

EXPERIMENTAL COMPARISON OF SYMBOLIC AND
SUBSYMBOLIC LEARNING

by

J. Wnek
R. S. Michalski

HEURISTICS, The Journal of Knowledge Engineering, Special Issue on Knowledge
Acquisition and Machine Learning, Vol.5, No.4, pp. 1-21, 1992.

Experimental Comparison of Symbolic and Subsymbolic Learning

Janusz Wnek and Ryszard S. Michalski*

ABSTRACT—This paper reports on three studies comparing symbolic and subsymbolic methods for concept learning from examples. The first study compared five learning methods, three representing symbolic learning paradigm—decision tree learning (C4.5), rule learning (AQ15), and constructive rule learning (AQ17-HCI), and the other two representing the subsymbolic paradigm—neural net learning using backpropagation (BpNet), and a classifier system employing genetic algorithm (CFS). All methods have been experimentally applied to learn several different DNF-type concepts (i.e., concepts representable by a simple DNF expression). The second study compared performance of a large number of learning programs on learning DNF-type concepts from data with and without noise, and a non-DNF-type "m-of-n" concept. The third study compared genetic algorithm based learning (GABIL and Adaptive GABIL) with decision tree learning (C4.5) and decision rule learning (AQ14) on twelve DNF-type concepts. In all studies, symbolic methods, particularly those applying constructive induction, outperformed subsymbolic methods in learning DNF-type concepts from data both without and with noise. In cases of learning non-DNF-type concepts, symbolic methods without constructive induction performed worse, but those with constructive induction matched the performance of neural network methods, and also gave comprehensive concept descriptions.

Introduction

In view of a rapidly growing interest in multistrategy learning systems, it is important to develop insights into the performance of diverse learning methods and paradigms, and to determine areas of their most desirable applicability. To this end, this paper describes three studies comparing the performance of symbolic and subsymbolic methods on a sample of learning problems. The first study involved a comparison of three symbolic and two subsymbolic methods. Symbolic methods were represented by C4.5—a decision tree learning program, AQ15—a decision rule learning program, and AQ17-HCI—a constructive decision rule learning program. Subsymbolic methods were represented by CFS—a genetic algorithm based classifier system, and BpNet—a neural network learning program using a backpropagation algorithm. Other studies involved the same programs or their different variants.

An important difference between symbolic and subsymbolic learning approaches lies in the cognitive aspects of the knowledge representation. Knowledge represented by logic-based rules or decision trees (especially when the latter are small) is relatively easy to comprehend and relate to human knowledge. This is not the case with knowledge represented by classifier systems

or neural networks. While for some applications it may not be important that the learned concept descriptions are understandable to people, e.g., in an adaptive controller of house temperature, in some applications, e.g., in expert systems for human disease diagnosis, business or military decision making, this requirement is essential.

Despite various attempts, there is no established universal measure of cognitive comprehensibility of concept representations [1]. Therefore, we will make a simplifying assumption that the comprehensibility of a concept can be roughly estimated by the number of rules needed to express it, or the number of disjuncts in an equivalent DNF expression. In this measure, called the *R-complexity* (*rule-complexity*) of a concept, elementary conditions in the rules representing a concept (or the components of disjuncts in DNF) are assumed to be simple conditions on attribute values. Based on this definition, one can distinguish between two general classes of concepts:

- 1) Concepts expressable by a simple DNF expression (using the a priori given attributes) or described by few rules; called *DNF-type* concepts.
- 2) Concepts requiring a complex DNF expression or many rules to describe them; called *non DNF-type*.

It is important to point out that concepts that have a long DNF expression using given attributes, may have a short

*The authors can be reached at the Cen. for AI, 4400 University Dr. Fairfax, VA 22020, Wnek@aic.gmu.edu or Michalski@aic.gmu.edu

DNF expression if these attributes are replaced by other attributes, or transformed into combined attributes through the process of constructive induction [2]. Thus, whether a given concept is DNF-type or not depends on the attributes (generally, descriptors—attributes or functions whose values characterize the entity) that are available for constructing a concept representation. In other words, the R-complexity is defined with regard to the assumed concept representation space.

All three studies compared several methods by applying them to learning the same class of DNF-type concepts. We found that concepts generated by human subjects who are asked to create classes of entities, and to express them linguistically, usually fall into such category. Given a concept representation, its R-complexity can thus be viewed as an approximate indication of the "cognitive" complexity of the concept. For representations other than rule-based, the R-complexity can be determined by converting them to logically equivalent sets of rules. When the description spaces are not too large, this can be done using the DIAV concept visualization method, outlined in a later section.

Presented studies follow several other efforts to compare learning methods and paradigms. For example, Fisher and McKusick [3] compared ID3 and a neural net using backpropagation algorithm on the problems of learning diagnostic rules for thyroid diseases, soybean plant diseases, and a few artificial problems. The comparison was based on the performance accuracy of descriptions as applied to testing examples and the training time. Their conclusion was that the neural net gave a better performance, but required a significantly longer training time and more training examples than ID3.

Mooney et al. [4] compared ID3 with perceptron and backpropagation algorithms using the domain of soybean diseases, chess-end games, audiological disorders, and the Nttalk data set. Their conclusion was that the accuracy of classifying new examples was about the same for all the three systems, but the neural net performed better than ID3 when there was noise in the data. Weiss and Kapouleas [5] compared ID3, predictive value maximization, neural net using backpropagation, and a few statistical methods. They found that the statistical classifiers performed consistently better in terms of accuracy in classifying testing examples.

Dietterich, Hild, and Bakiri [6] compared ID3 with a neural net using backpropagation on the task of text-to-speech mapping. Their major conclusion was that the neural net consistently outperformed ID3 in terms of the performance accuracy, and attributed this result to the capture of better statistical information by the neural net.

Bergadano et al. [7] compared POSEIDON (also called AQ16, it is an extended version of AQ15 using a *two-tiered* concept representation) with exemplar-based and decision tree (ASSISTANT) learning programs. Their study involved two real-world domains: labor contracts and U.S. congressional voting. In this study, descriptions learned by POSEIDON outperformed those produced by all other methods, both in terms of performance accuracy on new examples and in terms of the description's simplicity.

The first study presented here differs from the above studies in that it experimentally analyzes five different methods. It also compares the learned descriptions in terms of their *exact error rate*, rather than a statistical error estimate, and in terms of their R-complexity. In the study, the target and learned concepts were represented graphically by a novel technique of *diagrammatic visualization* [8]. This technique permits one to display an *image* of the target and learned concepts and an *error image* that identifies all errors.

The paper consists of seven sections. The second briefly describes the five learning systems used in the first study. Next, the methodology used to compare the methods is presented. Also described are training data and the concepts to be learned (five DNF-type concepts created by human subjects). The concepts are illustrated by the diagrammatic visualization technique (DIAV). Following this, results of experiments with the methods are described. Next, two related studies done by other research groups are summarized. The first one involved three types of problems, learning a DNF-type concept, learning a non-DNF-type concept, and learning a DNF-type concept from noisy data. The results were obtained using a large number of learning programs, that were grouped into four categories according to the representational paradigms used: decision tree, decision rules, neural networks, and Horn clauses. This study was coordinated by S. Thrun and involved several Artificial Intelligence laboratories both in USA and in Europe [9]. The second study applied a decision tree learning program, a decision rule learning program, and two genetic algorithm-based programs to learn twelve DNF-type concepts. The study was done by Spears and Gordon [10]. Results from comparison of learning systems are summarized and problems for further research indicated.

Learning Systems Involved in the First Study

As mentioned earlier, the symbolic paradigm is represented by a decision tree learning program, C4.5, and two rule learning programs, AQ15 and AQ17-HCI. The

subsymbolic paradigm is represented by a backpropagation neural network, BpNet, and a classifier system based on a genetic algorithm, CFS. These programs are widely known and well described in machine learning literature. To serve the tutorial purpose of this paper, we provide here a brief account of the basic algorithm underlying each program, and give references to the literature for readers interested in further details.

Decision Tree Learning Program C4.5

C4.5 learns concepts by building a decision tree that classifies supplied training examples of the concepts. Each interior node of the tree is assigned an attribute, while the leaf nodes are assigned concept names. A branch down from an interior node represents a value of the attribute assigned to the node. Any path from the root to a leaf in the tree can be viewed as a decision rule for the class assigned to that leaf.

The input to the algorithm consists of sets of training examples for different concepts (or decision classes). In the first step, the algorithm selects a random subset of training examples from each set (a "window"). Then, for each attribute it evaluates the information gain, i.e., the information gained if the attribute were chosen for testing. The attribute with the highest score is assigned to the root node of the tree. Branches from this node represent different values of this attribute. End-nodes of these branches (current leaves) are assigned subsets of examples in which the attribute takes the value associated with the given branch. If a subset contains examples of only one decision class, then the end-node becomes a leaf of the decision tree. For all other subsets, the algorithm is repeated, until all leaves in the tree are assigned single decision classes.

At this point, the created tree correctly classifies all examples in the window. Now the tree is used to classify remaining examples from the training set (outside the "window"). If the tree gives the correct answer for all examples, then the process terminates. If not, misclassified examples are added to the window and the process is repeated until the trial decision tree correctly classifies all examples not in the window. After the entire decision tree has been generated, C4.5 recursively examines each subtree to determine whether replacing it with a leaf or a branch would reduce the number of errors. An operation of replacing a subtree with a leaf or branch is called tree pruning. Both decision trees, unpruned and pruned, are evaluated on testing (unseen) examples.

