

Symbolic Interpretation of Artificial Neural Networks

Ismail A. Taha and Joydeep Ghosh, *Member, IEEE*

Abstract—Hybrid Intelligent Systems that combine knowledge-based and artificial neural network systems typically have four phases involving domain knowledge representation, mapping of this knowledge into an initial connectionist architecture, network training, and rule extraction, respectively. The final phase is important because it can provide a trained connectionist architecture with explanation power and validate its output decisions. Moreover, it can be used to refine and maintain the initial knowledge acquired from domain experts. In this paper, we present three rule-extraction techniques. The first technique extracts a set of binary rules from any type of neural network. The other two techniques are specific to feedforward networks, with a single hidden layer of sigmoidal units. Technique 2 extracts partial rules that represent the most important embedded knowledge with an adjustable level of detail, while the third technique provides a more comprehensive and universal approach. A rule-evaluation technique, which orders extracted rules based on three performance measures, is then proposed. The three techniques are applied to the iris and breast cancer data sets. The extracted rules are evaluated qualitatively and quantitatively, and are compared with those obtained by other approaches.

Index Terms—Rule extraction, hybrid systems, knowledge refinement, neural networks, rule evaluation.

1 INTRODUCTION

SEVERAL researchers have investigated the design of hybrid systems that combine expert and connectionist subsystems [40], [47], [7], [12], [11], [20], [48]. The typical result of a transformational type of hybridization [49] is a Knowledge-Based Neural Network (KBNN) system with theory refinement capabilities, usually involving four phases:

- 1) the rule-base representation phase, where initial domain knowledge is extracted and represented in a symbolic format (e.g., a rule-based system);
- 2) the mapping phase, where initial domain knowledge is mapped into an initial connectionist architecture;
- 3) the learning phase, where this connectionist architecture is trained by a set of domain examples, and
- 4) the rule-extraction phase, where the trained and, thus, modified connectionist architecture is mapped back into an updated rule-based system.

KBNNs attempt to exploit the complementary properties of knowledge-based and neural network paradigms to obtain more powerful and robust systems. HIA [40], KBANN [46], [26], RAPTURE [20], and KBCNN [7], [8] are examples of KBNN hybrid systems. Fig. 1 sketches typical components of a KBNN system that combines rule-based and connectionist paradigms.

Connectionist systems can alternatively be combined with fuzzy logic systems to obtain “NeuroFuzzy” hybrid

systems. In this setting the neural network subsystem is typically used to adapt membership functions of fuzzy variables [5], or to refine and extract fuzzy rules [42], [41], [18]. Neural networks have also been used for refinement of initial theories expressed in other knowledge representation schemes such as Bayesian Belief networks [29].

Extraction of symbolic rules from trained Artificial Neural Networks (ANNs) is an important feature of comprehensive hybrid systems, as it helps to:

- 1) Alleviate the knowledge acquisition problem and refine initial domain knowledge.
- 2) Provide reasoning and explanation capabilities.
- 3) Support cross-referencing and verification capabilities.
- 4) Alleviate the “catastrophic interference” problem of certain ANNs [37]. For models such as MLPs,¹ it has been observed that if a network originally trained on one task (data set) is subsequently trained on a different task (statistically different data set), then its performance on the first task degrades rapidly. In situations with multiple operating regimes, one can extract rules before the task or environment changes and, thus, obtain different rule sets for different environmental conditions. Together with a mechanism for detecting the current environment, this presents one solution to the “context discovery” and “context drift” problems.

Other uses of rule extraction include improving acceptability of the product, transfer of knowledge to more suitable form, and induction of scientific theories.

This paper proposes three rule-extraction techniques for KBNN hybrid systems. Also, it presents a simple

• I. Taha is with the Department of Computers and Operational Research, Military Technical College, Cairo, Egypt.
E-mail: ismail.taha@mailexcite.com.

• J. Ghosh is with the Department of Electrical and Computer Engineering, University of Texas, Austin, TX 78712-1084.
E-mail: ghosh@ece.utexas.edu.

Manuscript received 21 Sept. 1996; revised 13 Jan. 1998.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number 104323.

1. MLPs stands for Multi-Layered Perceptrons. This abbreviation is being used in this paper to signify a feedforward network, with a single hidden layer of sigmoidal units, unless specified otherwise.

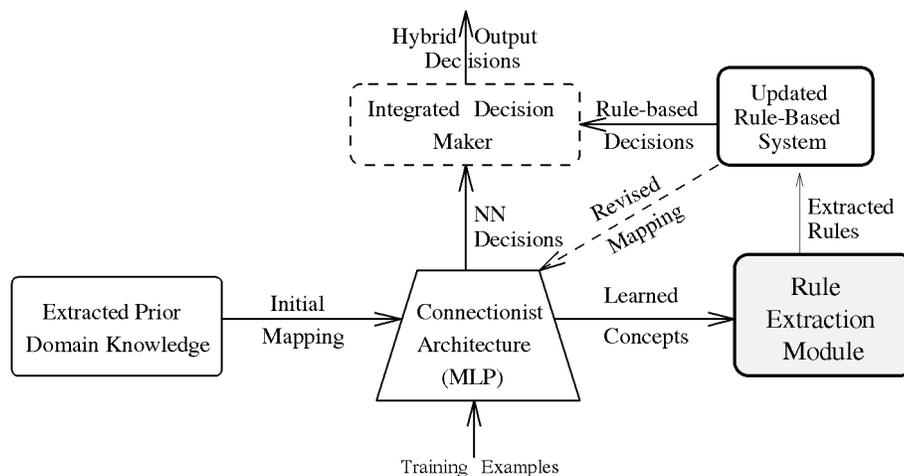


Fig. 1. Typical components of a KBNN system that integrates knowledge-based and connectionist paradigms.

rule-evaluation procedure that orders rules obtained by any extraction approach, according to some performance criteria. A qualitative evaluation of the three new techniques and a comparison with some other approaches is also provided. The next section illustrates key issues in extracting rules from trained neural networks and summarizes some of the existing rule-extraction techniques. Section 3 describes the proposed techniques and the rule-evaluation procedure. In Section 4, we present implementation results of these three techniques using an artificial problem, as well as the iris and breast cancer data sets. Section 5 compares the performance of rule sets extracted by our techniques with the rule sets extracted by some other approaches. In the concluding section, we comment on the different rule extraction techniques, summarize the significance of the proposed techniques, and point to future directions.

2 RULE EXTRACTION

2.1 Issues

Several issues should be carefully considered while designing a rule-extraction technique:

- 1) **Granularity of the explanation feature:** The level of detailed hypotheses and evidence that the system can provide with each of its output decisions.
- 2) **Comprehensiveness of the extracted rules:** In terms of the amount of embedded knowledge captured by them. This directly determines the *fidelity* of the extracted rules in faithfully representing the embedded knowledge.
- 3) **Comprehensibility:** Indicated by the number of rules and number of premises in each extracted rule from a trained network.
- 4) **Transparency of the extracted rules:** In terms of how well the decisions or conclusions can be explained.
- 5) **Generalization capability:** On test samples.
- 6) **Portability:** Capability of the rule-extraction algorithm to extract rules from different network architectures.
- 7) **Modifiability:** The ability of extracted rules to be updated when the corresponding trained network architecture is updated or retrained with different datasets.

- 8) **Theory refinement capability:** That can alleviate the knowledge acquisition bottleneck due to the incompleteness, inconsistency, and/or inaccuracy of initially acquired domain knowledge.
- 9) **Stability or robustness:** A measure of how insensitive the method is to corruptions in the training data or initial domain knowledge.
- 10) **Complexity and scalability:** Computational issues that are relevant for large datasets and rule bases.

These issues, in addition to others, should be used to measure the quality and performance of rules extracted from trained neural networks as discussed in Section 5. Note that these issues also depend on the rule representation, insertion and network training methods used. Also, it is difficult to simultaneously optimize all of the above criteria. For example, a very comprehensive technique may extract too many rules, with some of them having many premises, thus, degrading the robustness and comprehensibility of the resulting rule base.

2.2 Representative Rule Extraction Techniques

Research work in the area of extracting symbolic knowledge from trained ANNs has witnessed much activity recently. This subsection summarizes some of the existing approaches with emphasis on extracting rules from feed-forward (specifically, MLP) architectures. An excellent literature review of different rule-extraction approaches is a survey paper written by Andrews et al. [1]. (Also see the rule extraction home page at <http://www.fit.qut.au/~robert/rulex.html>.)

2.2.1 Link Rule Extraction Techniques

The methodology behind most of the techniques for rule extraction from MLPs can be summarized in these two main steps:

- 1) For each hidden or output node in the network, search for different combinations of input links whose weighted sum exceeds the bias of the current node.
- 2) For each of these combination generate a rule whose premises are the input nodes to this combination of links. All premises of a rule are conjuncted.

Either [27], KT [8], and Subset [45] are three notable rule-extraction algorithms in this category. Some of the problems of KT and Subset algorithms are:

- 1) The size of the search algorithm is $O(2^l)$ for a hidden/output node with fan-in = l ; assuming that network inputs are binary.
- 2) The algorithms extract a large set of rules, up to $\beta_p * (1 + \beta_n)$, where β_p and β_n are the number of subsets of positively weighted and negatively weighted links, respectively.
- 3) Some of the generated rules may be repetitive, as permutations of rule antecedents are not automatically taken care of.
- 4) There is no guarantee that all useful knowledge embedded in the trained network will be extracted.

However, the rules extracted from both algorithms are simple to understand. The size of the extracted rules can be limited by specifying the number of premises of the rules. Generally, the rules extracted by both KT and Subset algorithms are tractable especially in small application domains.

