ± 8.8	39.2 ± 13.5
Wine	95.9 ± 4.6	4.9 ± 0.3	92.4 ± 7.5	4.4 ± 0.7	92.9 ± 5.8	4.4 ± 0.9
Zoo	92.0 ± 7.5	6.2 ± 0.9	92.0 ± 7.5	6.2 ± 0.9	75.0 ± 12.9	33.4 ± 17.4

Table 3: Results of various distance metrics (Dice, Cosine and Jaccard coefficient metrics).

Dataset	Dice coefficient		Cosine coefficient		Jaccard coefficient	
	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>
2SP	56.8 ± 8.4	4.7 ± 2.7	56.8 ± 7.4	6.4 ± 10.0	55.3 ± 5.9	5.5 ± 3.2
Annealing	-	-	-	-	-	-
Audiology	-	-	-	-	-	-
Bridge	-	-	-	-	-	-
Cancer	-	-	-	-	-	-
CRX	-	-	-	-	-	-
Flag	-	-	-	-	-	-
Glass	66.2 ± 8.9	7.9 ± 4.6	68.6 ± 5.7	11.0 ± 5.4	66.2 ± 8.9	7.9 ± 4.6
Heart	-	-	-	-	-	-
HeartCle	-	-	-	-	-	-
HeartHun	-	-	-	-	-	-
HeartLB	-	-	-	-	-	-
HeartSwi	-	-	-	-	-	-
Hepatitis	-	-	-	-	-	-
Horse	-	-	-	-	-	-
Ionosphere	92.6 ± 3.9	5.4 ± 1.4	94.3 ± 5.0	5.5 ± 1.6	92.9 ± 3.7	5.8 ± 1.5
Iris	95.3 ± 6.7	3.1 ± 0.5	97.3 ± 3.3	4.0 ± 0.0	95.3 ± 6.7	3.1 ± 0.5
Liver	66.8 ± 5.8	6.5 ± 5.8	70.6 ± 6.2	6.9 ± 3.7	65.9 ± 5.3	6.2 ± 5.5
Monks-1	-	-	-	-	-	-
Monks-2	-	-	-	-	-	-
Monks-3	-	-	-	-	-	-
Pima	71.6 ± 2.8	13.0 ± 10.2	68.2 ± 5.9	8.0 ± 8.1	72.4 ± 3.0	11.4 ± 8.7
Promoters	-	-	-	-	-	-
Sonar	79.5 ± 7.2	8.2 ± 3.0	76.5 ± 8.1	4.8 ± 2.1	79.0 ± 7.7	6.9 ± 2.3
Soylarge	-	-	-	-	-	-
Soysmall	-	-	-	-	-	-
Vehicle	58.8 ± 3.9	20.2 ± 20.7	61.0 ± 3.3	17.3 ± 7.5	58.7 ± 3.9	20.1 ± 20.8
Votes	-	-	-	-	-	-
Vowel	69.8 ± 6.4	38.0 ± 8.3	57.3 ± 6.1	35.7 ± 12.1	69.6 ± 7.4	38.1 ± 8.6
Wine	94.3 ± 3.4	6.0 ± 3.6	83.5 ± 6.3	6.7 ± 4.0	80.6 ± 6.5	6.3 ± 4.1
Zoo	-	-	-	-	-	-

Table 3: Results of various distance metrics (Camberra and Attribute-based metrics).

