# Machine Learning for Information Retrieval: Neural Networks, Symbolic Learning, and Genetic Algorithms

#### **Hsinchun Chen**

University of Arizona, Management Information Systems Department, Karl Eller Graduate School of Management, McClelland Hall 430Z, Tucson, AZ 85721. E-mail: hchen@bpa.arizona.edu

Information retrieval using probabilistic techniques has attracted significant attention on the part of researchers in information and computer science over the past few decades. In the 1980s, knowledge-based techniques also made an impressive contribution to "intelligent" information retrieval and indexing. More recently, information science researchers have turned to other newer artificial-intelligence-based inductive learning techniques including neural networks, symbolic learning, and genetic algorithms. These newer techniques, which are grounded on diverse paradigms, have provided great opportunities for researchers to enhance the information processing and retrieval capabilities of current information storage and retrieval systems. In this article, we first provide an overview of these newer techniques and their use in information science research. To familiarize readers with these techniques, we present three popular methods: the connectionist Hopfield network; the symbolic ID3/ID5R; and evolution-based genetic algorithms. We discuss their knowledge representations and algorithms in the context of information retrieval. Sample implementation and testing results from our own research are also provided for each technique. We believe these techniques are promising in their ability to analyze user queries, identify users' information needs, and suggest alternatives for search. With proper user-system interactions, these methods can greatly complement the prevailing full-text, keywordbased, probabilistic, and knowledge-based techniques.

#### Introduction

In the past few decades, the availability of cheap and effective storage devices and information systems has prompted the rapid growth and proliferation of relational, graphical, and textual databases. Information collection and storage efforts have become easier, but effort required to retrieve relevant information has become significantly greater, especially in large-scale databases.

Received September 29, 1993; revised March 25, 1994; accepted June 1, 1994.

© 1995 John Wiley & Sons, Inc.

This situation is particularly evident for textual databases, which are widely used in traditional library science environments, in business applications (e.g., manuals, newsletters, and electronic data interchanges), and in scientific applications (e.g., electronic community systems and scientific databases). Information stored in these databases often has become voluminous, fragmented, and unstructured after years of intensive use. Only users with extensive subject area knowledge, system knowledge, and classification scheme knowledge (Chen & Dhar, 1990) are able to maneuver and explore in these textual databases.

Most commercial information retrieval systems still rely on conventional inverted index and Boolean querying techniques. Even full-text retrieval has produced less than satisfactory results (Blair & Maron, 1985). Probabilistic retrieval techniques have been used to improve the retrieval performance of information retrieval systems (Bookstein & Swanson, 1975; Maron & Kuhns, 1960). The approach is based on two main parameters, the probability of relevance and the probability of irrelevance of a document. Despite various extensions, probabilistic methodology still requires the *independence assumption* for terms and it suffers from difficulty of estimating term-occurrence parameters correctly (Gordon, 1988; Salton, 1989).

Since the late 1980s, knowledge-based techniques have been used extensively by information science researchers. These techniques have attempted to capture searchers' and information specialists' domain knowledge and classification scheme knowledge, effective search strategies, and query refinement heuristics in document retrieval systems design (Chen & Dhar, 1991). Despite their usefulness, systems of this type are considered performance systems (Simon, 1991)—they only perform what they were programmed to do (i.e., they are without learning ability). Significant efforts are often required to acquire knowledge from domain experts and to maintain and update the knowledge base.

A newer paradigm, generally considered to be the machine learning approach, has attracted attention of researchers in artificial intelligence, computer science, and other functional disciplines such as engineering, medicine, and business (Carbonell, Michalski, & Mitchell, 1983; Michalski, 1983; Weiss & Kulikowski, 1991). In contrast to performance systems, which acquire knowledge from human experts, machine learning systems acquire knowledge automatically from examples, that is, from source data. The most frequently used techniques include symbolic, inductive learning algorithms such as ID3 (Quinlan, 1979), multiple-layered, feed-forward neural networks such as backpropagation networks (Rumelhart, Widrow, & Lehr, 1986), and evolutionbased genetic algorithms (Goldberg, 1989). Many information science researchers have started to experiment with these techniques as well (Belew, 1989; Chen & Lynch, 1992; Chen et al., 1993; Gordon, 1989; Kwok, 1989).