Based on the shortcomings of the Subset algorithm, Towell and Shavlik [45] developed another rule-extraction algorithm called MofN, whose name the rule format that is used to represent the extracted rules:

*If ("at least" M of the following N premises are true) then
(the concept designated by the unit is true).*

The rationale behind MofN is to find a group of links that form an equivalence class whose members have similar effect (due to similar weight values) and can be used interchangeably with one another. MofN extracts rules from the KBANN trained network through six main procedures. For problems like "*promoter recognition in DNA nucleotides*" for which it is a natural fit, rules extracted by MofN are significantly superior to rules extracted by other symbolic approaches such as C4.5, Either, and LINUS [45]. NeuroRule is another rule-extraction approach that uses different combinations of weighted links to extract rules [36]. It extracts rules from networks after pruning their architectures and then discretizing their hidden units activation values. Another algorithm, by Howes and Crook [14], first puts limits on hidden nodes activation values to satisfy an activation on output nodes of at least 0.9. After this step the algorithm searches for input combinations that satisfy the predetermined (constrained) hidden nodes activation values. If found, the algorithm extracts a rule for each combination. While all the previously mentioned approaches are design for networks with binary inputs, this algorithm has a proposed, though not efficient, extension to continuous valued inputs.

We categorize all the approaches mentioned in this subsection as *Link Rule Extraction (LRE) techniques* because they all first search for weighted links that cause a node (hidden or output) to be "active." Then these combinations of weighted links are used to generate symbolic rules. Heuristic methods are commonly used in the LRE category to bound the search space for rules and to increase the

comprehensibility of the extracted rules. Some researchers use the term "*decompositional methods*" to refer to LRE type techniques [1], [8].

RuleNet [22] and RULEX [2] are two examples of "localized LRE" techniques, tailored to networks with localized hidden units. RULEX extracts rules from a Constrained Error Back-Propagation (CEBP) MLP network, similar to Radial Basis Function (RBF) networks. Each hidden node in this CEBP network is localized in a disjoint region of the training examples. A distinctive feature of RULEX is that it controls the search space through its network while other approaches use heuristic measures to do the same. RuleNet, on the other hand, uses the idea of adaptive mixture of local experts [15] to train a localized ANN then extracts binary rules in an LRE approach.

2.2.2 Black-Box Rule Extraction Techniques

Another class of rule-extraction approaches extracts rules from feedforward networks only by examining their input-output mapping behavior. An example of such a rule-extraction approach is the algorithm developed by Saito and Nakano to extract medical diagnostic rules from a trained network [31]. BRAINNE [33], Rule-extraction-as-learning [6], and DEDEC [44] are other examples of extracting rules by investigating the input-output mapping of a trained network. In this paper we refer to this class as the *Black-Box Rule Extraction (BRE)* category because rules are extracted regardless the type or the structure of the neural network. Another given name to this class of rule-extraction techniques is "*pedagogical*" approaches [2]. For example, DEDEC extracts rules by ranking the inputs of an ANN according to their importance (contribution) to the ANN outputs [44]. This ranking process is done by examining the weight vectors of the ANN, which puts DEDEC on the border between LRE and BRE techniques. The next step in DEDEC is to cluster these ranked inputs and use each cluster to generate a set of optimal binary rules that describes the functional dependencies between the attributes of this cluster and the outputs of the ANN. DEDEC has been implemented using a standard feedforward MLP and a Cascaded Correlation (CasCor) ANN. In spite of the LRE nature of its ranking procedure, DEDEC is classified as a BRE since its main theme is to extract rules based on the input-output mapping.

2.2.3 Extracting Fuzzy Rules from ANNs

Research in the area of Fuzzy Logic Neural Networks (FLNN) or "NeuroFuzzy systems" is concerned with combining neural networks and fuzzy logic. Some FLNN systems include a fuzzy rule-extraction module for refining fuzzy sets membership functions and explaining the trained neural network [13], [42], [41], [18].

2.2.4 Extracting Rules from Recurrent Networks

Recurrent networks have shown great success in representing finite state languages and deterministic finite state automata [10]. Omlin and Giles [25] have developed a heuristic algorithm to extract grammar rules in the form of Deterministic Finite-state Automata (DFA) from discrete-time

neural networks and specifically from second-order networks. Starting from a defined initial network state that represents the root of the search space, the DFA rule-extraction algorithm searches the equally partitioned output space of N state neurons in a breadth-first fashion. The authors claim that the DFA rules extraction algorithm improves network generalization performance based on the stability of the internal DFA representation.

3 PROPOSED RULE EXTRACTION APPROACHES

In this section, we introduce three different approaches to extracting rule bases from *trained neural networks*. The suitability of each approach depends on the network type, inputs, complexity, nature of application, the required quality of the extracted rules and some other factors as explained later. The first approach is a Black-Box Rule Extraction technique. The second and the third approaches belong to the Link Rule Extraction category. Subsequently, an evaluation procedure and rule-ordering algorithm is proposed.

3.1 First Approach (BIO-RE)

The first approach is a simple black box rule-extraction technique, that is surprisingly effective within its (relatively) narrow domain of applicability. It is named *Binarized Input-Output Rule Extraction* (BIO-RE) because it extracts binary rules from any neural network trained with “binary” inputs, based on its input-output mapping. If original inputs are not binary, they have to be binarized using (1).

$$y_i = \begin{cases} 1 & \text{if } x_i \geq \mu_i \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

where x_i is the value of original input X_i , μ_i is the mean value of X_i , and y_i is the corresponding binarized input.

See the outline of the BIO-RE Algorithm.

BIO-RE Algorithm

For a well trained neural network do:

- 1) Obtain the network output $O(Y) = \{o_j(Y) | o_j \in \{0, 1\}\}$ corresponding to each binary input pattern Y . If the number of input nodes is n then this conceptually requires 2^n input patterns. However, the problem specification may remove some combinations that will not occur.
 - 2) Generate a truth table by concatenating each input Y of step 1 and its corresponding output decision $O(Y)$ from the trained network. An output $o_j(Y)$ is set to 1 if its corresponding output node is active (above a threshold), otherwise it is 0.
 - 3) Generate the corresponding Boolean function (represented in the previously described binary rule format) from the truth table of step 2.
-

Any available Boolean simplification method can be used to perform step 3 of the BIO-RE algorithm (e.g., Karnaugh map, algebraic manipulation, or a tabulation method). We used Espresso² to generate the extracted rules [3]. Rules extracted by BIO-RE are represented in the format:

$If [Not] Input-Variable [And [Not] Input-Variable]^* \rightarrow Consequent_j$

where $[\cdot]$ is an optional term and $[\cdot]^*$ means that the term $[\cdot]$ can be repeated 0 or n times. In terms of final rules, an extracted rule “If Y_1 And Not Y_2 Then O_1 ” is rewritten as “If $X_1 > \mu_1$ And $X_2 \leq \mu_2$ Then O_1 ” where a “true” binary input (e.g., Y_1) is represented as “ $X_1 > \mu_1$ ” and a “negated” binary input variable (e.g., Y_2) is represented as: “ $X_2 \leq \mu_2$.” See Section 4.3 for examples.

The BIO-RE approach is suitable when the input/output variables are naturally binary or when binarization does not significantly degrade the performance. Also the input size (n) should be small. Given that the above conditions are satisfied, BIO-RE has some advantages:

- 1) It allows the use of available logic minimization tools.
- 2) Extracted rules are optimal and cannot be simplified any further. Hence, no rewriting procedure is required.
- 3) The set of rules extracted by BIO-RE is comprehensive. They are quite understandable since all premises of the extracted rules are conjuncted.
- 4) It does not require any information about the internal structure of the network, the type of network or the training algorithm used.

Section 4 presents experimental results obtained when the BIO-RE algorithm was tested on three problems.

3.2 Second Approach (Partial-RE)

While BIO-RE is efficient, it may not be suitable if the number of inputs is large and/or the forced binarizing of these inputs severely degrades network performance. These deficiencies motivate our second rule-extraction approach—a link rule-extraction method called *Partial-RE*. The idea underlying Partial-RE algorithm is to identify dominant link combinations. For each hidden or output node, j , the positive and negative incoming links are first sorted in descending order of weight values into two sets. Then, starting from the highest positive weight (say, i), the algorithm searches for individual incoming links that can cause the node j to be active regardless of other input links to this node. If such a link exists, it generates a rule:

$$" If Node_i \xrightarrow{cf} Node_j "$$

where cf represents the measure of belief in the extracted rule and is equal to the activation value of $node_j$ with this current combination of inputs. Values of certainty factors are computed by (3). If an afferent node i was found strong enough to activate node j , then this node is marked and cannot be used in any further combinations when inspecting node j . Partial-RE continues checking subsequent

2. Espresso is a software package for logic design [30].

weights in the positive set until it finds one that cannot activate the current node j by itself.

It is important to mention that Partial-RE assumes that all inputs have the same range, so that their effect on the hidden layer is simply determined by the weights. Therefore, the original input features may need to be scaled using (2).

$$z_i = \frac{1.0}{1.0 + e^{-\left(\frac{x_i - \mu_i}{2.0\sigma_i}\right)}} \quad (2)$$

where, $z_i \in (0, 1)$ is the corresponding scaled input value of the original input value x_i and σ_i is the standard deviation of input feature X_i . In (2), σ_i is multiplied by "2" to provide a wider distribution of input X_i (a range of $\mu_i \pm 2\sigma_i$ will contain approximately 95 percent of the X_i measurements [23], if X_i is normally distributed).