Dataset	Camberra		Attribute-based	
	<i>Accuracy</i>	<i>Hidden</i>	<i>Accuracy</i>	<i>Hidden</i>
2SP	*	*	63.7 ± 9.0	16.4 ± 9.2
Annealing	-	-	-	-
Audiology	-	-	-	-
Bridge	-	-	-	-
Cancer	-	-	-	-
CRX	-	-	-	-
Flag	-	-	-	-
Glass	*	*	65.7 ± 8.7	22.7 ± 9.0
Heart	-	-	-	-
HeartCle	-	-	-	-
HeartHun	-	-	-	-
HeartLB	-	-	-	-
HeartSwi	-	-	-	-
Hepatitis	-	-	-	-
Horse	-	-	-	-
Ionosphere	*	*	92.6 ± 4.3	8.5 ± 2.9
Iris	95.3 ± 6.0	3.1 ± 0.3	92.6 ± 4.3	8.5 ± 2.9
Liver	*	*	72.9 ± 5.1	21.5 ± 27.3
Monks-1	-	-	-	-
Monks-2	-	-	-	-
Monks-3	-	-	-	-
Pima	*	*	74.7 ± 3.9	39.5 ± 27.9
Promoters	-	-	-	-
Sonar	77.5 ± 11.2	4.5 ± 2.3	78.5 ± 6.3	5.3 ± 3.4
Soylarge	-	-	-	-
Soysmall	-	-	-	-
Vehicle	*	*	56.9 ± 5.2	76.1 ± 23.7
Votes	-	-	-	-
Vowel	*	*	50.2 ± 6.2	48.7 ± 28.1
Wine	95.3 ± 5.8	4.2 ± 0.8	97.1 ± 4.0	5.5 ± 1.7
Zoo	-	-	-	-

Table 4: Comparison of generalization accuracy between various algorithms. **DistAl** is the best results from Table 3 and **NN** is the best results obtained by nearest neighbor algorithms in [57].

Dataset	DistAl	NN
Annealing	96.6	96.1
Audiology	66.0	77.5
Bridge	63.0	60.6
Cancer	97.8	95.6
CRX	87.7	81.5
Flag	65.8	58.8
Glass	70.5	72.4
Heart	86.7	83.0
HeartCle	85.3	80.2
HeartHun	85.9	81.3
HeartLB	80.0	71.5
HeartSwi	94.2	93.5
Hepatitis	84.7	82.6
Horse	86.0	76.8
Ionosphere	94.3	92.6
Iris	97.3	96.0
Liver	72.9	63.5
Monks-1	90.9	77.1
Monks-2	100	97.5
Monks-3	99.1	100
Pima	76.3	71.9
Promoters	88.0	92.4
Sonar	83.0	87.0
Soylarge	81.0	92.2
Soysmall	97.5	100
Vehicle	65.4	70.9
Votes	96.1	95.2
Vowel	69.8	99.2
Wine	97.1	97.8
Zoo	96.0	98.9

Table 5: Comparison of neural network pattern classifiers constructed using the entire set of features against those constructed using the best (in accuracy) GA-selected subset for datasets from UCI Repository. *Features* is the number of features used, *Accuracy* is the generalization accuracy obtained, and *Hidden* is the number of hidden neurons generated in the neural networks.