The C4.5 program [11] is a derivative of the inductive decision tree (ID3) program [12]. In addition to decision

tree generating, C4.5 is able to convert an unpruned decision tree into sets of generalized (pruned) decision rules. A tree is converted to rules by forming a rule corresponding to each path from the root of the tree to each of its leaves. All rules are then examined and some of them are generalized (pruned) by dropping conditions. Next, rules for each class are considered separately and redundant rules are sifted out. For uncovered examples, a default class is assigned.

We have tested all three representations learned by C4.5 using default parameter setting, i.e., the best tree was selected out of ten generated from the same training set, attributes were selected according to "gainratio" (ratio of information gain to potential information) criterion, etc. As expected, since the training examples did not have noise, on average, unpruned decision trees performed best in terms of predictive accuracy. Pruned decision trees and decision rules were simpler but more erroneous. Thus, we here report the results obtained for unpruned decision trees only.

Rule Learning Program AQ15

The AQ15 program generates concept descriptions from concept examples. The descriptions are in the form of decision rules expressed in an attributional logic calculus, called *variable-valued logic* VL1 [13]. A distinct feature of this representation is that it employs, in addition to standard logic operators, also the *internal disjunction* operator (a disjunction of values of the same attribute), which can significantly simplify rules involving multivalued attributes. The program can optimize the rules according to user-defined (or default) preference criteria, such as the overall simplicity or the evaluation and/or storage cost of the rules. The main procedure of AQ15 is based on the AQ algorithm that builds a concept description from a set of positive and negative examples (e.g., [13]). Below is a simplified version of the AQ algorithm:

- Randomly select a *seed* example from the set of positive training examples of the concept to be learned.
- Generate a set of most general rules (a *star*) that cover the seed and do not cover any negative examples. (This operation employs the *extension against* generalization operator [1].)
- Select the "best rule" from the star (according to the assumed preference criteria), and remove examples covered by this rule from the set of training examples.

- If the set of training examples does not become empty, return to the first Step. Otherwise, the obtained set of rules constitutes a complete and consistent concept description with respect to the input examples.

The algorithm is repeated for each concept to be learned. It is biased toward finding a conjunctive description of a concept (a single rule), because if such a description exists for the given set of examples, it will be found in the very first step. (The description will be a member of the first star generated.) The AQ15 program has various parameters whose default values can be changed by a user according to the requirements of the domain. In the experiments reported here, the preferred criterion was to minimize both the number of rules describing the concepts, and the number of conditions in them. Since training examples did not have noise, no rule truncation was performed. For further details see [7, 14].

Constructive Rule Learning Program AQ17-HCI

AQ17-HCI represents a recent major advance in the development of the AQ-based series of inductive learning programs, specifically, the above-described AQ15 system. The main new feature is an incorporation of a method for *hypothesis-driven constructive induction* (HCI). Constructive induction, as introduced by Michalski in 1978 [15], addresses the problem of changing the representation space so that it is more adequate for the learning problem at hand. This involves creating new attributes (or descriptors) that better characterize the concepts to be learned than the original descriptors. The last few years have witnessed an increasing interest in constructive induction methods, since they can produce concept descriptions that are more accurate and/or simpler than those produced by the traditional *selective* induction methods [2, 16, 17].

The HCI method generates problem-relevant descriptors by analyzing consecutively created inductive hypotheses [2]. Below is a brief description of the algorithm used in AQ17-HCI. For the sake of simplicity, it is assumed that the training set consists of subsets of positive examples E^+ , and negative examples, E^- . If a training set consists of subsets representing different concepts, then E^+ represents the subset of training examples for the concept under consideration, and the union of the remaining subsets plays the role of E^- .

HCI Algorithm

1. Randomly divide each of the training sets, E^+ and E^- ,

into two subsets: the primary set P and the secondary set S . $E^+ = P^+ \cup S^+$, $E^- = P^- \cup S^-$. (The primary training set is used for the representation space modifications. Both primary and secondary sets are used for final rules generation).

2. For each concept, induce the most specific (ms) cover of the set P , against the set P . (Such a cover is denoted $COV_{ms} P/P$.) It represents the set of the most specific rules that together characterize examples in P , but no examples in P^- .
3. Evaluate the performance of the rules on the secondary training set, S . If the performance exceeds a predefined threshold, or all changes in the representation space were exhausted, go to Step 8.
4. Analyze the rules in order to identify possible changes in the representation space.
5. Change the representation space by removing irrelevant attribute values or attributes, or by adding new attribute values or attributes.
6. Modify the training set of examples, E , according to the changes in the representation space.
7. Go to Step 2.
8. For each concept, induce a set of the most specific rules from all positive examples against all negative examples, i.e., a cover $COV_{ms} (E^+/E^-)$, and the most general cover of negative examples against positive examples $COV_{mg} (E^-/E^+)$.
9. Build final concept descriptions by generalizing the most specific positive rules against the most general negative rules, i.e., $COV_{mg} (COV_{ms} (E^+/E^-) / COV_{mg} (E^-/E^+))$.

The AQ17-HCI program has two important moved features. The first one is an ability to change the representation space using HCI. This means that the method uses additional knowledge transmutations allowing abstraction and concretion, apart from inductive generalization and specialization [18]. The second feature is an extended generalization heuristic, employed in steps 8 and 9, that additionally generalizes the most specific rules characterizing the training set. Instead of producing the most specific rules or most general rules (consistent with the input examples), the obtained rules represent an intermediate degree of generalization. This extension was proposed by Wnek [19].

Neural Net Program BpNet

A neural network is a structure of interconnected processing units that define a problem representation space. They can be of three types: input, output, and hidden. From the classification task point of view, the input units provide means for representing instances of a concept, while the output units denote the concept. The hidden units are optional and can be organized in layers. They provide communication links between input and output units in the task of translating input example into output classification. The units are interconnected and each connection has its weight. The connections are directed from input units, through hidden, to output units. Some networks, however, may have feedback loops. Each connection is subject to changes in the learning process.

Backpropagation, as originally introduced by Rumelhart, Hinton and Williams [20], is a learning algorithm for feed-forward networks (networks in which the interconnections form no feedback loops) based on gradient minimization. We consider a network of units in which a weighted sum of the inputs is performed, the result of this sum (also called the activation level of the unit) is being fed through a non-linear element, with a sigmoid input-output function.

Learning by backpropagation involves two phases. During the first phase an example is presented and propagated forward through the network to compute the output values o_n for each output unit. These outputs are then compared with the target values t_n , resulting in output errors e_n for each unit. The second phase involves a backward pass through the network (analogous to the initial forward pass) during which the error message is passed to each unit in the network and the appropriate weight changes are made.

The two phases are repeated until the network reaches stable equilibrium, i.e., the overall error reaches a predefined level. The output error for a given training example is given by

$$e_n = o_n - t_n$$

where o_n and t_n are the output and the target values of the output unit n . The total squared error for that input example is

$$E = \sum_{n \in U} e_n^2$$

where U denotes the set of output units.

Thus, learning by backpropagation corresponds to gradient minimization of the average squared error. The average is computed over all examples in a given training set. The

BpNet program is an implementation of the backpropagation algorithm [21].

Classifier System CFS

A classifier system is a parallel rule-based (production) system. The rules, called classifiers, have the same form, so it is easy to determine whether a condition part of a rule is satisfied. Since the rules can be active simultaneously, complex situations are expressed by combinations of rules. The classifiers can be modified by a general-purpose learning system.

Classifier system learning and classification is done in cycles. In each cycle, an input example is translated into a message that has the same form as the condition part of the rules. The message is compared with all rules. All matched rules compete with each other in order to become active and to produce new messages. The new messages can either store some intermediate information and then be used in the next cycle, or produce a classification, if matched with the system's effectors. The classification is compared with the target class of the example and payoff is distributed among active rules. It means that given a set of classifiers, the behavior of the classifier can be changed by changing the strengths associated with those classifiers. In particular, the strength of classifiers that cause system's correct behavior can be increased; and the strength of classifiers that deteriorate the performance of the system can be decreased.

The learning process can be supplemented by utilizing genetic algorithms, a class of adaptive search techniques. "Genetic algorithms derive their name from the fact that they are loosely based on models of genetic change in a population of individuals. These models consist of three basic elements: (1) a Darwinian notion of "fitness," which governs the extent to which an individual can influence future generations; (2) a "mating operator," which produces offspring for the next generation; and (3) "genetic operators," which determine the genetic makeup of offspring from the genetic material of the parents" [22].

Classifier systems were first introduced by Holland and Reitman [23, 24]. The shell for the classifier system used in the experiments was developed by Riolo [25]. The CFS package of subroutines and data structures is domain independent and provides routines to perform the major cycle of the classifier system. The CFS system was run in the stimulus-response mode, i.e., without generating internal messages. Training cycles were repeated fifty times for each example. Payoff for correct and incorrect answer was set to 6 and -1, respectively, with a full payoff paid to all active classifiers. Final classification was

produced by two effectors. The CFS package uses more than 150 control parameters. The population size of sixty classifiers, number of training cycles, payoff, and about 20 other parameters were determined experimentally. The remaining parameters were set to default values.