If more detailed rules are required (i.e., the comprehensibility measure $p > 1$), then Partial-RE starts looking for combinations of two unmarked links starting from the first (maximum) element of the positive set. This process continues until Partial-RE reaches its terminating criteria (maximum number of premises in rule = p). Also, it looks for negative weights such their not being active allows a node in the next layer to be active, and extracts rules in the format:

$$\text{If Not Node}_g \xrightarrow{cf} \text{Node}_j.$$

Moreover, it looks for small combinations of positive and negative links that can cause any hidden/output node to be active, to extract rules such as:

$$\text{If Node}_i \text{ And Not Node}_g \xrightarrow{cf} \text{Node}_j,$$

where the link between node i and j is positive and between g and j is negative. After extracting all rules, a rewriting procedure takes place. Within this rewriting procedure any premise that represents an intermediate concept (i.e., a hidden unit) is replaced by the corresponding set of conjuncted input features that causes it to be active. Final rules are written in the format:

$$\text{" If } X_i \geq \mu_i \text{ And } X_g \leq \mu_g \xrightarrow{cf} \text{Consequent}_j \text{"}$$

See Table 2 and Table 6 for examples.

Partial-RE can be used efficiently in applications where the main objective of extracting rules from trained neural networks is to study the primary parameters that cause specific output decisions to be taken. In such cases, the cost of implementing the Partial-RE is low compared to the MofN algorithm if a small number of premises per rule is enough. By extracting only certain rules with small number of premises per rule we are reducing the combinatorial nature of the rule-extraction process into one that is polynomial in n , where n is the number of input features. Partial-RE examines small subsets S_{j_s} of incoming links to a hidden or output node j , and extracts a rule if

$$\sum_{w_{ji} \in S_{j_s}} (w_{ji}x_i - \theta_j) \geq \Delta,$$

where w_{ji} is the weight value of the link between input x_i and hidden/output node j , θ_j is the threshold value of the node j , and Δ is a small positive value (between 0.1 and 0.3) called *certainty parameter*, whose value is chosen based on how "certain" we want the extracted rules to be. Partial-RE is conceptually close to subset-type algorithms [8], [45], not surprising since all are LRE algorithms. Note however that while subset algorithms reduce the size of the search space (all subsets of a given size) through branch-and-bound, Partial-RE achieves this through sorting by weight sizes, which is less computationally demanding. Partial-RE also avoids the rewriting procedure involved in subset algorithms, and is able to produce soft rules with associated measures of belief or certainty factors. The adjustable parameters Δ and p (which determines the number of premises in a rule) provides flexibility to the Partial-RE algorithm. Several other differences in implementational details are described in [38].

Partial-RE is easily parallelizable, as nodes can be inspected concurrently. Experimental results show that Partial-RE algorithm is suitable for large size problems, since extracting all possible rules is NP-hard and extracting only the most effective rules is a practical alternative.

3.3 Third Approach (Full-RE)

Like the Partial-RE approach, Full-RE falls in the LRE category. It is notable because:

- It extracts rules with certainty factors from trained feedforward networks whose cells have monotonically increasing activation functions.
- It can extract rules from networks trained with continuous, normal, and binary inputs. Therefore, there is no restriction on the values that any input feature can take. This capability makes Full-RE a "universal" extractor.

For each hidden node j , Full-RE first generates *intermediate* rules in the format:

$$\text{If} \left[(w_{j1} \cdot X_1 + w_{j2} \cdot X_2 + \dots + w_{jn} \cdot X_n) \geq \theta_j + \Delta \right] \xrightarrow{cf} \text{Consequent}_j$$

where θ_j and Δ are the same as defined for Partial-RE.

Note that a range of X_i values may satisfy an intermediate rule, and one would want to determine a suitable extremum value in such a range. To make this tractable, each input range is discretized into a small number of values that can be subsequently examined. Thus, each input feature $X_i \in (a_i, b_i)$ is discretized using k intervals:

$$(D_i \in \{d_{i,0} = a_i, d_{i,1}, \dots, d_{i,k-1}, d_{i,k} = b_i\}),$$

where $d_{i,l-1}$ and $d_{i,l}$ are the lower and upper boundary values of interval l of input X_i , respectively. Different discretization approaches can be exploited to compute discretization boundaries of input features X_i s [39], [16]. Full-RE uses the Chi2 [19] algorithm,³ a powerful discretization tool, to

3. We are thankful to Liu and Setiono for making their Chi2 source code available to us.

compute discretization boundaries of input features. When Full-RE finds more than one discretization value of an input X_i that satisfies the intermediate rule (i.e., the rule has more than one feasible solution), it chooses the minimum or the maximum of these values based on the *sign* of the corresponding link. If w_{ji} is negative then Full-RE chooses the maximum discretization value of X_i , otherwise it chooses the minimum value. However, all selected discretization values should satisfy the left-hand side (the inequality) of the intermediate rule and the boundary constraints of all input features of this inequality.

After discretization, the following Linear Programming (LP) problem is formulated:

$$\text{Minimize } w_{j1}D_1 + w_{j2}D_2 + \dots + w_{jn}D_n$$

such that

$$w_{j1}D_1 + w_{j2}D_2 + \dots + w_{jn}D_n \geq \theta_j + \Delta,$$

and

$$D_i \in \{a_i, d_{i,1}, \dots, d_{i,k-1}\} \forall i, i = 1 \dots n.$$

Full-RE solves this LP problem by selecting the discretization boundaries $X_S(D_S)$ that determine the feasible solution(s) of the intermediate rule and satisfy all given constraints. Values of single input features that can satisfy this LP problem regardless of other input features can be found easily by substituting X_S of positive weights to nodes j by their minimum values (a_S) and X_S of negative weights by their maximum values (b_S). Higher combinations are found similarly by examining different combinations of discretized inputs and finding the edges of the feasible solution surface of the LP problem. Note that any linear programming tools can alternatively be used to solve this standard LP problem.⁴ As an example, assume that a feasible solution is $x_1 = e_1$ and $x_2 = e_2$, and that the weights from input X_1 and X_2 on node j are positive and negative, respectively. Then, Full-RE extracts the following rule:

$$\text{" If } X_1 \geq e_1 \text{ And } X_2 \leq e_2 \xrightarrow{cf} h_j \text{"}$$

where $a_i \leq e_i \leq b_i$, and e_i s are determined by the discretization process or the LP tool.

Full-RE computes certainty factors of extracted rules based on (3).

$$cf = \begin{cases} \frac{1}{1 + \exp(\sum_{i=1}^n w_{ji} \cdot x_i - \theta_j - \Delta)} & \text{if act(j) is a sigmoid} \\ \min(1, \sum_{i=1}^n w_{ji} \cdot x_i - \theta_j - \Delta) & \text{if act(j) is a linear threshold} \\ 1 & \text{if act(j) is a hard limiting and} \\ & \sum_{i=1}^n w_{ji} \cdot x_i \geq \theta_j + \Delta \end{cases} \quad (3)$$

Extracted rules between hidden and output nodes are represented in the same format of Partial-RE, e.g.,

$$\text{If } h_1 \text{ And } h_2 \xrightarrow{cf} O_k.$$

Full-RE replaces each hidden node (h_j) in the previous rule by the left-hand side of the rule(s) whose right-hand side is h_j . The general format of final rules extracted by the Full-RE is:

If Simple – Boolean – Expression

$$[\text{AND Simple – Boolean – Expression}]^* \xrightarrow{cf}$$

Consequent_j

where

Simple – Boolean – Expression ::=

Variable Operator Constant,

Operator ::= > | < | ≥ | ≤

The * means that the term [AND Simple-Boolean-Expression] may be repeated (up to $n - 1$ times), the symbol | stands for an alternation (i.e., *operator* can take any of the four Boolean operators) and *cf* represents the certainty factor computed by (3) for each extracted rule. The final rules extracted by Full-RE have the same format as Partial-RE except that each μ_i is replaced by one of the discretization boundaries (say, $d_{i,j}$) selected by Full-RE as described earlier. See Table 3 and Table 7 for examples. Note that there is no restriction on the number of premises in the final rules extracted by the Full-RE. The only limitation applied is on the number of premises of rules between nodes of adjacent layers (e.g., number of premises in intermediate rules between input and hidden nodes or between hidden and output nodes). Note that when input features are binary, the discretization step is no longer required.

3.4 Rule Evaluation

The suitability of rule evaluation criteria depends on the application domain. For example, in image understanding, one may be primarily interested in finding relational visual structures [4]. Also, one can define metrics for the complete rule set [17] or for individual rules. We will do the former in Section 5. In this section, we evaluate individual rules based on the following primary goals.

- 1) Find the order of the extracted rules that maximizes their performance on the available data set, assuming an inference engine that inspects rules sequentially.
- 2) Test the fidelity of the extracted rule-based system (i.e., its capability to mimic the embedded knowledge in the trained network).
- 3) Measure how much knowledge is left unextracted from the internal structure of the trained network.
- 4) Identify cases where the extracted rule-based system surpasses the trained neural network and vice versa. This analysis helps in the process of integrating and combining the output decisions of the two subsystems.

The simplest and most popular inference engines examine rules in a predetermined sequential order, and make a decision based on the first fireable rule encountered. Alternatively, an inference engine can identify all fireable rules. The latter inference engine is considered more powerful

4. We also used Mathematica to find out feasible solutions.

then the former because it provides the system user with all possible output decisions and hence more choices can be examined.