In this article, we aim to review the prevailing machine learning techniques and to present several sample implementations in information retrieval to illustrate the associated knowledge representations and algorithms. Our objectives are to bring these newer techniques to the attention of information science researchers by way of a comprehensive overview and discussion of algorithms. We review the probabilistic and knowledge-based techniques and the emerging machine learning methods developed in artificial intelligence (AI). We then summarize some recent work adopting AI techniques in information retrieval (IR). After the overview, we present in detail a neural network implementation (Hopfield network), a symbolic learning implementation (ID3 and ID5R), and a genetic algorithms implementation. Detailed algorithms, selected IR examples, and preliminary testing results are also provided. A summary concludes the study.

# Information Retrieval Using Probabilistic, Knowledge-Based, and Machine Learning Techniques

In classical information retrieval models, relevance feedback, document space modification, probabilistic techniques, and Bayesian inference networks are among the techniques most relevant to our research. In this section, we first summarize important findings in these areas and then present some results from knowledge-based systems research in information retrieval. However, our main purpose will be to present research in machine learning for information retrieval. Similarities and differences among techniques will be discussed.

#### Relevance Feedback and Probabilistic Models in IR

One of the most important and difficult operations in information retrieval is to generate queries that can suc-

cinctly identify relevant documents and reject irrelevant documents. Since it is often difficult to accomplish a successful search at the initial try, it is customary to conduct searches iteratively and reformulate query statements based on evaluation of the previously retrieved documents. One method for automatically generating improved query formulations is the well-known relevancefeedback process (Ide, 1971; Ide & Salton, 1971; Rocchio, 1971; Salton, 1989). A query can be improved iteratively by taking an available query vector (of terms) and adding terms from the relevant documents, while subtracting terms from the irrelevant documents. A single iteration of relevance feedback usually produces improvements of from 40% to 60% in search precision (Salton, 1989). A similar approach can also be used to alter the document representation. Document-vector modification changes and improves document indexes based on the user relevance feedback of relevant and irrelevant documents (Brauen, 1971). Using such a technique, the vectors of documents previously retrieved in response to a given query are modified by moving relevant documents closer to the query and at the same time moving irrelevant documents away from the query. While the relevance feedback procedure is efficient and intuitively appealing, it does not attempt to analyze characteristics associated with the relevant and irrelevant documents to "infer" what concepts (terms) are most appropriate for representing a given query (or queries).

In probabilistic information retrieval, the goal is to estimate the probability of relevance of a given document to a user with respect to a given query. Probabilistic assumptions about the distribution of elements in the representations within relevant and irrelevant documents are required. Using relevance feedback from a few documents, the model can be applied to estimate the probability of relevance for the remaining documents in a collection (Fuhr & Buckley, 1991; Fuhr & Pfeifer, 1994; Gordon, 1988). To simplify computation, an assumption is usually made that terms are distributed independently (Maron & Kuhns, 1960). Fuhr and his coworkers discussed probabilistic models as an application of machine learning. They presented three different probabilistic learning strategies for information retrieval. First, the classical binary independence retrieval model (Robertson & Sparck Jones, 1976; Yu & Salton, 1976) implemented a query-oriented strategy. In the relevance feedback phase, given a query, relevance information was provided for a set of documents. In the application phase, this model can be applied to all documents in the collection, but only for the same initial query. The second document-oriented strategy collected relevance feedback data for a specific document from a set of queries (Maron & Kuhns, 1960). The parameters derived from these data can be used only for the same document, but for all queries submitted to the system. Neither of these strategies can be generalized to all documents and for all

queries. Fuhr et al. proposed a third, feature-oriented strategy. In query-oriented and document-oriented strategies, the concept of abstraction was adopted implicitly by regarding terms associated with the query or the document, instead of the query or document. In this featureoriented strategy, abstraction was accomplished by using features of terms (e.g., the number of query terms, length of the document text, the with-document frequency of a term, etc.) instead of terms themselves. The featureoriented strategy provides a more general form of probabilistic learning and produces bigger learning samples for estimation; but the disadvantage is the heuristics required to define appropriate features for analysis. After transforming terms into features, Fuhr et al. (1990) adopted more sophisticated general-purpose statistical and machine learning algorithms such as regression methods and the decision-tree building ID3 algorithm (Quinlan, 1986) for indexing and retrieval. In summary, by using features of terms instead of terms, Fuhr et al. were able to derive larger learning samples during relevance feedback. The general-purpose analytical techniques of regression methods and ID3 they adopted are similar to the techniques to be discussed in this article.