Methodology

The testing domain in this study is the world of robot-like figures in the EMERALD system [26] (Emerald is a large-scale system integrating several different learning programs for the purpose of education and research in machine learning. It was developed at the Center for Artificial Intelligence at George Mason University. An earlier version, ILLIAN, was developed by the second author and his collaborators at the University of Illinois at Urbana-Champaign and is demonstrated at the Boston Museum of Science). For simplicity's sake, the robots are described by just six multivalued attributes (Figure 1-A). The attributes are Head Shape, Body Shape, Smiling, Holding, Jacket Color, and Tie, and can have 3, 3, 2, 3, 4, and 2 values, respectively. Consequently, the size of event space (the space of all possible robot descriptions) is $3 \times 3 \times 2 \times 3 \times 4 \times 2 = 432$. The space of

all possible concepts in this space is $2^{432} - 1 (\approx 10^{143})$. Undergraduate computer science students, unfamiliar with machine learning, were asked to create five concepts characterizing subsets of imaginary robots from a predefined set of sixteen robots in the EMERALD system. Below are descriptions of the five concepts ("Target concepts") used in the experiments (the numbers in parentheses represent the total numbers of positive and negative examples, respectively):

- C1: *Head is round and jacket is red or head is square and is holding a balloon* (84 positive, 348 negative)
- C2: *Smiling and is holding balloon or smiling and head is round* (120 positive, 312 negative)
- C3: *Smiling and not holding sword* (144 positive, 288 negative)
- C4: *Jacket is red and is wearing no tie or head is round and is smiling* (117 positive, 315 negative)
- C5: *Smiling and holding balloon or sword* (144 positive, 288 negative)

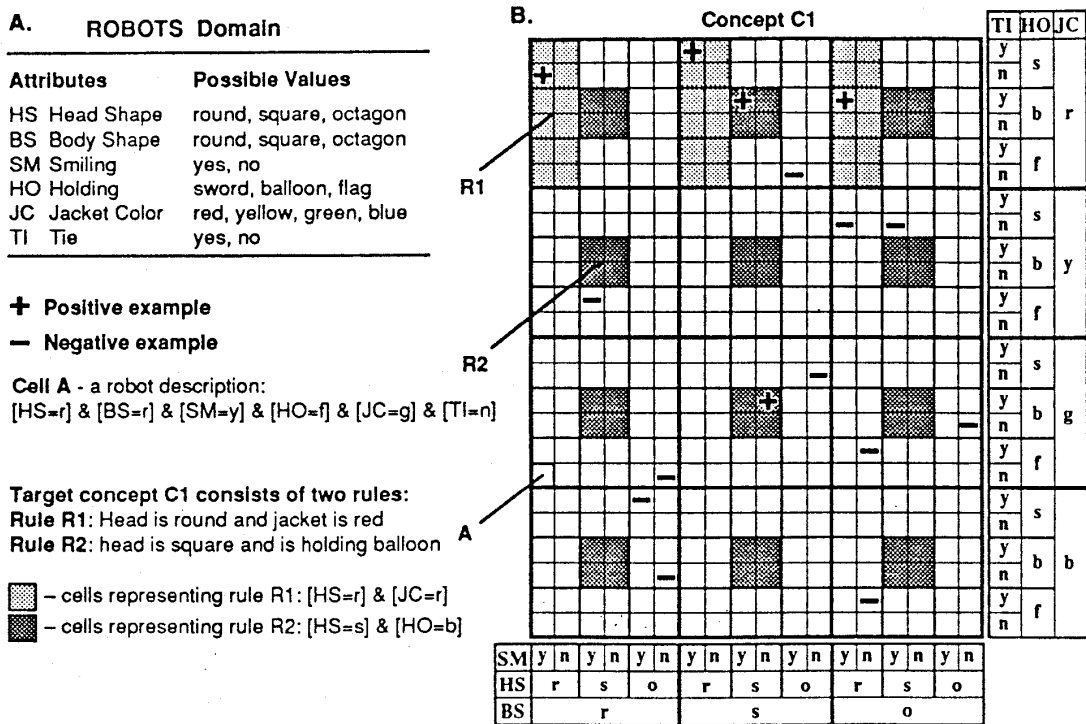


Figure 1. The Target Concept C1 and the Initial Training Examples in the ROBOTS Domain

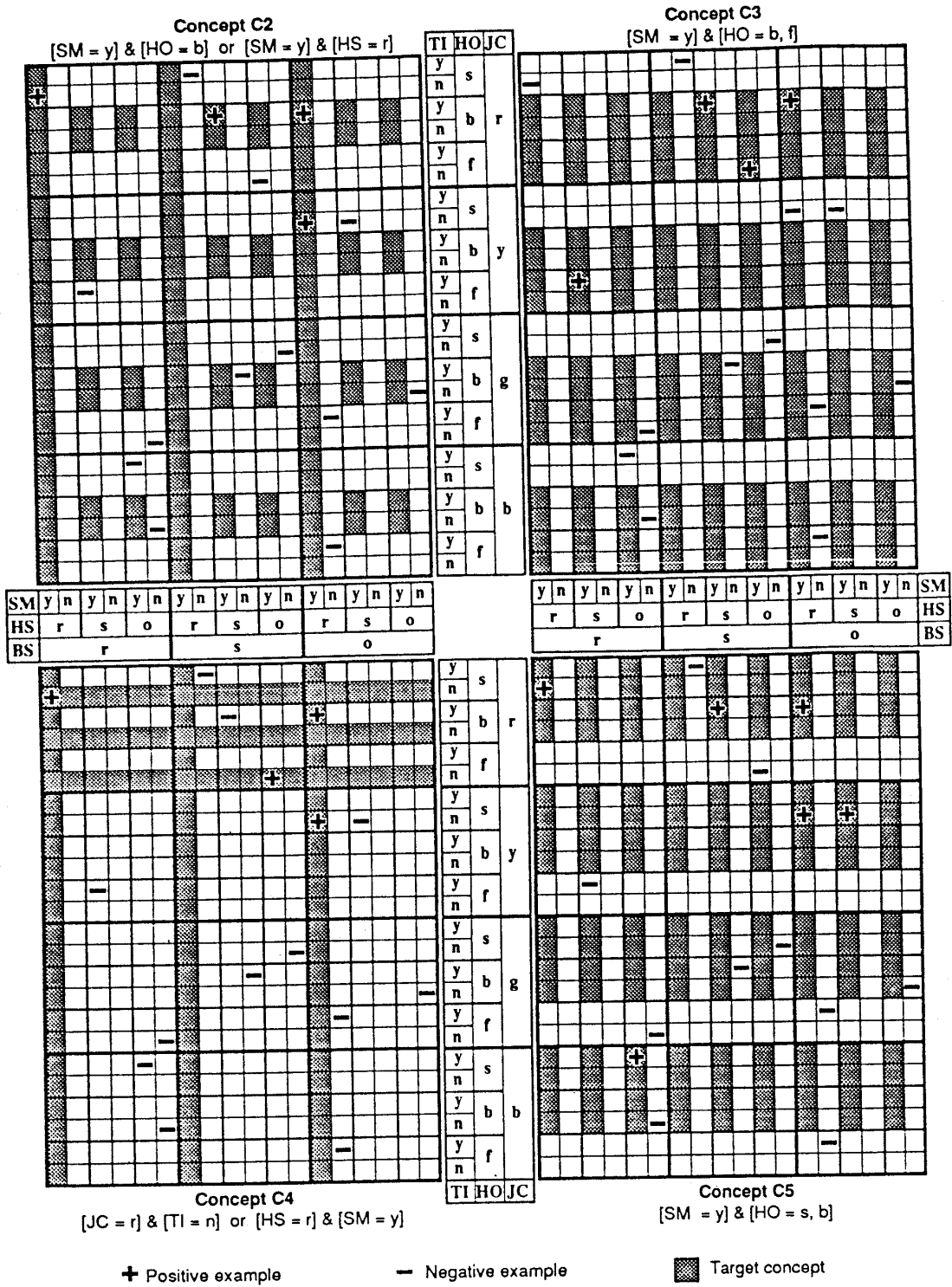


Figure 2. A Visualization of the Target Concepts C2, C3, C4 and C5, and the Initial Training Examples for Each Concept

Each such concept represents a partitioning of the event space into robots that belong to the concept (positive examples) and those that do not (negative examples). Based on the concepts C1-C5 the students generated initial sets of training examples used in Experiment 1. Each initial training set consisted of approximately 5 positive examples and 11 negative examples (6% of all possible positive examples and 3% of all possible negative examples). The remaining sets for Experiments 2-5 were generated by adding to the initial set an appropriate number of randomly generated examples. The average total sizes of the training sets in Experiments 2-5 are following: 12 positive and 30 negative (10% and 10%), 18 positive and 30 negative (15% and 10%), 30 positive and 30 negative (25% and 10%), 122 positive and 30 negative (100% and 10%). The additional experiments were performed in order to observe the convergence of the learned concepts to the target concepts.

The concepts C1-C5 are presented graphically in Figures 1-B and 2 using a method for *diagrammatic visualization*. This method employs a *General Logic Diagram* (GLD) which is a planar representation of a multi-dimensional space spanned over multivalued discrete attributes. (The system DIAV implementing the visualization method permits one to directly display description spaces up to 10^6 cells—e.g., about twenty binary attributes.) Larger spaces can also be displayed, but their representations have to be projected to subspaces [8, 13]. Each cell in the diagram represents a combination of the attribute values, e.g., a concept example. For example, the cell A in Figure 1-B represents the following robot description:

Head Shape = round, Body Shape = round, SMiling = yes, HOLDing = flag, Jacket Color = green, TTe = no

Positive and negative training examples are marked with + and -, respectively. Concepts are represented as sets of cells. The concept C1 can be described by two rules:

R1: Head Shape is round *and* Jacket Color is red

R2: Head Shape is square *and* is HOLDing balloon

The rules are represented in the diagram by shaded areas marked R1 and R2.