For both types of inference engines, a predetermined order of the extracted rules plays an important role in determining which rule is going to be fired or the order in which all fireable rules are considered. Since the network does not directly provide information on rule ordering, a rule-evaluation procedure that can rank order the extracted rules is crucial. We achieve this based on three performance measures:

- 1) **The soundness measure:** This measures how many times each rule is correctly fired, as one goes through the available data set. A rule is correctly fired if all its *premises* are satisfied and its *consequent* matches the target decision. The *soundness measure* of an extracted rule represents the ability of this rule to correctly interpret the output decisions of the trained network. Note that the soundness measure does not depend on the rule order.
- 2) **The completeness measure:** This measures attached to a rule represents how many *distinct* times this rule is correctly fired (i.e., how many unique patterns are correctly identified/classified by this rule and not by any other extracted rule that is inspected by the inference engine before this rule). Thus this measure depends on the order in which the extracted rules are applied. For example, a rule with soundness measure > 0 can have zero completeness measure if all the input patterns covered by this rule are also covered by at least one of the preceding rules.
- 3) **The false-alarm measure:** This measures how many times a rule is misfired over the available data set. When considered for application, a rule is *misfired* if all its premises are satisfied but its consequent does not match the target output. The value of this measure also depends on the order of rule application and the mechanism of the inference engine.

3.4.1 Rule Ordering Algorithm

Finding the optimal ordering of extracted rules is a combinatorial problem. So we have developed the following “greedy” algorithm to order any set of extracted rules, based on the three performance measures. The rule-ordering algorithm first creates a list L that contains all extracted rules. This list is divided into two lists, a head list (L_h) and a tail list (L_t), where L_h is the list of all ordered rules and L_t is the list of all remaining (unordered) rules.⁵ Initially, L_h is empty and L_t includes all the extracted rules. A performance criteria is used to select one rule from L_t to be moved to the end of L_h , and the process continues until L_t is null.

See the steps of the Rule-Ordering Algorithm.

Rule-Ordering Algorithm

- 1) Initialize $L_h = \{ \}$, $L_t = \{ \text{all extracted rules} \}$.
 - 2) **WHILE** $L_t \neq \{ \}$, **DO**
 - a) Fire all rules in L_h in order.
 - b) Compute the completeness and false-alarm measures for each rule in L_t using the available data set.
 - c) **IF** \exists a rule with zero false-alarm **THEN** This rule is moved from L_t to the end of L_h .⁶ **ELSE** Among all rules in L_t select the one with the highest (Completeness—False-alarm) measure; add this rule to the end of L_h , delete it from L_t .
 - d) **IF** \exists any rule in L_t with a zero completeness measure then remove this rule from L_t .
 - 3) **END DO.**
-

In this paper, all rules extracted by our approaches are ordered using the above algorithm. Also, the measures attached to all extracted rules assume that an inference engine that fires only one rule per input (namely, the first fireable rule) is used. Note that the three performance measures along with a rule’s certainty factor (if applicable) can be used to form a composite measure (in an application dependent manner) which specifies the importance of the extracted rules.

An important issue that needs to be addressed here is: “Should one discard rules with low soundness, completeness, and/or high false-alarm measure(s)?” An example of such a rule is R_{10} in Table 5. For small datasets, we might still retain such rules at the bottom of the application ladder in the hope of better generalization, as they are part of the overall characteristics of the corresponding trained network. In cases where available datasets are representative, the answer depends on the application nature. For example, in medical applications where one is interested in high detection rates, rules with low soundness, completeness, and/or high false-alarm measures may still be kept. In other applications, like Automatic Target Recognition (ATR), one may only retain rules with low false-alarm rates to reduce the chances of “friendly fire.”

4 IMPLEMENTATION AND PERFORMANCE EVALUATION

4.1 Data Sets

We applied all three rule-extraction techniques to three problems:

- 1) An artificial rule-based system which has six rules relating four binary inputs and four binary outputs.
- 2) Iris database, a simple classification problem which contains 50 examples each of classes Iris Setosa, Iris Versicolor, and Iris Virginica [24]. These 150 instances were partitioned into a training set of size 89 and a testing set of size 61. Each input pattern has four

6. If \exists more than one rule with zero false-alarm, **THEN** select the one with the highest completeness measure out of these rules to be moved from L_t to the end of L_h .

5. That is, the ordering of rules in L_t has no effect.

continuous input features: I_1 = Sepal-length, I_2 = Sepal-width, I_3 = Petal-length, and I_4 = Petal-width.

- 3) Breast-Cancer data set which has nine inputs and two output classes [21], [24]. The input features are: X_1 = Clump Thickness, X_2 = Uniformity of Cell Size, X_3 = Uniformity of Cell Shape, X_4 = Marginal Adhesion, X_5 = Single Epithelial Cell Size, X_6 = Bare Nuclei, X_7 = Bland Chromatin, X_8 = Normal Nucleoli, and X_9 = Mitoses. All nine inputs are continuous and in the range $[1, 10]$. Each of the 683 available instances is labeled as Benign (444 instances) or Malignant. These instances are divided into a training set of size 341 and a test set of size 342.

Other popular data sets that have been used as benchmarks for rule-extraction approaches are the Monk [43], Mushroom [32], and the DNA promoter [47] datasets. All three of these data sets inputs are symbolic/discrete by nature. Since we want to test more general problems that may include continuous valued variables, Iris and Breast-Cancer were preferred for our initial experiments.

4.2 Methodology

We now describe some of the important procedures followed to perform the experimental work presented in this paper.

- 1) **Training procedure:** In all experiments, an MLP network is trained using the backpropagation algorithm with momentum as well as a regularization term P which adds

$$\frac{-2\lambda w_{jk} w_0^2}{w_0^2 + w_{jk}^2}$$

to the weight update term in the backpropagation equation [9]. Cross-validation is used for the stopping criteria.

- 2) **Network architectures and data reduction:** For the iris problem, an MLP with four input, six hidden, and three output nodes is used for the three experiments but trained with different data sets each time, as described later. For the breast-cancer classification problem, we reduced the dimensionality of the input space from nine to six inputs. This has been done by removing X_4 , X_5 , and X_6 , as these inputs correspond to the lowest three eigenvalues of the covariance matrix of the original input space. The remaining six input features are then used for training and testing an MLP with nine hidden and two output nodes.
- 3) **Network initialization:** For the artificial problem, the six initial rules are used by the Node Link Algorithm [40], to initialize a network of four input, six hidden, and four output nodes. For both the iris and breast-cancer data sets, there is no prior knowledge, so the corresponding networks are initialized randomly.
- 4) **Input representation:** Inputs of the artificial problem are naturally binary, so there was no required mapping. Since the input features of both iris and breast-cancer problems are continuous, a binarized and a

normalized version of these two data sets were computed and then used for training/testing the network architectures subsequently used for BIO-RE and Partial-RE, respectively.

- 5) **Extraction techniques and networks labeling:** BIO-RE is used to extract rules from networks trained with binary/binarized input patterns. For iris and breast-cancer problems, these networks are labeled *Iris-Bin* and *Cancer-Bin*, respectively. Partial-RE is used to extract rules from the networks trained with normalized input patterns (labeled *Iris-Norm* and *Cancer-Norm*). Full-RE uses the original datasets of both problems to train the corresponding networks. These two networks are labeled *Iris-Cont* and *Cancer-Cont*.
- 6) **Default class rule:** A comprehensive rule-extraction approach is one that extracts rules to cover all input-output mapping cases. In some cases, achieving such a goal is hard, and it may be convenient to use default rules to cover the input-output mapping cases that cannot be covered by the extracted rule-base. Such default rules make the set of extracted rules complete, but they do not provide any interpretation of why the default action was done other than “None of the extracted rules could be fired.” If a default rule is used, its output (consequent) can be chosen to minimize the false-alarm rate and to maximize the correct classification rate. Note that the default rule may only fire when none of the extracted rules can be fired.

4.3 Experimental Results

4.3.1 An Artificial Binary Problem

This experiment is designed to test the soundness and completeness of the three rule-extraction techniques. The initial rules are as follows:

- Rule 1: $If A \text{ And } B \xrightarrow{0.8} O_1$
- Rule 2: $If B \text{ And } C \text{ And } D \xrightarrow{0.7} O_2$
- Rule 3: $If \text{ Not } C \xrightarrow{0.6} O_3$
- Rule 4: $If \text{ Not } A \text{ And } D \xrightarrow{0.7} O_4$
- Rule 5: $If B \text{ And } D \xrightarrow{0.7} O_1$
- Rule 6: $If D \xrightarrow{0.8} O_1$

where A , B , C , and D are binary inputs, and O_s ($i = 1, \dots, 4$) are binary consequents. After using the Node Link Algorithm [40] to map these six rules into an initial network with four input, six hidden, and four output nodes, the following two experiments were performed.

Experiment 1: The objective of this experiment is to check whether the three approaches are able to extract the original rules from the mapped network. Therefore, the network was not trained before the extraction procedures were applied.