The use of Bayesian classification and inference networks for information retrieval and indexing represents an extension of the probabilistic models (Maron & Kuhns, 1960; Turtle & Croft, 1990). The basic inference network consists of a document network and a query network (Turtle & Croft, 1990, 1991; Tzeras & Hartmann, 1993) that is intended to capture all of the significant probabilistic dependencies among the variables represented by nodes in the document and query networks. Given the prior probabilities associated with the documents and the conditional probabilities associated with the interior nodes, the posterior probability associated with each node in the network can be computed using Bayesian statistics. The feedback process in a Bayesian inference network is similar to conventional relevance feedback and the estimation problems are essentially equivalent to those observed in probabilistic models. Tzeras and Hartmann (1993) showed that the network can be applied for automatic indexing in large subject fields with encouraging results, although it does not perform better than the probabilistic indexing technique described in Fuhr et al. (1990). Turtle and Croft (1991) showed that, given equivalent document representations and query forms, the inference network model performed better than conventional probabilistic models.

Although relevance feedback and probabilistic models exhibit interesting query or document refinement capabilities, their abstraction processes are based on either simple addition/removal of terms or probabilistic assumptions and principles. Their learning behaviors are very different from those developed in symbolic machine learning, neural networks, and genetic algorithms. In the following two subsections, we will first re-

view knowledge-based information retrieval, and then provide an extensive discussion of the recent machine learning paradigms for information retrieval.

## Knowledge-Based Systems in IR

Creating computer systems with knowledge or "intelligence" has long been the goal of researchers in artificial intelligence. Many interesting knowledge-based systems have been developed in the past few decades for such applications as medical diagnosis, engineering troubleshooting, and business decisionmaking (Hayes-Roth & Jacobstein, 1994). Most of these systems have been developed based on the manual knowledge acquisition process, a significant bottleneck for knowledge-based systems development. A recent approach to knowledge elicitation is referred to as "knowledge mining" or "knowledge discovery" (Frawley, Pietetsky-Shapiro, & Matheus, 1991; Pietetsky-Shapiro, 1989). Grounded on various AI-based machine learning techniques, the approach is automatic and it acquires knowledge or identifies patterns directly from examples or databases. We review some important work in knowledge-based systems in IR and learning systems in IR, respectively, in the next two subsections.

There have been many attempts to capture information specialists' domain knowledge, search strategies, and query refinement heuristics in document retrieval systems design. Some of such systems are "computer-delegated," in that decisionmaking has been delegated to the system and some are "computer-assisted," wherein users and the computer form a partnership (Buckland & Florian, 1991). Because computer-assisted systems have been shown to be more adaptable and useful for search tasks than computer-delegated systems, many knowledge-based systems of this type have been developed for IR over the past decade.

CoalSORT (Monarch & Carbonell, 1987), a knowledge-based system, facilitates the use of bibliographic databases in coal technology. A semantic network, representing an expert's domain knowledge, embodies the system's intelligence. PLEXUS, developed by Vickery and Brooks (1987), is an expert system that helps users find information about gardening. Natural language queries are accepted. The system has a knowledge base of search strategies and term classifications similar to a thesaurus. EP-X (Smith et al., 1989) is a prototype knowledgebased system that assists in searching environmental pollution literature. This system makes extensive use of domain knowledge, represented as hierarchically defined semantic primitives and frames. The system interacts with users to suggest broadening or narrowing operations. GRANT, developed by Cohen and Kjeldsen (1987), is an expert system for finding sources of funding for given research proposals. Its search method—constrained spreading activation in a semantic networkmakes inferences about the goals of the user and thus finds information that the user has not explicitly requested but that is likely to be useful. Fox's CODER system (Fox, 1987) consists of a thesaurus that was generated from the Handbook of Artificial Intelligence and Collin's Dictionary. In CANSEARCH (Pollitt, 1987) a thesaurus is presented as a menu. Users browse and select terms for their queries from the menu. It was designed to enable doctors to search the MEDLINE medical database for cancer literature. The "Intelligent Intermediary for Information Retrieval" ( $I^3R$ ), developed by Croft and Thompson (1987), consists of a group of "experts" that communicate via a common data structure, called the blackboard. The system consists of a user model builder, a query model builder, a thesaurus expert, a search expert (for suggesting statistics-based search strategies), a browser expert, and an explainer. The IOTA system, developed by Chiaramella and Defude (1987), includes natural language processing of queries, deductive capabilities (related to user modeling, search strategies definition, use of expert and domain knowledge), management of full-text documents, and relevance evaluation of answers. Chen and Dhar's (1991) METACAT incorporates several human search strategies and a portion of the Library of Congress Subject Headings (LCSH) for bibliographic search. The system also includes a branch-and-bound algorithm for an automatic thesaurus (LCSH) consultation process.