The diagrammatic visualization method permits one to display the target and learned concepts, individual steps in a learning process, and the errors in learning. The set of cells representing the target concept (the concept to be learned) is called *target concept image* (T). The set of cells representing the learned concept is called *learned concept image* (L). The areas of the *target concept* not

covered by the *learned concept* represent *errors of omission* ($T \setminus L$), while the areas of the learned concept not covered by the target concept represent *errors of commission* ($L \setminus T$). The union of both types of errors represents the *error image*. In the diagrams, errors are marked by slanted lines.

Figure 3 explains the meaning of various cases in concept visualization. Concept images are represented in the diagrams by shaded areas, e.g., Figure 3-A. If the target and learned concepts are visualized in the same diagram, then the shaded areas represent learned concept (Figure 3-B). Error image is represented by slanted areas. It is easy to distinguish between errors of omission and errors of commission. Since errors of commission are part of a learned concept, corresponding areas on the diagram are both shaded and slanted. Errors of omission are not part of the learned concept, thus the corresponding slanted areas remain white in the background. The location of the target concepts is implicitly indicated by correctly learned concept and errors of omission. The parts of the target concept that were correctly learned are shaded only.

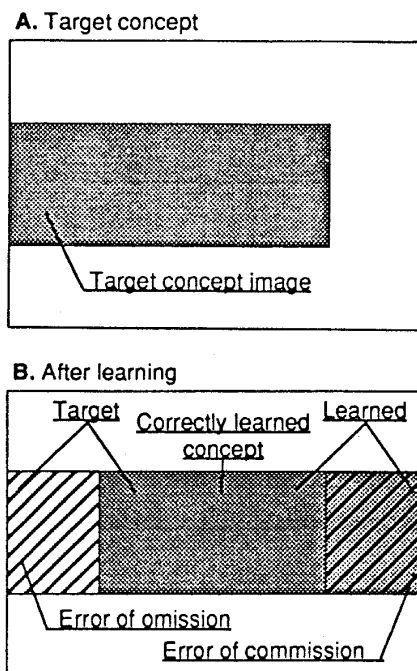


Figure 3. Interpretation of Various Areas in the Diagrammatic Visualization

The descriptions learned by the methods were compared in terms of the *exact error rate*. *Exact error rate* is the

ratio between exact error and the size of event space. It is measured as a function of the number of training examples. *Exact error* is defined as the total number of errors of omission and errors of commission, or equivalently the cardinality of the set-difference between the union and the intersection of the target and learned concepts.

$$Exact_error_rate = \frac{Exact_error}{\#Event_space}$$

where

$$Exact_error = \#[(T \setminus L) \cup (L \setminus T)] = \#[(T \cup L) \setminus (T \cap L)]$$

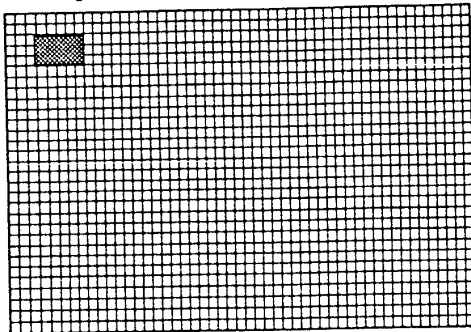
where

($T \setminus L$) — error of omission,
 ($L \setminus T$) — error of commission

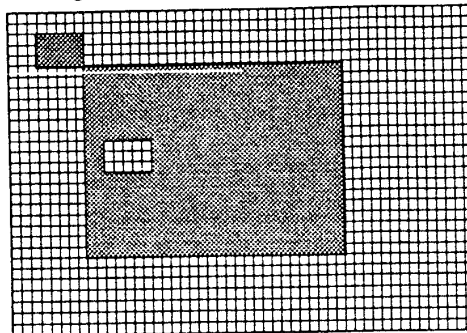
Equation 1.

There are many ways to define error rates in order to characterize learning capabilities of a system. Here are three assumptions related to the Equation 1. Firstly, for simplicity, we do not make any distinction between errors of omission and errors of commission, which may be important in some real-world domains. Secondly, Equation 1 indicates that if the event space were large and the target concept were relatively small, then the error rate would always be small, and thus, not sensitive to learning errors. Figure 4 illustrates two cases where Equation 1 would give the same error rate for different learning results. In the first case (Figure 4-A), a system did not learn any part of the target concept and still maintained 2% error rate, or, in other words, it was 98% accurate! Since the correct performance of a system is artificially increased by the system's performance on non-examples of a concept, the error estimating method is subjected to the Hempel's paradox [27, 28].

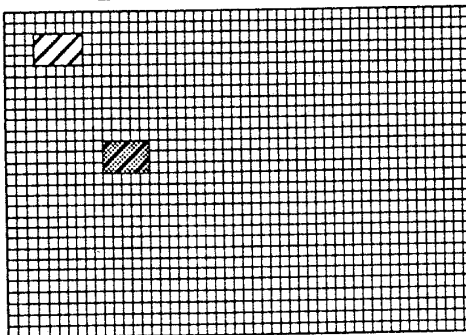
A. Target concept A — "small concept"



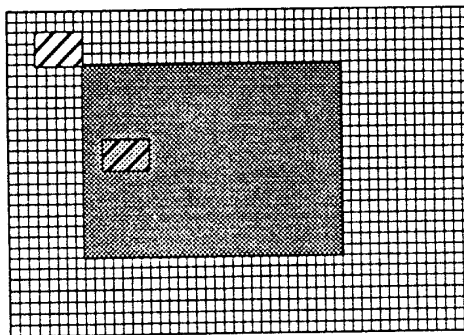
B. Target concept B — "large concept"



A'. Concept A is not learned at all:
 Exact_error_rate = (2%)



B'. Large portion of concept B is learned:
 Exact_error_rate = (2%)










 Errors of omission
  Errors of commission
  Correctly learned concept
 +  Target concept
  +  Learned concept

Figure 4. The Dependence of the Error Rate Definition on the Size of the Target Concept

In the second case (Figure 4-B), where the target concept occupies relatively large portion of the event space, the 2% error rate intuitively reflects the true performance of a system. In our experiments, the description space is small and the target concepts cover approximately 30%.

The third assumption related to Equation 1 is that in order to get complete insight into the performance of the tested methods, we used all examples from the event space to test the performance. In other studies, training examples might have been excluded from the testing phase.

In addition to the exact error rate, we used *R-complexity* (rule complexity) measure of a method performance. The *R-complexity* of a concept representation is defined by the number of conjunctive statements (rules) in the minimal DNF expression that is logically equivalent to the given representation. Since finding such a minimal DNF expression for any given representation may be difficult (it is generally an NP-hard task), we use an estimate of the *R-complexity*. For a method that learns a rule-based representation, the number of rules generated by the method is taken as such an estimate. For example, the *R-complexity* of the C1-C5 target concepts is 2, 2, 1, 2 and 1, respectively. For a decision tree learning method, the *R-complexity* is estimated by the number of leaves in the tree, since each leaf corresponds to a rule. For neural nets and classifiers, the *R-complexity* is estimated by determining the number of conjunctive statements needed to re-express the learned concept as a DNF expression.

Experiments in the Robots Domain

Representations Learned

Figure 5 presents an example of representations learned by each method. In the figure, the representations were learned in Experiment 1 from 6% of all positive and 3% of all negative examples of the target concept C1:

Head is round and jacket is red or head is square and is holding a balloon.

A Decision Tree Generated by C4.5

Figure 5-A shows the best, unpruned decision tree selected out of ten different trees generated from the training set. The learned concept is described using two attributes: Jacket Color and Head Shape. The learned concept can be read as follows:

IF Jacket Color is red and Head Shape is round or
Jacket Color is red and Head Shape is square or
Jacket Color is green and Head Shape is square

THEN C1

The exact error rate is 16.7% and the *R-complexity* of this tree is 3. After pruning, the tree is reduced to a root labeled \sim C1. Such a tree classifies all examples as not belonging to concept C1 and thus produces 84 omission errors (19.4% error rate). The *R-complexity* of this tree is 1. The third representation learned by C4.5 is decision rules obtained from the unpruned decision tree (Figure 5-A). After rule pruning and simplification, the final outcome consists of two rules:

- 1) IF Head Shape is octagonal THEN \sim C1
- 2) DEFAULT CLASS is \sim C1

These decision rules are equivalent to the pruned decision tree and produce the same error. The *R-complexity* is 2.

A Decision Rule Generated by AQ15

The method generated one rule. It consists of three conditions. Each condition tests values of one attribute. The internal disjunctions simplify the rule (Figure 5-B).

A Decision Rule Generated by AQ17-HCI

The rule learned by AQ17-HCI was exactly equivalent to the target concept (Figure 5-C). It was generated in a transformed, smaller description space. Figure 6 shows steps in learning concept C1 by AQ17-HCI. The input to the method is a set of training examples in the original representation space as shown in diagram A (the diagram also shows the target concept). The method divides the training set into primary and secondary examples and employs the AQ15 learning algorithm to induce rules from the primary set of training examples (diagram B). Since the performance test on the secondary training set is not satisfactory, the representation space is reduced to contain relevant attributes only, i.e., those attributes that are present or significant in the induced hypothesis. The method changed ROBOTS original representation space by removing three irrelevant attributes: Body Shape, SMiling, and THe (diagram C). The new representation space implied changes in the event space so the number of training examples was decreased by 1. It is due to the fact that two positive examples, E1 and E2, from the original event space, have the same description in the new event space.