The results of applying the three rule-extraction techniques to the generated (but not trained) network are as follows:

TABLE 1
RULES EXTRACTED FROM NETWORK "IRIS-BIN" BY BIO-RE TECHNIQUE

Rule No.	Rule Body	Iris Class	Soundness Measure	Completeness Measure	False-Alarm Measure
1	If $I_2 \leq 3.0$ $I_3 \geq 3.7$ and $I_4 \leq 1.2$	Versicolor	10/50	10/50	0/150
2	If $I_1 \leq 5.8$ and $I_3 \leq 3.7$ and $I_4 \leq 1.2$	Setosa	50/50	50/50	5/150
3	If $I_1 \geq 5.8$ and $I_3 \geq 3.7$ and $I_4 \geq 1.2$	Virginica	47/50	47/50	24/150
4	If $I_1 \leq 5.8$ and $I_2 \leq 3.0$ and $I_4 \geq 1.2$	Versicolor	15/50	11/50	3/150
Overall Performance %				118/150	32/150
				78.67%	21.33%

TABLE 2
RULES EXTRACTED FROM NETWORK "IRIS-NORM" BY PARTIAL-RE TECHNIQUE

Rule No.	Rule Body	Iris Class	Certainty Factor	Soundness Measure	Completeness Measure	False-Alarm Measure
1	If $I_2 \geq 3.0$ and $I_3 \leq 3.7$	Setosa	0.98	48/50	48/50	0/150
2	If $I_1 \geq 5.8$ and $I_3 \geq 3.7$ and $I_4 \geq 1.2$	Virginica	0.99	47/50	47/50	27/150
3	If $I_1 \leq 5.8$ and $I_2 \leq 3.0$ and $I_3 \geq 3.7$	Versicolor	0.74	18/50	16/50	3/150
4	If $I_1 \leq 5.8$ and $I_2 \leq 3.0$ and $I_4 \geq 1.2$	Versicolor	0.72	15/50	1/50	0/150
5	Default Class	Versicolor	1.0	6/50	6/50	2/150
Overall Performance %					118/150	32/150
					78.67%	21.33%

- BIO-RE extracts the same set of binary rules but without certainty factors.
- Partial-RE with $p = 2$ (i.e., maximum two conditions per rule) extracts all five original rules with two conditions or less. On increasing p to 3, Rule #2 was also extracted. The certainty factors attached to each output decision were approximately the same as the original rules.
- Full-RE extracts the same six original rules from the untrained network.

Experiment 2: Based on the original rules, 2^4 binary patterns were generated. After training the previous network, we applied the three approaches to the final network architecture (i.e., the adapted one):

- Both BIO-RE and Full-RE extract the same six original rules.
- Partial-RE extracts all the six rules plus an extra one:
Rule #7: If B And $D \rightarrow O_2^{0.74}$. This rule was extracted when $p = 3$.

4.3.2 Iris Classification

Table 1, Table 2, and Table 3, present the ordered rules extracted by BIO-RE, Partial-RE, and Full-RE techniques, respectively, from their corresponding networks trained on the Iris data set. They also present the corresponding measures for each extracted rule as generated by the

rule-evaluation procedure. Table 4 provides a summary of the performance of each rule-extraction technique and compares it with the performance of the corresponding trained network. It shows that binarizing or scaling input patterns of the iris problem degrades the performance of the trained networks ("Iris-Bin" and "Iris-Norm") as well as the corresponding rules extracted from these two networks. Also, it shows the remarkable performance of the rules extracted from network "Iris-Cont" by Full-RE.

Note that:

- 1) The numeric values compared with input features I_j in the rules extracted by both BIO-RE and Partial-RE represent the mean (μ_j) of these input feature (see the rule bodies in Table 1 and Table 2). This coarse thresholding is largely responsible for the (relatively) poor performance of the two networks and subsequently of the extracted rules.
- 2) In Table 3, a numeric value that is compared to an input feature I_i in the rule body represents one of the critical discretization boundaries of that feature which was selected by Full-RE.
- 3) For rules examined later (e.g., Rule 4 in Table 2), completeness may be much less than soundness, because some instances where these rules would fire correctly have already been covered by other preceding rules.
- 4) Full-RE leads to three simple rules that classify the iris data set very well.

TABLE 3
RULES EXTRACTED FROM NETWORK "IRIS-CONT" BY FULL-RE TECHNIQUE

Rule No.	Rule Body	Iris Class	Certainty Factor(cf)	Soundness Measure	Completeness Measure	False-Alarm Measure
1	If $I_3 \leq 2.1$	Setosa	0.99	50/50	50/50	0/150
2	If $I_3 \leq 5.1$ and $I_4 \leq 1.7$	Versicolor	0.97	49/50	49/50	3/150
3	If $I_3 \geq 4.8$	Virginica	0.98	47/50	47/50	1/150
Overall Performance					146/150	4/150
%					97.33%	2.67%

TABLE 4
PERFORMANCE COMPARISON BETWEEN THE SETS OF EXTRACTED RULES AND THEIR CORRESPONDING TRAINED NETWORKS FOR THE IRIS PROBLEM

		Neural Network		Extracted Rules	
		ratio	% match	ratio	% match
Binarized Network (Iris-Bin)	Training	66/89	74.16	67/89	75.28
	Testing	43/61	70.49	51/61	83.61
	Overall	109/150	72.67	118/150	78.67
Normalized Network (Iris-Norm)	Training	84/89	94.38	69/89	77.53
	Testing	56/61	91.80	49/61	80.33
	Overall	140/150	93.33	118/150	78.67
Continuous Network (Iris-Cont)	Training	87/89	97.75	83/86	96.63
	Testing	59/61	96.72	60/61	98.36
	Overall	146/150	97.33	146/150	97.33

4.3.3 Breast-Cancer Classification

For the breast-cancer classification problem, Table 5, Table 6, and Table 7 present three sets of ordered rules extracted by the three rule-extraction techniques, along with the corresponding performance measures. Table 8 provides an overall comparison between the extracted rules and their corresponding trained networks. It shows that the three techniques were successfully used with approximately the same performance regardless of the nature of the training and testing data sets used for each network. Also, it shows that binarizing and scaling breast cancer data set did substantially degrade performance. Since the original input features of the breast cancer problem have the same range (1-10) and are not very skewed, by binarizing and/or scaling them we did not change their nature much.

4.4 Discussion

The implementation results of Section 4.3 indicate that:

- 1) All rules extracted by the three techniques are sound.
- 2) Partial-RE is sound but not complete. Its completeness depends on the chosen degree of comprehensibility (p).
- 3) Full-RE is sound and complete and can extract exact rules. Full-RE extracts accurate rules to represent the generalization capability of the Rules extracted by Full-RE are much more comprehensible than those extracted by the BIO-RE and Partial-RE. This is likely due to the finer gradation of the input space allowed by this algorithm.
- 4) Binarizing or normalizing continuous features may degrade the accuracy of the extracted rules as well as the

generalization capability of the corresponding trained neural network. See the first six rows of Table 4.

- 5) Full-RE was tested several times on different networks initialized randomly each time and trained with different sets of training patterns. Each time the set of extracted rules are similar except for the values of the certainty factors. This indicates that Full-RE is more accurate and can extract rules based on more combinations of input features, not just the most effective features, see Table 3 and Table 7.
- 6) Although BIO-RE and Partial-RE were used to extract rules from networks trained with binarized and normalized input features, they were still able to extract "certain" rules that may be adequate in some application examples. See Table 5 and Table 6.
- 7) Although some of the extracted rules have a low firing rate on the available data set, they were extracted to represent the generalization capability of the trained network on unseen data. Also, they were extracted to cover all training and testing data sets and hence increase the completeness of the extracted set of rules. Examples of such rules are: R_1 and R_4 of Table 1, R_4 of Table 2, and R_2 - R_5 of Table 5.

5 COMPARATIVE PERFORMANCE EVALUATION

Since both iris and breast-cancer problems have continuous input features, Full-RE is naturally suited for them. There is no need to prune the trained network since Full-RE is capable of extracting rules from MLPs of any size. In this section, we compare the performance of the extracted rules from the iris and the breast-cancer databases with the rules extracted by both NeuroRule [35] and

TABLE 5
RULES EXTRACTED FROM NETWORK "CANCER-BIN" BY BIO-RE TECHNIQUE

Rule No.	Rule Body	B-Cancer Class	Soundness Measure	Completeness Measure	False-Alarm Measure
1	If $X_3 \leq 3.0$ and $X_7 \leq 3.3$ and $X_8 \leq 2.7$ and $X_9 \leq 1.5$	Benign	391/444	391/444	2/683
2	If $X_1 \leq 4.1$ and $X_3 \leq 3.0$ and $X_7 \leq 3.3$ and $X_9 \leq 1.5$	Benign	317/444	8/444	0/683
3	If $X_1 \leq 4.1$ and $X_3 \leq 3.0$ and $X_8 \leq 2.7$ and $X_9 \leq 1.5$	Benign	316/444	7/444	0/683
4	If $X_1 \leq 4.1$ and $X_3 \leq 3.0$ and $X_7 \leq 3.3$ and $X_8 \leq 2.7$	Benign	316/444	7/444	0/683
5	If $X_1 \leq 4.1$ and $X_7 \leq 3.3$ and $X_8 \leq 2.7$ and $X_9 \leq 1.5$	Benign	314/444	5/444	0/683
6	If $X_1 \geq 4.1$ and $X_3 \geq 3.0$	Malignant	200/239	199/239	15/683
7	If $X_3 \geq 3.0$ and $X_7 \geq 3.3$	Malignant	187/239	27/239	2/683
8	If $X_3 \geq 3.0$ and $X_8 \geq 2.7$	Malignant	187/239	3/239	0/683
9	If $X_1 \geq 4.1$ and $X_7 \geq 3.3$	Malignant	167/239	7/239	1/683
10	If $X_1 \geq 4.1$ and $X_9 \geq 1.5$	Malignant	100	1	3
11	Default Class	Benign	5/444	5/444	0/239
Total For Benign Rules %				423/444 95.27%	2/683 0.29%
Total For Malignant Rules %				237/239 99.16%	21/683 3.07%
Overall Performance %				660/683 96.63%	23/683 3.37%