The National Library of Medicine's thesaurus projects are probably the largest-scale effort that uses the knowledge in existing thesauri. In one of the projects, Rada and Martin (Martin & Rada, 1987; Rada et al., 1989) conducted experiments for the automatic addition of concepts to MeSH (Medical Subject Headings) by including the CMIT (Current Medical Information and Terminology) and SNOMED (Systematized Nomenclature of Medicine) thesauri. Access to various sets of documents can be facilitated by using thesauri and the connections that are made among thesauri. The Unified Medical Language System (UMLS) project is a longterm effort to build an intelligent automated system that understands biomedical terms and their interrelationships and uses this understanding to help users retrieve and organize information from machine-readable sources (Humphreys & Lindbergh, 1989; Lindbergh & Humphreys, 1990; McCray & Hole, 1990). The UMLS includes a Metathesaurus, a Semantic Network, and an Information Sources Map. The Metathesaurus contains information about biomedical concepts and their representation in more than ten different vocabularies and thesauri. The Semantic Network contains information about the types of terms (e.g., "disease," "virus," etc.) in the Metathesaurus and the permissible relationships among these types. The Information Sources Map contains information about the scope, location, vocabulary,

and access conditions of biomedical databases of all kinds.

Another important component of information retrieval is user modeling capability, which is a unique characteristic of reference librarians. During the user-librarian consultation process, the librarian develops an understanding of the type of user being dealt with on the basis of verbal and nonverbal clues. Usually, the educational level of the user, the type of question, the way the question is phrased, the purpose of the search, and the expected search results all play major roles in helping the librarian determine the needs of the user. The librarian, in essence, creates models of the user profile and the task requirements during the consultation process.

User modeling has played a crucial role in applications such as question-answering systems, intelligent tutoring systems, and consultation systems (Appelt, 1985; Chen & Dhar, 1990; Sleeman, 1985; Swarthout, 1985; Zissos & Witten, 1985). An intelligent interface for document retrieval systems must also exhibit the user-modeling capability of experienced human intermediaries. Daniels proposed a frame-based representation for a user model and rules for interacting with the users. She has shown that user modeling is a necessary function in the presearch information interaction (Daniels, 1986). Rich's Grundy system builds models of its users, with the aid of stereotypes, and then uses those models to guide it in its task, suggesting novels that people may find interesting (Rich, 1979a, 1979b, 1983). IR-NLI II (Brajnik, Guida, & Tasso, 1988) incorporates user modeling into a domain-independent bibliographic retrieval expert system. A user model is built based on the user's amount of domain knowledge and search experience.

Despite successes in numerous domains, the development process for knowledge-based systems is often slow and painstaking. Knowledge engineers or system designers need to be able to identify subject and classification knowledge from some sources (usually some domain experts) and to represent the knowledge in computer systems. The inference engines of such systems, which mainly emulate human problem-solving strategies and cognitive processes (Chen & Dhar, 1991), may not be applicable across different applications.

After examining the potential contribution of knowledge-based techniques (natural language processing and expert systems, in particular) to the information retrieval and management tasks, Sparck Jones (1991) warned that it is important not to overestimate the potential of such techniques for IR. She argued that for really hard tasks we will not be able to replace humans by machines in the foreseeable future and many information operations are rather shallow, linguistic tasks, which do not involve elaborate reasoning or complex knowledge. However, she believed AI can contribute to specialized systems and in situations where users and systems complement each other (i.e., computer-assisted systems).

Learning Systems: Neural Networks, Symbolic Learning, and Genetic Algorithms

Unlike the manual knowledge acquisition process and the linguistics-based natural language processing technique used in knowledge-based systems design, learning systems rely on algorithms to extract knowledge or identify patterns in examples or data. Various statistics-based algorithms have been developed by management scientists and have been used extensively over the past few decades for quantitative data analysis. These algorithms examine quantitative data for the purposes of (Parsaye et al., 1989): (1) clustering descriptors with common characteristics, for example, nearest neighbor methods, factor analysis, and principal components analysis; (2) hypothesis testing for differences among different populations, for example, t-test and analysis of variance (ANOVA); (3) trend analysis, for example, time series analysis; and (4) correlation between variables, for example, correlation coefficient, discriminant analysis. and linear/multiple regression analysis (Freund, 1971; Montgomery, 1976). These analysis techniques often rely on complex mathematical models, stringent assumptions, or special underlying distributions. The findings are then presented in mathematical formulas and parameters.