E1: (round, round, yes, sword, red, no)

E2: (round, square, yes, sword, red, yes)

Although such an abstracted problem is simpler for

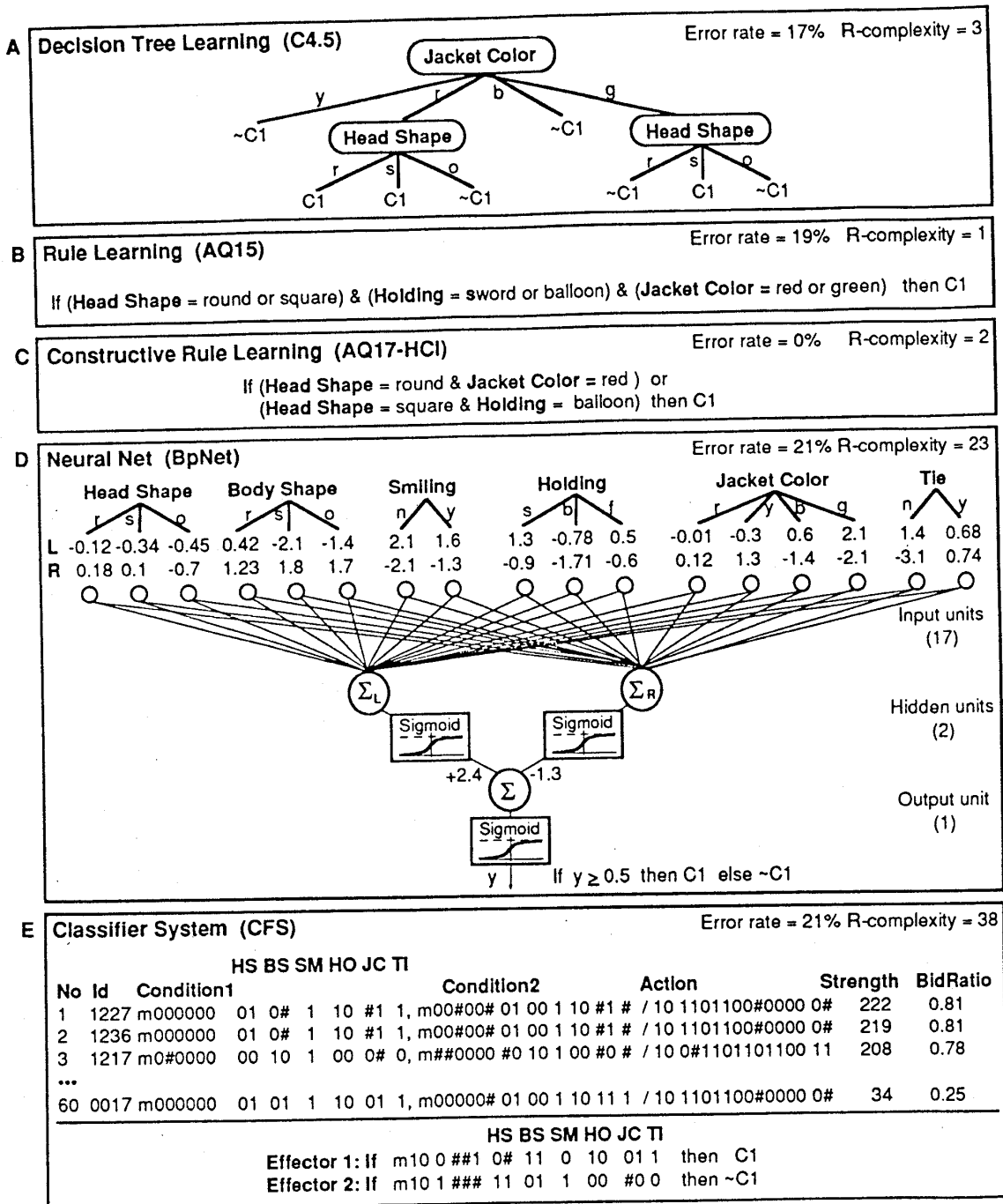
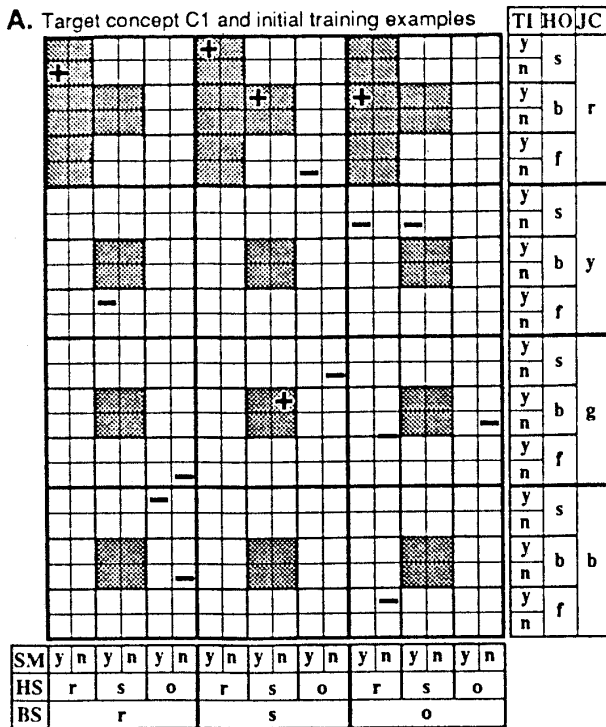


Figure 5. Representations of the Concept C1 Learned by Different Methods (From the Initial Set of Examples Consisting of 6% of Positive and 3% of Negative Examples)

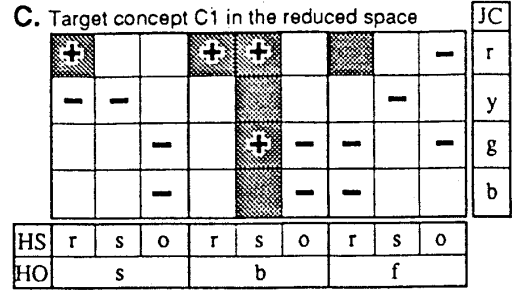
A. Target concept C1 and initial training examples



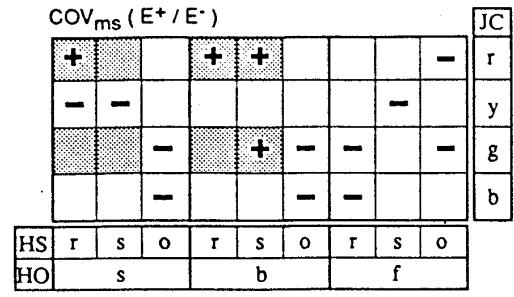
B. Concept C1 as learned from the set of primary training examples

-

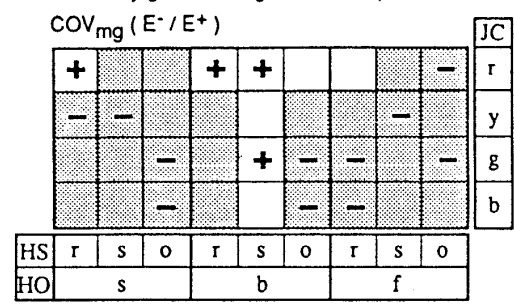
C. Target concept C1 in the reduced space



D. Maximally specific "positive" concept:



E. Maximally general "negative" concept:



F. Final concept learned (no errors)

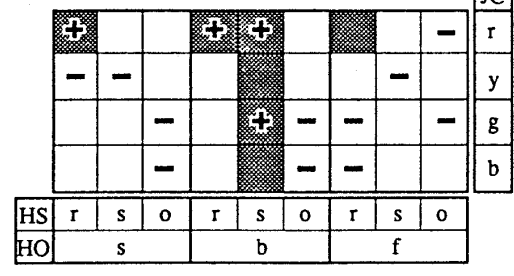


Figure 6. Steps in Learning Concept C1 by AQ17-HCI

learning, the resulting hypothesis is still not accurate (diagram D). At this point, the training data set seems to be insufficient to allow proper learning. The lacking information can, however, be induced while taking into consideration both positive and negative hypotheses. Figure 6, diagrams D and E show two covers, $COV_{m_+}(E^+/E^-)$ and $COV_{m_-}(E^+/E^-)$, that were generated using all initial training examples. AQ17-HCI generalized the positive concept description against the negative concept, and, by this means, improved the learned concept. The concept C1 was learned precisely (Figure 6, diagram F).

A Neural Net Generated by BpNet

Figure 5-D shows an architecture of the neural net used in the experiments. There were seventeen input units, all having either value 0 or 1, corresponding to attribute-values. All input units had connections to two hidden units. The number of hidden units was determined experimentally. The two hidden units were connected to one output unit. The network was trained by the BpNet backpropagation algorithm until it reached a root mean square error below 0.0007. The figure also shows final connection weights for the concept C1. Due to the space limitations, the connections from the input units to the left hidden unit and the right hidden unit are specified in the rows marked L and R, respectively. The weights from the hidden units to the output unit are 2.4 and -1.3. An input example is classified as a class member if it is translated into output value greater or equal 0.5.

Classifiers Generated by CFS

Each line in Figure 5-E represents one classifier in the following format: No, Id, Classifier, Strength, and BidRatio [25]. The total population for representing the concept consists of 60 classifiers. Each of the classifiers (condition-action rules) is in the following form: condition1, condition2 / action. Each condition consists of a string of a fixed length (16), built from the tertiary alphabet {0, 1, #}. A condition string with prefix "m" is matched by any message that has 0's and 1's in exactly the same positions as the 0's and 1's in the condition string. The # in the condition is considered as a "wildcard" symbol that can match a 0 or a 1. A classifier's condition-part is satisfied when both of its conditions are matched. When the condition-part of a classifier is satisfied, the classifier becomes active, i.e., its action-part produces one output message. The # in the action part has different meanings. Here it plays a role in the "pass through" in the sense that whenever it occurs in the action part, the corresponding bit in a message is passed through into the produced message. The messages generated by active classifiers are compared to effectors in order to produce

final classification. In Figure 5-E, *BidRatio* is a number between one and zero that is a measure of the classifier's specificity, i.e., how many different messages it can match. *Strength* is meant to be a measure of a classifier's "usefulness" to the system, thus the higher a classifier's strength, the more it bids.