TABLE 6
RULES EXTRACTED FROM NETWORK "CANCER-NORM" BY PARTIAL-RE TECHNIQUE

Rule No.	Rule Body	B-Cancer Class	Certainty Factor(cf)	Soundness Measure	Completeness Measure	False-Alarm Measure
1	If $X_2 \leq 3.0$ and $X_3 \leq 3.0$ and $X_7 \leq 3.3$	Benign	0.99	412/444	412/444	6/683
2	If $X_1 \leq 4.1$ and $X_3 \leq 3.0$ and $X_7 \leq 3.3$	Benign	0.99	324/444	1/444	0/683
3	If $X_2 \geq 3.0$ and $X_3 \geq 3.0$	Malignant	0.99	222/239	219/239	15/683
4	If $X_1 \geq 4.1$ and $X_7 \leq 3.3$ and $X_8 \geq 2.7$	Malignant	0.99	137/239	8/239	0/683
5	If $X_1 \leq 4.1$ and $X_2 \leq 3.0$ and $X_7 \leq 3.3$	Benign	0.84	327/444	4/444	0/683
6	If $X_1 \geq 4.1$ and $X_2 \geq 3.0$	Malignant	0.99	198/239	2/239	0/683
7	If $X_1 \leq 4.1$ and $X_2 \leq 3.0$ and $X_3 \leq 3.0$	Benign	0.84	333/444	9/444	1/683
8	If $X_1 \geq 4.1$ and $X_3 \geq 3.0$	Malignant	0.99	200/239	3/239	2/683
9	If $X_2 \leq 3.0$ and $X_3 \leq 3.0$ and $X_8 \leq 2.7$	Benign	0.99	409/444	1/444	0/683
Total For Benign Rules %					427/444 96.17%	7/683 1.02%
Total For Malignant Rules %					232/239 97.07%	17/683 2.49%
Overall Performance %					659/683 96.49%	24/683 3.51%

C4.5 rules algorithms [36], since these methods have previously been applied on the same databases. Moreover, they both extract comprehensive rules with relatively high correct classification rate, as reported in [35]. For iris problem, we also compare the set of rules extracted by Full-RE with the corresponding set of rules extracted by KT algorithm [8].

Before analyzing the extracted rules, we summarize the computational complexity of NeuroRule and C4.5 rules.

- **NeuroRule:** As a starting point of the NeuroRule algorithm, a large number (100) of fully connected MLPs are generated. Before training any of these 100 networks, each input feature is discretized into one of

TABLE 7
RULES EXTRACTED FROM NETWORK "CANCER-CONT" BY FULL-RE TECHNIQUE

Rule No.	Rule Body	B-Cancer Class	Certainty Factor	Soundness Measure	Completeness Measure	False-Alarm Measure
1	If $X_1 < 8$ and $X_3 < 3$	Benign	0.96	394/444	394/444	5/683
2	If $X_2 \geq 2$ and $X_7 \geq 3$	Malignant	0.83	227/239	223/239	18/683
3	If $X_1 < 8$ and $X_7 < 3$	Benign	0.75	300/444	27/444	1/683
4	If $X_1 \geq 8$	Malignant	0.89	123/239	9/239	1/683
5	If $X_1 < 8$ and $X_2 < 2$	Benign	0.79	369/444	4/444	1/683
Total For Benign Rules					425/444	7/683
%					95.72%	1.02%
Total For Malignant Rules					232/239	19/683
%					97.07%	2.78%
Overall Performance					657/683	26/683
%					96.19%	3.81%

TABLE 8
PERFORMANCE COMPARISON BETWEEN THE SETS OF EXTRACTED RULES AND THEIR CORRESPONDING TRAINED NETWORKS FOR THE BREAST-CANCER PROBLEM

		Neural Network		Extracted Rules	
		ratio	% match	ratio	% match
Binarized Network (Cancer-Bin)	Training	333/341	97.65	331/341	97.07
	Testing	317/342	92.69	329/342	96.20
	Overall	650/683	95.17	660/683	96.63
Normalized Network (Cancer-Norm)	Training	329/341	96.48	331/341	97.07
	Testing	325/342	95.03	328/342	95.91
	Overall	654/683	95.75	659/683	96.49
Continuous Network (Cancer-Cont)	Training	334/341	97.95	330/341	96.77
	Testing	331/342	96.78	327/342	95.61
	Overall	665/683	97.36	657/683	96.19

n values, which are then converted into binary ones using a thermometer coding. Although this discretization step helps simplify the last step of the rule-extraction process, it has the major drawback of greatly increasing the number of input nodes and hence the complexity of the required network architecture. For example, the network architecture generated by NeuroRule has 39 input nodes for the iris problem and 91 input nodes for the breast cancer problem, while the corresponding networks used by Full-RE have four and six input nodes, respectively. This increase in input size makes the network more difficult to train, both in terms of training time and performance.

NeuroRule employs a pruning procedure after the training phase to reduce the complexity of the generated network architectures. The pruning process continues until network performance drops to 95 percent of original performance. This process is applied to the 100 MLPs. The rule-extraction procedure starts by choosing the *best one* out of the 100 pruned networks (the one with the highest performance). NeuroRule extracts rules by clustering the remaining hidden nodes activation values and then checking which input combination can make each hidden (and later output) node active.

The power of NeuroRule lies in its pruning and clustering techniques. In its pruning phase, NeuroRule removes input nodes. For example, the best pruned architecture for the iris problem is a network of four input, two hidden, and three output nodes. For breast-cancer problem the best pruned network has six input, one hidden, and two output nodes. The pruning and clustering processes lead to substantial overhead, but since the resulting network architectures from the pruning step are very small, the rule-extraction process becomes easy.

- **C4.5 rules:** C4.5 rules was used by the authors of NeuroRule to extract rules from the iris and breast-cancer databases for comparison reasons. Like ID3, C4.5 rules [28] generate decision tree rules based on the available input samples. Therefore, the complexity is moderate, but the performance of the rules generated by C4.5 rules is highly affected by the noise level in the available data samples [8].

5.1 Comparison Using Iris Data Set

The rules extracted by the Full-RE techniques for the iris problem were given in Table 3. The rules extracted by NeuroRule for the same problem are:

- Rule 1: If $I_3 \leq 1.9$, then Iris Setosa
- Rule 2: If $I_3 \leq 4.9$ and $I_4 \leq 1.6$, then Iris Versicolor
- Rule 3: Default Rule (Iris Virginica)

TABLE 9
CORRECT CLASSIFICATION RATE (PERCENT) OF THE RULE SETS EXTRACTED BY DIFFERENT TECHNIQUES

		Full-RE	NeuroRule	C4.5rules	KT
Iris	with default rule	97.33	98.00	96.00	97.33
	without default rule	97.33	64.67	96.00	97.00
Breast Cancer	with default rule	96.19	97.21	97.21	N/A
	without default rule	96.19	63.10	94.72	N/A

The corresponding rules extracted by C4.5 rules are:

- Rule 1: If $I_3 \leq 1.9$, then Iris Setosa
- Rule 2: If $I_3 \geq 1.9$ and $I_4 \leq 1.6$, then Iris Versicolor
- Rule 3: If $I_4 \geq 1.6$, then Iris Virginica
- Rule 4: Default Rule (Iris Setosa)

The corresponding rules extracted by KT approach are:

- Rule 1: If $I_3 \leq 2.7$, then Iris Setosa
- Rule 2: If $I_3 \leq 5.0$ and $I_3 > 2.7$ and $I_4 \leq 1.6$ and $I_4 > 0.7$, then Iris Versicolor
- Rule 3: If $I_3 > 5.0$, then Iris Virginica
- Rule 4: If $I_4 > 1.6$, then Iris Virginica
- Rule 5: If $I_2 > 3.1$ and $I_3 > 2.7$ and $I_3 \leq 5.0$, then Iris Versicolor

From the methodologies and results, we note that:

- 1) **Completeness:** Both Full-RE and KT extracted complete sets of rules that cover all cases, and so no default rule was required. However, a default rule is essential for both NeuroRule (due to its pruning step) and for C4.5.
- 2) **Comprehensibility:**
 - a) **Number of rules:** Both Full-RE and NeuroRule extract three rules while KT extracts five rules and C4.5 extracts four rules.
 - b) **Number of premises per rule:** Except KT, the maximum number of conditions per rule for all other techniques is two. KT extracted rules with a maximum of four conditions per rule.
- 3) **Performance:** Since iris is a simple classification problem, all techniques performed well. In fact, all of them were able to show that the Setosa class is linearly separable from the other two classes. Moreover, rule extracted by all of them showed that *Petal-length* is the most dominant input feature (see row 1 and 2 in Table 9).
- 4) **Certainty factors:** Rules extracted by the Full-RE provide a certainty factor attached with each extracted rule, unlike the other approaches.

5.2 Comparison Using Breast Cancer Data Set

For the breast cancer database, the rules extracted by the Full-RE from a simple MLP architecture (six input, six hidden, and two output nodes) are presented in Table 7. The rules extracted by NeuroRule from the best among the pruned 100 MLP network architectures (six inputs, one hidden, and two output nodes) are [36], [34]:

- Rule 1: If $X_1 < 7.0$ and $X_2 < 8.0$ and $X_3 < 3.0$ and $X_8 < 9.0$, then Benign

- Rule 2: If $X_1 < 7.0$ and $X_2 < 8.0$ and $X_3 < 3.0$ and $X_6 < 9.0$, then Benign
- Rule 3: If $X_2 < 8.0$ and $X_3 < 3.0$ and $X_6 < 3.0$ and $X_8 < 9.0$, then Benign
- Rule 4: Default Rule (Malignant).