Dataset	All Attributes			GA-selected Subset		
	<i>Features</i>	<i>Accuracy</i>	<i>Hidden</i>	<i>Features</i>	<i>Accuracy</i>	<i>Hidden</i>
Annealing	38	96.6±2.0	12.1 ± 2.4	21.0 ± 3.1	99.5 ± 0.9	11.1 ± 2.9
Audiology	69	66.0±9.7	24.7 ± 4.8	36.4 ± 3.5	83.5 ± 8.2	27.4 ± 5.6
Bridges	11	63.0 ± 7.8	5.2 ± 3.3	5.6 ± 1.5	81.6 ± 7.6	17.6 ± 12.4
Cancer	9	97.8 ± 1.2	2.9 ± 1.2	5.4 ± 1.4	99.3 ± 0.9	5.7 ± 2.9
CRX	15	87.7 ± 3.3	7.7 ± 6.9	8.0 ± 2.1	91.5 ± 2.8	12.5 ± 7.6
Flag	28	65.8 ± 9.5	9.1 ± 6.2	14.0 ± 2.6	78.1 ± 7.8	11.2 ± 6.5
Glass	9	70.5 ± 8.5	9.8 ± 6.9	5.5 ± 1.4	80.8 ± 5.0	14.5 ± 6.6
Heart	13	86.7 ± 7.6	5.7 ± 4.4	7.2 ± 1.6	93.9 ± 3.8	7.5 ± 3.9
HeartCle	13	85.3 ± 2.7	3.4 ± 1.1	7.3 ± 1.7	92.9 ± 3.6	7.6 ± 4.2
HeartHun	13	85.9 ± 6.3	5.0 ± 2.9	7.0 ± 1.2	93.0 ± 4.0	7.1 ± 3.7
HeartSwi	13	94.2 ± 3.8	2.2 ± 0.6	6.6 ± 1.7	98.3 ± 3.3	3.7 ± 1.5
HeartVa	13	80.0 ± 7.4	5.1 ± 2.6	7.1 ± 1.7	91.0 ± 5.7	8.5 ± 3.0
Hepatitis	19	84.7 ± 9.5	6.2 ± 4.0	9.2 ± 2.3	97.1 ± 4.3	8.1 ± 2.8
Horse	22	86.0 ± 3.6	5.3 ± 4.5	11.1 ± 2.3	92.6 ± 3.4	9.5 ± 4.1
Ionosphere	34	94.3 ± 5.0	5.5 ± 1.6	17.3 ± 3.5	98.6 ± 2.4	7.5 ± 2.4
Liver	6	72.9 ± 5.1	21.5 ± 27.3	4.1 ± 0.7	77.8 ± 4.0	25.9 ± 24.3
Pima	8	76.3 ± 5.1	8.1 ± 4.9	3.8 ± 1.5	79.5 ± 3.1	20.8 ± 21.2
Promoters	57	88.0 ± 7.5	2.2 ± 0.4	28.8 ± 3.3	100 ± 0.0	2.7 ± 1.0
Sonar	60	83.0 ± 7.8	6.4 ± 2.7	30.7 ± 3.7	97.2 ± 2.9	7.2 ± 3.0
Soylarge	35	81.0 ± 5.6	20.2 ± 3.2	19.4 ± 2.7	92.8 ± 5.9	23.3 ± 4.3
Vehicle	18	65.4 ± 3.5	23.7 ± 5.0	9.1 ± 1.7	68.8 ± 4.3	36.2 ± 18.2
Votes	16	96.1 ± 1.5	3.2 ± 1.5	8.9 ± 1.8	98.8 ± 1.2	4.0 ± 1.8
Vowel	10	69.8 ± 6.4	38.0 ± 8.3	6.5 ± 1.2	78.4 ± 3.8	41.5 ± 7.7
Wine	13	97.1 ± 4.0	5.5 ± 1.7	6.7 ± 1.6	99.4 ± 2.1	5.9 ± 2.1
Zoo	16	96.0 ± 4.9	6.1 ± 1.1	9.3 ± 1.6	100 ± 0.0	6.2 ± 1.1

Table 6: Comparison of neural network pattern classifiers constructed using the entire set of features against those constructed using the best GA-selected subset in document classification.

Dataset	All Attributes			GA-selected Subset		
	<i>Features</i>	<i>Accuracy</i>	<i>Hidden</i>	<i>Features</i>	<i>Accuracy</i>	<i>Hidden</i>
Abstract1	790	89.0±9.4	3.7 ± 3.5	393.7 ± 12.9	97.6 ± 4.7	5.0 ± 1.9
Abstract2	790	84.0±12.0	9.5 ± 7.0	393.8 ± 14.6	94.4 ± 7.3	11.6 ± 8.2
Reuters1	1568	91.6±2.9	48.8 ± 14.4	786.1 ± 19.1	94.9 ± 2.5	65.4 ± 13.3
Reuters2	435	88.5±10.5	6.2 ± 2.0	218.3 ± 9.7	97.5 ± 4.7	10.6 ± 5.0
Reuters3	1440	96.4±1.6	19.1 ± 3.7	715.4 ± 20.3	98.7 ± 1.0	39.7 ± 10.9