Summary of Results

Figures 6 and 7 present the results of learning concept C1 by the five learning systems using diagrammatic visualization. In comparison to the representations in Figure 5, these diagrams give a uniform image of the learning results. From the diagrams, one can easily determine learning accuracy (correct vs. error areas—black vs. shaded areas) and interpret the errors (why certain areas were covered or not). Most importantly, one can generate rules equivalent to the learned representation and determine a R-complexity of the description. This feature is especially useful for subsymbolic systems that do not have easily understood knowledge representation, as shown in Figure 5.

The final concept description learned by AQ17-HCI exactly matches the target concept and, thus, there are no slanted areas in diagram F, in Figure 6. The other four methods did not learn the concept C1 precisely; however, all the methods were consistent with the training examples (Figure 7). The error rate level is almost even for all of them (about 20%), but one can note differences in their generalization patterns. The symbolic methods yield regular, rectangular covers, as opposed to irregular covers of subsymbolic methods.

Tables 1 and 2 summarize the results of all the experiments. For each learning program, the final result in Experiment 1 is an average over results from learning the five concepts from their initial training sets (column 1). In the remaining experiments, since additional examples were generated randomly, the testing was repeated 10 times for each concept. Consequently, for each learning program, the result is an average from 50 learning sessions (cols. 2-5). Pairs (a,b) in the top row of the tables denote the percentage of positive and negative examples used in experiments.

Table 1 shows the average exact error rate of the descriptions learned in five experiments, and Figure 8 presents corresponding learning curves. The error rate of the CFS-generated descriptions was much higher than that of the other descriptions, and what is most surprising, it did not improve much with the growth of the training sets. Differences between decision tree learning (C4.5), neural network (BpNet), and decision rule learning (AQ15) are

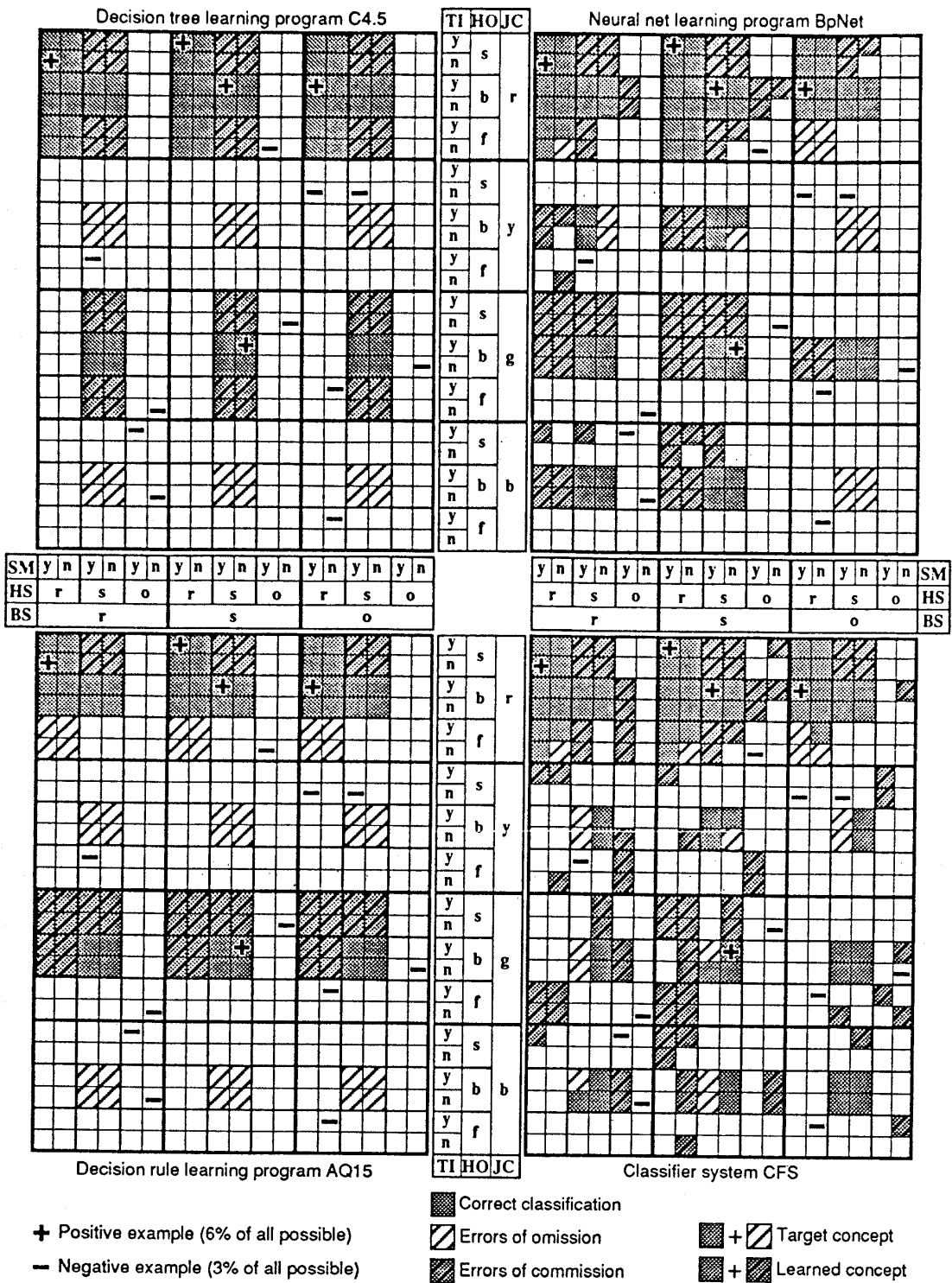


Figure 7. Concepts Learned by Different Methods in the Relation to the Target Concept C1

	Percentage of Training Examples Used in Each Experiment				
	Experiment 1 (6%, 3%)*	Experiment 2 (10%, 10%)*	Experiment 3 (15%, 10%)*	Experiment 4 (25%, 10%)*	Experiment 5 (100%, 10%)*
Genetic Alg. (CFS)	21.3%	20.3%	22.5%	19.7%	16.3%
Neural Nets (BpNet)	9.7%	6.3%	4.7%	7.8%	4.8%
Decision Trees (C4.5)	9.7%	8.3%	1.3%	2.5%	1.6%
Decision Rules (AQ15)	22.8%	5.0%	4.8%	1.2%	0.0%
Decision Rules (AQ17-HCI)	4.8%	1.2%	0.0%	0.0%	0.0%

Table 1. The Average Error Rate of Learned Descriptions

* In each (x%, y%), x denotes percentage of positive training examples selected from all possible positive examples, and y—negative training examples selected from all possible negative examples

	Percentage of Training Examples Used in Each Experiment				
	Experiment 1 (6%, 3%)*	Experiment 2 (10%, 10%)*	Experiment 3 (15%, 10%)*	Experiment 4 (25%, 10%)*	Experiment 5 (100%, 10%)*
Genetic Alg. (CFS)	49	45	51	48	41
Neural Nets (BpNet)	35	26	12	22	12
Decision Trees (C4.5)	3.1	2.8	2.5	2.5	2.5
Decision Rules (AQ15)	2.6	2.2	2	1.6	1.6
Decision Rules (AQ17-HCI)	2.4	2.0	1.6	1.6	1.6

Table 2. Numbers of Rules Representing Concepts Learned by Different Methods (R-complexity)

* In each (x%, y%), x denotes percentage of positive training examples selected from all possible positive examples, and y—negative training examples selected from all possible negative examples

relatively small, although only AQ15 precisely learned all concepts in Experiment 5. The 4.8% average error rate of the BpNet-generated concepts was primarily due to an inadequate learning of concepts C1 and C4. Also, decision trees generated by C4.5 produced some error even when 100% positive examples were given. For AQ15-generated descriptions, errors in experiments 2-4 were primarily due to errors in learning concept C1. The results of the constructive rule learning program AQ17-HCI show that

all the concepts were precisely learned when the program was given 15% positive and 10% negative examples. In all experiments, AQ17-HCI generated the best performing descriptions.

One interesting finding is that increasing the number of training examples in experiments 1 to 5 resulted in only a slight improvement in the performance of the CFS-generated descriptions (from 21.3% to 16.3%). Other

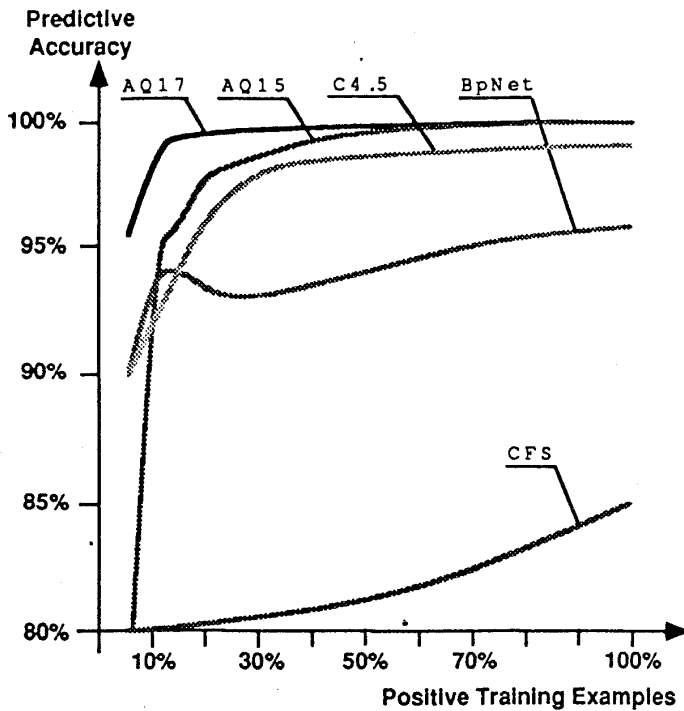


Figure 8. Learning Curves for Concepts in the ROBOTS Domain for the Fixed Number of 10% of Negative Examples

interesting findings are that even with 100% positive examples, the neural net, the genetic algorithm, and to a smaller degree the decision tree method, did not learn the concepts precisely. The CFS classifier system does not seem to be well-suited for DNF-type problems. To further test this finding, this paper reports results involving other genetic algorithm based systems.