The corresponding rules extracted by C4.5 rules are [36]:

- Rule 1: If $X_1 < 7.0$ and $X_2 < 8.0$ and $X_3 < 3.0$, then Benign
- Rule 2: If $X_1 < 7.0$ and $X_2 < 2.0$, then Benign
- Rule 3: If $X_2 \geq 5.0$, then Malignant
- Rule 4: If $X_6 \geq 9.0$, then Malignant
- Rule 5: If $X_1 \geq 7.0$, then Malignant
- Rule 6: If $X_4 \geq 4.0$, then Malignant
- Rule 7: Default Rule (Benign).

Comparing these three sets of extracted rules, we observed:

- 1) **Completeness:** NeuroRule did not extract any rule for class Malignant. Both NeuroRule and DT (C4.5) have a default rule while rules extracted by Full-RE cover all cases and hence there was no need for a default rule. Default rules are undesirable because they cannot provide a symbolic interpretation of the decision other than “*because none of the above occurred.*”
- 2) **Comprehensibility:**
 - a) **Number of rules:** Number of rules extracted by Full-RE is five and by DT (C4.5) is seven. NeuroRule extracted only four rules, as it used a default rule to cover all cases of class Malignant, which were applied to a highly pruned network with only one hidden node.
 - b) **Number of premises per rule:** For Full-RE, the maximum number of conditions per extracted rule is two. All rules extracted by NeuroRule have four conditions, while those extracted by DT have a maximum of three conditions per rule. Thus the rules extracted by Full-RE are more comprehensible than those extracted by the other two techniques.
- 3) **Performance:** The performance of the rules extracted by all the three techniques are very high and they all achieve very low misclassification rate (see row 3 of Table 9). When default rules are removed, the performance of NeuroRule drops dramatically (see row 4 of Table 9). Authors of NeuroRule reported that by choosing different trained network architectures they extracted different rules. In one case, only two rules were extracted, one of which is a default rule. In another experiment, NeuroRule extracts only three rules (one of which is also a default rule). In both experiments, the achieved completeness measure is

TABLE 10
A QUALITATIVE COMPARISON OF DIFFERENT RULE-EXTRACTION TECHNIQUES

	BIO-RE	Partial-RE	Full-RE	NeuroRule	KT or Subset	MofN
Provides CF	No	Yes	Yes	No	No	No
May need a default rule	Yes	Yes	No	Yes	Yes	Yes
Works for						
1.Binary inputs	Yes	Yes	Yes	Yes	Yes	Yes
2.Normalized inputs	No	Yes	Yes	No	No	No
3.Continuous inputs	No	No	Yes	No	No	No
Complexity	Very Low	Low	Med	Very High	Med	High
Additional overheads	No	No	No	Yes	No	Yes

approximately 95 percent. However, we did not observe any effect on the rules extracted by the Full-RE due to changing the initialization of the “Cancer-Cont” network or when we used different input samples for training and testing. This indicates that Full-RE rules are quite stable so long the network is trained reasonably well.

- 4) **Certainty factors:** Rules extracted by the Full-RE provide a certainty factor attached with each extracted rule while NeuroRule and C4.5 rules do not. Note that KT was not used to extract rules from the breast cancer problem.

Table 9 compares the classification rates obtained using the rules extracted by the four technique (Full-RE, NeuroRule, C4.5 rules, and KT), for the iris and the breast-cancer databases, while Table 10 presents a qualitative comparison between our three techniques and some other notable rule-extraction techniques from trained neural networks (NeuroRule, KT, Subset, MofN). Note that C4.5 rules was not included in the comparative study presented by Table 10 because it extracts rules from decision trees and not from trained networks like the other approaches.

6 CONCLUSIONS

In this paper, we introduced three new rule-extraction techniques. The suitability of each approach depends on the network type and architecture, complexity, the application nature, inputs, and the required transparency level. All three methods are able to extract meaningful rules for the well known Iris database and Wisconsin breast cancer diagnosis database, where no preexisting rules are available. The extracted rules compare favorably with other reported implementation results. The proposed techniques are less complex, and the rules extracted by them are efficient, comprehensible and powerful.

The ordering of extracted rules has to be determined while designing the inference engine. The neural network does not provide any (direct) information on this issue, and other KBNN researchers have not reported on this aspect. We developed a simple greedy rule-evaluation procedure

and an algorithm that can order rules extracted by any rule-extraction algorithm, with a goal of maximizing performance and minimizing error rates of the extracted rules over available data. We also presented a qualitative comparison of some key issues involved in the process of extracting rules from trained networks by different approaches.

It is important to mention that obtaining all possible combinations of rules is NP-hard and a feasible alternative is often to extract key rules that cover most of the embedded knowledge. More progress is needed in determining when an adequate set of rules has been extracted. Another important issue that needs to be investigated is how the outputs of both the rule extraction and the trained ANN modules can be integrated to provide more robust decisions, and how the extracted rules can be used for knowledge refinement and truth maintenance of domain knowledge. The stability or robustness of rule-extraction methods in presence of noise or outliers in the data should also be explored further.

ACKNOWLEDGMENTS

This research was supported, in part, by ARO Contracts No. DAAH04-94-G-0417 and No. DAAH04-95-10494, and by ATP Grant No. 442. Ismail Taha was also supported by the Egyptian government’s PhD fellowship program. A preliminary version of this paper appears in the *Proceedings of ANNIE 1996*, held in St. Louis, Missouri, in November 1996.

Joydeep Ghosh benefited from the superb NIPS 1996 workshop on rule extraction, organized by Robert Andrews and Joachim Diederich, and from the comments of Lee Giles, Paul Munro, Mark Craven, and other workshop participants. Discussions on theory refinement with Ray Mooney and Xiaowei Wang were very helpful, as was the input provided by the anonymous referees.

REFERENCES

- [1] R. Andrews, J. Diederich, and A. Tickle, "A Survey and Critique of Techniques for Extracting Rules from Trained Artificial Neural Networks," *Knowledge-Based Systems*, vol. 8, no. 6, pp. 373-389, Dec. 1995.
- [2] R. Andrews and S. Geva, "Inserting and Extracting Knowledge from Constrained Error Backpropagation Networks," *Proc. Sixth Australian Conf. Neural Networks*, Sydney, Australia, 1995.
- [3] R. Brayton, G. Hachtel, C. McMullen, and A. Sangiovanni-Vincentelli, *Logic Minimization Algorithm for VLSI Synthesis*. Kluwer, 1984.
- [4] T. Caelli and W. Bischof, "The Role of Machine Learning in Building Image Interpretation Systems," *J. Artificial Intelligence and Pattern Recognition*, vol. 11, p. 143-168, 1996.
- [5] R. Chalio, R.A. McLauchlan, D.A. Clark, and S.I. Omar, "A Fuzzy Neural Hybrid System," *Proc. IEEE Int'l Conf. Neural Networks*, vol. 3, pp. 1,654-1,657, Orlando, Fla., 1994.
- [6] M.W. Craven and J.W. Shavlik, "Using Sampling and Queries to Extract Rules from Trained Neural Networks," *Machine Learning: Proc. 11th Int'l Conf.*, pp. 37-45, 1994.
- [7] L.M. Fu, "Knowledge-Based Connectionism for Revising Domain Theories," *IEEE Trans. Systems, Man, and Cybernetics*, vol. 23, no. 1, pp. 173-182, 1993.
- [8] L.M. Fu, *Neural Networks in Computer Intelligence*. McGraw-Hill, 1994.
- [9] J. Ghosh and K. Tumer, "Structural Adaptation and Generalization in Supervised Feed-Forward Networks," *J. Artificial Neural Networks*, vol. 1, no. 4, pp. 431-458, 1994.
- [10] C.L. Giles et al., "Extracting and Learning an Unknown Grammar with Recurrent Neural Networks," S.J. Hanson, J.E. Moody, and R.P. Lippmann, eds., *Advances in Neural Information Processing Systems—4*. San Mateo, Calif.: Morgan Kaufmann, 1992.
- [11] C.W. Glover, M. Silliman, M. Walker, and P. Spelt, "Hybrid Neural Network and Rule-Based Pattern Recognition System Capable of Self-Modification," *Proc. SPIE, Application of Artificial Intelligence*, vol. 8, pp. 290-300, 1990.
- [12] J.A. Hendler, "Marker-Passing Over Microfeatures: Towards a Hybrid Symbolic/Connectionist Model," *Cognitive Science*, vol. 13, pp. 79-106, 1989.
- [13] S. Horikawa, T. Furuhashi, and Y. Uchikawa, "On Fuzzy Modeling Using Fuzzy Neural Networks with Back-Propagation Algorithm," *IEEE Trans. Neural Networks*, vol. 3, no. 5, pp. 801-806, Sept. 1992.
- [14] P. Howes and N. Crook, "Rule Extraction from Neural Networks," R. Andrews and J. Diederich, eds., *Rules and Networks: Proc. Rule Extraction from Trained Artificial Neural Networks Workshop*, pp. 60-67, Queensland Univ. of Technology, Neurocomputing Research Center, Apr. 1996.
- [15] R.A. Jacobs, M.I. Jordan, S.J. Nowlan, and G.E. Hinton, "Adaptive Mixtures of Local Experts," *Neural Computation*, vol. 3, pp. 78-88, 1991.
- [16] R. Kerber, "ChiMerge: Discretization of Numeric Attributes," *Proc. 10th Nat'l Conf. Artificial Intelligence AAAI*, pp. 123-128, July 1992.
- [17] R.H. Klinkenberg and D.C. St. Clair, "Rule Set Quality Measures for Inductive Learning Algorithms," *Intelligent Eng. Systems Through Artificial Neural Networks ANNIE*, vol. 6, pp. 161-168. ASME Press, Nov. 1996.
- [18] C.T. Lin and C.S.G. Lee, "Neural-Network-Based Fuzzy Logic Control and Decision System," *IEEE Trans. Computers*, vol. 40, no. 12, pp. 1,320-1,326, Dec. 1991.
- [19] H. Liu and R. Setiono, "Chi2: Feature Selection and Discretization of Numeric Attributes," *Proc. Seventh Int'l Conf. Tools with Artificial Intelligence*, pp. 388-391, Nov. 1995.
- [20] J.J. Mahoney and R.J. Mooney, "Combining Connectionist and Symbolic Learning to Refine Certainty Factor Rule Bases," *Connection Science*, vol. 5, nos. 3-4, pp. 339-364, 1993.
- [21] O.L. Mangasarian and H.W. Wolberg, "Cancer Diagnosis via Linear Programming," *SIAM News*, vol. 23, no. 5, pp. 1-18, 1990.
- [22] C. McMillan, M.C. Mozer, and P. Smolensky, "The Connectionist Scientist Game: Rule Extraction and Refinement in a Neural Network," *Proc. 13th Ann. Conf. Cognitive Science Soc.*, 1991.
- [23] W. Mendenhall, *Introduction to Probability and Statistics*, fifth ed. Wadsworth, 1979.
- [24] P.M. Murphy and D.W. Aha, "UCI Repository of Machine Learning Database," technical report, Dept. of Computer Science, Univ. of California, 1992.
- [25] C.W. Omlin and C.L. Giles, "Extraction of Rules from Discrete-Time Recurrent Neural Networks," *Neural Networks*, vol. 9, no. 1, pp. 41-52, 1996.
- [26] D.W. Optiz and J.W. Shavlik, "Heuristically Expanding Knowledge-Based Neural Network," *Proc. 13th Int'l Joint Conf. Artificial Intelligence*, pp. 512-517, 1993.
- [27] D. Ourston and R.J. Mooney, "Changing the Rules: A Comprehensive Approach to Theory Refinement," *Proc. Eighth Nat'l Conf. Artificial Intelligence*, pp. 815-820. AAAI Press, 1990.
- [28] J.R. Quinlan, *C4.5: Programs for Machine Learning*, San Mateo, Calif.: Morgan Kaufman, 1992.
- [29] S. Ramachandran, "Theory Refinement of Bayesian Networks with Hidden Variables," PhD thesis, Dept. of Computer Science, Univ. of Texas at Austin, Dec. 1997.
- [30] R. Ruddle and A. Sangiovanni-Vincentelli, "Espresso-MV: Algorithms for Multiple-Valued Logic Minimization," *Proc. Cust. Int'l Circ. Conf.*, Portland, Ore., May 1985.
- [31] K. Saito and R. Nakano, "Medical Diagnostic Expert System Based on PDP Model," *Proc. IEEE Int'l Conf. Neural Networks*, vol. 1, pp. 255-262, 1988.
- [32] J.S. Schlimmer, "Concept Acquisition Through Representational Adjustment," PhD thesis, Dept. of Information and Science, Univ. of California at Irvine, May 1996.
- [33] S. Sestito and T. Dillon, "Automated Knowledge Acquisition of Rules with Continuously Valued Attributes," *Proc. 12th Int'l Conf. Expert Systems and Their Applications (AVIGNON)*, pp. 645-656, May 1995.
- [34] R. Setiono, "Extracting Rules from Pruned Neural Networks for Breast Cancer Diagnosis," *Artificial Intelligence in Medicine*, vol. 8, no. 1, pp. 37-51, Feb. 1996.
- [35] R. Setiono and H. Liu, "Understanding Neural Networks via Rule Extraction," *Proc. 14th Int'l Joint Conf. Artificial Intelligence (IJCAI)*, pp. 480-485, 1995.
- [36] R. Setiono and H. Liu, "Symbolic Representation of Neural Networks," *Computer*, pp. 71-77, Mar. 1996.
- [37] N.E. Sharkey and A.J.C. Sharkey, "Understanding Catastrophic Interference in Neural Networks," Technical Report CS-94-4, Dept. of Computer Science, Sheffield, U.K., 1994.
- [38] I. Taha, "A Hybrid Intelligent Architecture for Revising Domain Knowledge," PhD thesis, Military Technical College, Cairo, 1997.
- [39] I. Taha and J. Ghosh, "A Hybrid Intelligent Architecture for Refining Input Characterization and Domain Knowledge," *Proc. World Congress on Neural Networks (WCNN)*, vol. 2, pp. 284-287, July 1995.
- [40] I. Taha and J. Ghosh, "Hybrid Intelligent Architecture and Its Application to Water Reservoir Control," *Int'l J. Smart Eng. Systems*, vol. 1, pp. 59-75, 1997.
- [41] H. Takagi and I. Hayashi, "NN-Driven Fuzzy Reasoning," J.C. Bezdek and S.K. Pal, eds., *Fuzzy Models for Pattern Recognition*, pp. 496-512. IEEE Press, 1992.
- [42] E. Tazaki and N. Inoue, "A Generation Methods for Fuzzy Rules Using Neural Networks with Planar Lattice Architecture," *Proc. IEEE Int'l Conf. Neural Networks*, vol. 3, pp. 1,743-1,748, Orlando, Fla., 1994.
- [43] S.B. Thrun, J. Bala, E. Bloedorn, B. Cheng, I. Bratko, S. Dzeroski, K. De-Jong, S. Fahlman, D. Fisher, R. Hamann, K. Kaufman, S. Keller, I. Kononenko, J. Kreuziger, R. Michalski, T. Mitchell, P. Pachowicz, Y. Reich, H. Vafaie, K. Van de Welde, W. Wenzel, J. Wnek, and J. Zhang, "The Monk's Problem: A Performance Comparison of Different Learning Algorithms," Technical Report CMU-CS-91-197, Carnegie Mellon Univ., Dec. 1990.
- [44] A.B. Tickle, M. Orłowski, and J. Diederich, "DEDEC: A Methodology for Extracting Rules from Trained Artificial Neural Networks," R. Andrews and J. Diederich, eds., *Rules and Networks: Proc. Rule Extraction from Trained Artificial Neural Networks Workshop*, pp. 90-102, Neurocomputing Research Center, Queensland Univ. of Technology, Apr. 1996.
- [45] G.G. Towell and J.W. Shavlik, "The Extraction of Refined Rules from Knowledge-Based Neural Networks," *Machine Learning*, vol. 13, no. 1, pp. 71-101, 1993.
- [46] G.G. Towell and J.W. Shavlik, "Knowledge-Based Artificial Neural Networks," *Artificial Intelligence*, vol. 70, nos. 1-2, pp. 119-165, 1994.

- [47] G.G. Towell, J.W. Shavlik, and M.O. Noordwier, "Refinement of Approximate Domain Theories by Knowledge-Based Artificial Neural Network," *Proc. Eighth Nat'l Conf. Artificial Intelligence*, pp. 861-866, 1990.
- [48] A.B. Tickle, R. Andrews, M. Golea, and J. Diederich, "The Truth Will Come to Light: Directions and Challenges in Extracting the Knowledge Embedded within Trained Artificial Neural Networks," *IEEE Trans. Neural Networks*, vol. 9, no. 6, pp. 1,057-1,068, 1998.
- [49] K. McGarry, S. Wertmer, and J. MacIntyre, "Hybrid Neural Systems: From Simple Coupling to Fully Integrated Neural Networks," *Neural Computing Surveys*, vol. 2, pp. 62-93, 1999.



Ismail A. Taha received his BS and MS degrees from the Military Technical College (MTC), Cairo, Egypt, in electrical engineering and computer science, both with top honors and MTC fellowship. He was a lecturer at MTC until 1993, when he joined the graduate program at the University of Texas at Austin with an Egyptian Government Graduate Fellowship. After completing his PhD in May 1997 and a stint at Trilogy Development Group in Austin, Dr. Taha rejoined the faculty of MTC as an associate professor and graduate

coordinator. Dr. Taha has several refereed publications in the areas of intelligent systems and databases, including one that received the best application paper award (first runner up) at ANNIE 1995.



Joydeep Ghosh was educated at the Indian Institute of Technology at Kanpur (BTech, 1983) and the University of Southern California (MS, PhD, 1988). He is currently a full professor with the Department of Electrical and Computer Engineering at the University of Texas at Austin, where he holds the Endowed Engineering Foundation Fellowship. He directs the Laboratory for Artificial Neural Systems (LANS), where his research group is studying the theory and applications of adaptive pattern recognition and mul-

tilearner systems. Dr. Ghosh has published more than 125 refereed papers and edited six books. He received the 1992 Darlington Award for the Best Paper in the areas of CAS/CAD, besides "best conference paper" citations for five neural network papers. Dr. Ghosh served as the general chairman for the SPIE/SPSE Conference on Image Processing Architectures, Santa Clara, California, held in February 1990, as conference cochair of Artificial Neural Networks in Engineering (ANNIE) from 1993 through 1996, and in the program committee of several conferences on neural networks and parallel processing. He is an associate editor of *Pattern Recognition*, *IEEE Transactions on Neural Networks*, *Neural Computing Surveys*, and the *International Journal of Smart Engineering Design*, and as the letters editor for *IEEE Transactions on Neural Networks*. He was a plenary speaker for ANNIE in 1997. He is a member of the IEEE.