Table 2 gives the average R-complexity of the descriptions learned from different training sets. This measure gives a clear division between symbolic and subsymbolic methods. The symbolic methods generate approximately ten times simpler descriptions than the subsymbolic methods. The results in this table are correlated with the results in Table 1. Methods that perform better in terms of predictive accuracy consistently yield simpler concept descriptions.

As mentioned earlier, target concepts were generated by human subjects, and therefore the study favored methods that use symbolic representations, as such representations are more closely related to human representations. Studying how systems learn such human-generated concepts is important for applications where knowledge that needs to be acquired is in such forms, and/or applications where the knowledge learned needs to be understood by human experts. There are problem domains in which these factors are not relevant. Next, two studies present a wider range of problems involving both DNF-type and non-DNF-type concepts.

Second Study: The MONK's Problems

This study reports results from a performance comparison of different learning algorithms on three problems defined in the ROBOTS domain [9]. The so called MONK's problems address three machine learning problems. Problem 1 is a DNF-type problem. Next is a "m-of-n," non-DNF-type problem. The concept to be learned requires a very complex DNF expression to describe it in terms of the available attributes. Problem 3 is DNF-type, but the learning data set contains noise.

Problem M1:

Head shape is the same as the body shape, or color of the jacket is red.

Training set contains 124 randomly selected examples. No noise.

Problem M2:

Exactly two of the six given attributes take their first value.

For example, if attributes Head Shape and Body Shape take value *round*, which is the first value in their value set, then no other attribute may take the first value in its value set. Training set contains 189 randomly selected examples. No noise.

Problem M3:

Jacket is green and holding a sword, or jacket is not blue and body is not octagonal.

Training set contains 122 randomly selected examples. 5% noise in the data.

The tested algorithms fall into 4 categories:

- Neural Networks (Backpropagation, Cascade Correlation)
- Decision Trees (ID3, Assistant Professional, ID5R, IDL, ID5R-hat, TDIDT, PRISM)

- Decision Rules (AQ17-DCI, -HCI, -FCLS, AQ14-NT, AQ15-GA, AQR, CN2)
- Inductive Logic Programming-ILP (mFOIL).

Table 3 shows all reported results [9]. No one classifier based on genetic algorithms was tested as a separate program. In the AQ15-GA program, genetic algorithms are used in conjunction with AQ15. Genetic algorithms are used to explore the space of all subsets of a given attribute set, and AQ15 is used to build concept descriptions. This multistrategy approach improves performance accuracy of the symbolic learning system while learning the M3 problem.

PARADIGM	PROGRAM	PREDICTION ACCURACY		
		DNF-type (no noise)	Non-DNF (m-of-n)	DNF-type (noise)
Neural Nets	* (1) Backpropagation	100%	100%	93%
	(1) Cascade Correlation	100%	100%	97%
Decision Trees	(2) Assistant Professional	100%	81%	100%
	* (3) ID3	99%	68%	94%
	(3) ID3 (no windowing)	83%	69%	96%
	(3, 4) ID5R	82%	69%	95%
	(4) ID5R-hat	90%	66%	—
	(4) IDL	97%	66%	—
	(4) TDIDT	76%	67%	—
(5) PRISM	86%	73%	90%	
Decision Rules	(6) AQ14-NT	100%	77%	100%
	(3) AQR	96%	80%	87%
	(3) CN2	100%	69%	89%
	(6) AQ15	100%	77%	84%
	(6) AQ15-GA	100%	87%	100%
	¶ (6) AQ17-DCI	100%	100%	97%
	* ¶ (6) AQ17-FCLS	100%	93%	97%
	¶ (6) AQ17-HCI	100%	93%	100%
	(6) AQ17	100%	100%	100%
Inductive Logic P. (2)	mFOIL	100%	69%	100%

Table 3. Summary of Results for MONK's Problems

* Programs compared in the first study. ¶ Constructive induction programs. Experiments were performed at the following laboratories: 1) School of Computer Science, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA; 2) AI Laboratory, Josef Stefan Institute, Ljubljana, Slovenia; 3) Institute for Real-Time Computer Control Systems and Robotics, and University of Karlsruhe, Karlsruhe, Germany; 4) Artificial Intelligence Laboratory, Vrije Universiteit Brussel, Brussels, Belgium; 5) AI-Lab, Institute for Informatics, University of Zurich, Switzerland; 6) Center for Artificial Intelligence, George Mason University, Fairfax, Virginia, USA.

Problem M1 is of similar complexity to C1-C5 ROBOTS problems, and it is easily learned by decision rule algorithms, AQ15 and AQ17-HCI. ID3 cannot learn concept M1 precisely; however, M1 was learned by another decision tree algorithm: Assistant Professional neural nets and decision trees paradigms, one can find programs that correctly learned descriptions (Cascade Correlation, Assistant Professional).

Problems M2 and M3 are difficult for selective decision rule and decision tree algorithms AQ15, and ID3. The learned descriptions have either high R-complexity (problem M2) or contain rules that cover noisy examples (problem M3). These problems were not learned as well by a hybrid of decision rules and decision trees, i.e., decision lists (CN2 algorithm, [29]). This suggests that techniques other than those implemented in these programs are required to solve this kind of problem.

The hypothesis-driven constructive induction method

implemented in AQ17-HCI changes the representation space by narrowing and/or expanding the initial set of attributes. The method analyzes inductive hypotheses generated by a selective program and removes and/or generates new attributes. The new attributes are patterns found either in conditions or in the rules. This is sufficient to solve the M3 problem. Problem M2, however, still remains hard. A solution lies in another type of change in the representation, i.e., attribute generation based on combining initial attributes using logical and/or algebraic operators [30]. An initial integration of constructive induction methods was done in the AQ17 program. The multistrategy constructive induction program AQ17 learned all three MONK's problems.

Third Study: The nDmC Learning Problems

This study uses another artificial domain to test twelve DNF concepts and is based on the experiments conducted

Paradigm (Program)	PREDICTION ACCURACY			
	1DmC	2DmC	3DmC	4DmC
Genetic Alg. (GABIL)	96%	93%	90%	89%
(Adptv GABIL) ¶	97%	96%	95%	94%
Decision Trees (C4.5) *	98%	95%	89%	84%
Decision Rules (AQ14) *	99%	97%	96%	95%

Table 4. Prediction Accuracy in the Four DNF Categories
* Programs compared in the first study. ¶ Multistrategy learning program.

Paradigm (Program)	CONVERGENCE (No examples needed to achieve 95% accuracy)			
	1DmC	2DmC	3DmC	4DmC
Genetic Alg. (GABIL)	94	169	151	167
(Adptv GABIL) ¶	63	83	84	88
Decision Trees (C4.5) *	96	135	209	206
Decision Rules (AQ14) *	33	52	102	105

Table 5. Convergence to 95% in the Four DNF Categories
* Programs compared in the first study. ¶ multistrategy learning program.

by Spears and Gordon [10]. The experiments involved learning concepts in a designed domain defined by 4 nominal attributes, each having 4 distinct values. The representation space consisted of 256 examples (vectors of attribute values). There were twelve DNF-type target concepts, differing from each other in the number of rules (the number of disjunctions), and in the number of conditions in disjunctions (the conjunctions in the rules). All twelve concepts can be characterized by the formula $nDmC$, in which n , the number of disjunctions, varied from 1 to 4, and m , the number of conjunctions, varied from 1 to 3. Spears and Gordon first compare three learning methods. Two symbolic methods represented by C4.5—a decision tree learning program and AQ14—a decision rule learning program. Subsymbolic methods were represented by GABIL—Genetic Algorithms Batch Incremental Learner. They conclude that AQ14 is the best performer and use some of AQ's strategies to improve GABIL. The resulting multistrategy system, Adaptive GABIL, is finally evaluated using the same problems.

Tables 4 and 5 show the results from testing the systems according to the prediction accuracy and the convergence criteria. The prediction accuracy is an average over all values on a learning curve. The convergence criterion is the number of events seen before a 95% prediction accuracy is maintained [31]. The results in the tables were averaged for each DNF category over three cases ($m = 1..3$).

Problems labeled 1DmC and 2DmC are similar to problems C1-C5 defined in ROBOTS domain, as far as the R-complexity of descriptions is concerned. In learning such problems, the symbolic learning program AQ14 outperformed the other three programs, both in terms of predictive accuracy and convergence to 95%. For the remaining problems, 3DmC and 4DmC, AQ14 maintains the best prediction accuracy. However, the Adaptive GABIL algorithm that combines symbolic and subsymbolic strategies outperformed the decision tree learning algorithm.

Summary and Future Work

From the multistrategy learning point of view, it is important that capabilities and limitations of different learning strategies and paradigms are well understood. The goal of this study was to make experiments that would help to develop insights into the performance of diverse learning approaches on selected classes of learning problems.

One finding is that symbolic methods outperformed subsymbolic methods in learning DNF-type problems. We

found that the performance accuracy of symbolic methods was high, the convergence to the target concept fast, the learned descriptions have matched or closely matched the target concepts, and were easy to understand. In addition, preliminary results show that the symbolic methods performed very well with DNF-type problems with noise, which contradicts a sometimes expressed belief that neural nets are particularly good for such problems, and symbolic methods are not. The most surprising result, however, was that symbolic methods employing constructive induction performed on par with neural nets on learning non-DNF-type concepts, such as "m-of-n." For such problems, neural nets were supposed to be superior, as the problems are easily representable by such nets. Multistrategy methods, such as those implemented in AQ17 family and Adaptive GABIL, although at an early stage of development, have already shown an improved performance over monostrategy methods.

The performance of the programs was analyzed using a *diagrammatic visualization* system, DIAV. This system, working on line with a learning program, turned out to be a very useful tool for visualizing learned and target concepts, comparing the learned concepts, and presenting errors in learning (an "error image"). The method was also exceptionally useful for visualizing concepts learned by subsymbolic methods, and comparing them with concepts learned by symbolic methods, such as neural net learning and genetic algorithm learning. Concept images helped comprehending knowledge encoded in a neural network or in a population of classifiers. In addition, the visualization method enabled us to determine the R-complexity of the concepts learned by the subsymbolic methods.

Among important topics for the future is the application of the methods to a wider range of non DNF-type problems, such as learning a text-to-speech mapping [6, 32], and to randomly generated problems, in order to evaluate an overall performance of the methods. Future research might also compare the performance of the methods in learning from noisy data, from inconsistent examples, and in learning imprecisely defined or flexible concepts [7].

Acknowledgments

The authors thank Ken Kaufman, Gheorghe Tecuci and Jianping Zhang for comments on earlier drafts of this paper. They are grateful to Jayshree Sarma and Ashraf A. Wahab for performing experiments with backpropagation algorithm BpNet and genetic algorithm learning program CFS, Ross Quinlan for providing a copy of C4.5 system, and to Rick Riolo for providing a copy of his CFS-C software.

This research was done in the Center for Artificial Intelligence at George Mason University. The Center's research is supported in part by the National Science Foundation under the grant No. IRI-9020226, in part by the Defense Advanced Research Projects Agency under the grant administered by the Office of Naval Research No. N00014-91-J-1854, in part by the Office of Naval Research under the grant No. N00014-91-J-1351, and in part by the grant No. F49620-92-J-0549, administered by the Air Force Office of Scientific Research.

References

1. Michalski, R.S. 1983. A Theory and Methodology of Inductive Learning. In *Machine Learning: An Artificial Intelligence Approach*, ed. R.S. Michalski, J.G. Carbonell and T.M. Mitchell. Los Altos, CA: Tioga/Morgan Kaufmann.
2. Wnek, J. and Michalski, R.S. 1991. Hypothesis-Driven Constructive Induction in AQ17 - A Method and Experiments. In *Proceedings of IJCAI-91 Workshop on Evaluating and Changing Representation in Machine Learning*. Sydney, Australia.
3. Fisher, D.H. and McKusick, K.B. 1989. An Empirical Comparison of ID3 and Backpropagation. In *Proceedings of IJCAI-89*. Detroit, MI: Morgan Kaufmann.
4. Mooney, R.J., Shavlik, J., Towell, G., and Gove, A. 1989. An Experimental Comparison of Symbolic and Connectionist Learning Algorithms. In *Proceedings of IJCAI-89*. Detroit, MI: Morgan Kaufmann.
5. Weiss, S.M. and Kapouleas, I. 1989. An Empirical Comparison of Pattern Recognition, Neural Nets, and Machine Learning Classification Methods. In *Proceedings of IJCAI-89*. Detroit, MI: Morgan Kaufmann.
6. Dietterich, T.G., Hild, H. and Bakiri, G. 1990. A Comparative Study of ID3 and Backpropagation for English Text-to-Speech Mapping. In *Proceedings of the 7th International Conference on Machine Learning*. Austin, TX: Morgan Kaufmann.
7. Bergadano, F., Matwin, S., Michalski, R.S., and Zhang, J. 1992. Learning Two-tiered Descriptions of Flexible Concepts: The POSEIDON System. *Machine Learning* 8, 5-43.
8. Wnek, J. and Michalski, R.S. 1992. A Diagrammatic Visualization of Learning Processes. To appear in *Reports of Machine Learning and Inference Laboratory*, George Mason University.
9. Thrun, S.B., et al. 1991. The MONK's Problems: A Performance Comparison of Different Learning Algorithms. *Computer Science Reports, Carnegie Mellon University*, CMU-CS-91-197, Carnegie Mellon University.
10. Spears, W.M. and Gordon, D.F. 1991. Adaptive Strategy Selection for Concept Learning. In *Proceedings of the First International Workshop on Multistrategy Learning*. Harpers Ferry, WV: Center for AI, George Mason University.
11. Quinlan, J.R. 1989. Documentation and User's Guide for C4.5. *unpublished*.
12. Quinlan, J.R. 1986. Induction of Decision Trees. *Machine Learning* 1, no.1, pp. 81-106.
13. Michalski, R.S. 1973. AQVAL/1 - Computer Implementation of a Variable-Valued Logic System VLI and Examples of its Application to Pattern Recognition. In *Proceedings of the First International Joint Conference on Pattern Recognition*. Washington, D.C.
14. Michalski, R.S., Mozetic, I., Hong, J. and Lavrac, N. 1986. The Multi-Purpose Incremental Learning System AQ15 and Its Testing Application to Three Medical Domains. In *Proceedings of AAAI-86*. Philadelphia, PA: Morgan Kaufmann.
15. Michalski, R.S. 1978. Pattern Recognition as Knowledge-Guided Computer Induction. *Technical Report, 927*, Computer Science Dept., University of Illinois, Urbana.
16. Pagallo, G. and Haussler, D. 1990. Boolean Feature Discovery in Empirical Learning. *Machine Learning* 5, no.1, pp. 71-99.
17. Rendell, L.A. and Seshu, R. 1990. Learning Hard Concepts Through Constructive Induction: Framework and Rationale. *Computational Intelligence* 6, pp. 247-270.
18. Michalski, R.S. 1993. Inferential Learning Theory: Developing Foundations for Multistrategy Learning. In *Machine Learning: A Multistrategy Approach, Vol IV*. ed. R.S. Michalski and G. Tecuci. San Mateo, CA: Morgan Kaufmann.
19. Wnek, J. 1992. Version Space Transformation with Constructive Induction: The VS* Algorithm. *Reports of Machine Learning and Inference Laboratory*, MLI-92-01, Center for AI, George Mason University.
20. Rumelhart, D.E., Hinton, G.E., and Williams, R.J.

1986. Learning Internal Representations by Error Propagation. In *Parallel Distributed Processing*, ed. D.E. Rumelhart and J.L. McClelland. Cambridge, MA: The MIT Press.
21. McClelland, J. and Rumelhart, D. 1988. *Explorations in Parallel Distributed Processing*. Cambridge, MA: MIT Press.
22. DeJong, K.A. 1990. The Genetic Algorithm Approach to Machine Learning. In *Machine Learning: An Artificial Intelligence Approach*, ed. Y. Kodratoff and R.S. Michalski. Palo Alto, CA: Morgan Kaufmann.
23. Holland, J.H. and Reitman, J.S. 1978. Cognitive Systems Based on Adaptive Algorithms. In *Pattern-directed Inference Systems*, ed. D.A. Waterman and F. Hayes-Roth. New York: Academic Press.
24. Holland, J.H. 1986. Escaping Brittleness: The Possibilities of General Purpose Learning Algorithms Applied to Parallel Rule-Based Systems. In *Machine Learning: An Artificial Intelligence Approach*, ed. R.S. Michalski, J.G. Carbonell and T.M. Mitchell. Los Altos, CA: Morgan Kaufmann.
25. Riolo, R.L. 1988. CFS-C: A Package of Domain Independent Subroutines for Implementing Classifier Systems in Arbitrary, User-Defined Environments. *Technical Report*, Logic of Computers Group, Division of Computer Science and Engineering, University of Michigan.
26. Kaufman, K.A., Michalski, R.S. and Schultz, A.C. 1989. EMERALD 1: An Integrated System of Machine Learning and Discovery Programs for Education and Research. User's guide MLI-89-12, AI Center, George Mason University.
27. Hempel, C.G. 1965. *Aspects of Scientific Explanation*. New York: The Free Press.
28. Kodratoff, Y. 1993. Induction and the Organization of Knowledge. In *Machine Learning: A Multistrategy Approach, Vol. IV*, ed. R.S. Michalski and G. Tecuci. San Mateo, CA: Morgan Kaufmann.
29. Clark, P. and Niblett, T. 1989. The CN2 Induction Algorithm. *Machine Learning* 3, pp. 261-284.
30. Bloedorn, E. and Michalski, R.S. 1992. Data-driven Constructive Induction in AQ17-DCI: A Method and Experiments. *Reports of Machine Learning and Inference Laboratory*, MLI-92-03, Center for AI, George Mason University.
31. Valiant, L.G. 1984. A Theory of the Learnable. *Communications of the ACM* 27, pp. 1134-1142.
32. Sejnowski, T.J. and Rosenberg, C.R. 1987. Parallel Networks that Learn to Pronounce English Text. *Complex Systems* 1, pp. 145-168.