Learning Systems: An Overview. The symbolic machine learning technique, the resurgent neural networks approach, and evolution-based genetic algorithms provide drastically different methods of data analysis and knowledge discovery (Chen et al., in press; Fisher & McKusik, 1989; Kitano, 1990; Mooney et al., 1989; Weiss & Kapouleas, 1989; Weiss & Kulikowski, 1991). These techniques, which are diverse in their origins and behaviors, have shown unique capabilities for analyzing both qualitative, symbolic data and quantitative, numeric data. We provide below a brief overview of these three classes of techniques, along with a representative technique for each class.

• Symbolic learning and ID3: Symbolic machine learning techniques, which can be classified based on such underlying learning strategies as rote learning, learning by being told, learning by analogy, learning from examples, and learning from discovery (Carbonell, Michalski, & Mitchell, 1983), have been studied extensively by AI researchers over the past two decades. Among these techniques, learning from examples, a special case of inductive learning, appears to be the most promising symbolic machine learning technique for knowledge discovery or data analysis. It induces a general concept description that best describes the positive and negative examples. Examples of algorithms which require both positive and negative examples are Quinlan's (1983) ID3 and Mitchell's (1982) Version Space. Some algorithms are batch-oriented, such as

Stepp and Michalski's CLUSTER/RD algorithm (Stepp & Michalski, 1987) and ID3; but some are incremental, such as Utgoff's ID5R (Utgoff, 1989). Many algorithms create a hierarchical arrangement of concepts for describing classes of objects, including Lebowitz' UNIMEM (Lebowitz, 1987), Fisher's COBWEB (Fisher & McKusick, 1989) and Brieman's CART (Brieman et al., 1984). Most of the symbolic learning algorithms produce production rules or concept hierarchies as outputs. These representations are easy to understand and their implementation is typically efficient (especially when compared with neural networks and genetic algorithms).

Among the numerous symbolic learning algorithms which have been developed over the past 15 years, Quinlan's ID3 decision-tree building algorithm and its descendants (Quinlan, 1983, 1986) are popular and powerful algorithms for inductive learning. ID3 takes objects of a known class, described in terms of a fixed collection of properties or attributes, and produces a decision tree incorporating these attributes that correctly classifies all the given objects. It uses an information-economics approach aimed at minimizing the expected number of tests to classify an object. Its output can be summarized in terms of IF–THEN rules.

• Neural networks and backpropagation: The foundation of the neural networks paradigm was laid in the 1950s and this approach has attracted significant attention in the past decade due to the development of more powerful hardware and neural algorithms (Rumelhart, Widrow, & Lehr, 1994). Nearly all connectionist algorithms have a strong learning component. In symbolic machine learning, knowledge is represented in the form of symbolic descriptions of the learned concepts, for example, production rules or concept hierarchies. In connectionist learning, on the other hand, knowledge is learned and remembered by a network of interconnected neurons, weighted synapses, and threshold logic units (Lippmann, 1987; Rumelhart, Hinton, & McClelland, 1986). Learning algorithms can be applied to adjust connection weights so that the network can predict or classify unknown examples correctly. Neural networks have been adopted in various engineering, business, military, and biomedical domains (Chen et al., 1994; Lippmann, 1987; Simpson, 1990; Widrow, Rumelhart, & Lehr, 1994). For example, Hopfield networks have been used extensively in the area of global optimization and search (Hopfield, 1982: Tank & Hopfield, 1987); Kohonen networks have been adopted in unsupervised learning and pattern recognition (Kohonen, 1989). For a good overview of various artificial neural systems, readers are referred to Lippmann (1987).

Among the numerous artificial neural networks that have been proposed recently, backpropagation networks have been extremely popular for their unique learning capability (Widrow et al., 1994). Backpropagation networks (Rumelhart, 1986) are fully connected, layered, feed-forward models. Activations flow from the input layer through the hidden layer, then to the output layer. A backpropagation network typically

starts out with a random set of weights. The network adjusts its weights each time it sees an input-output pair. Each pair is processed at two stages, a forward pass and a backward pass. The forward pass involves presenting a sample input to the network and letting activations flow until they reach the output layer. During the backward pass, the network's actual output is compared with the target output and error estimates are computed for the output units. The weights connected to the output units are adjusted to reduce the errors (a gradient descent method). The error estimates of the output units are then used to derive error estimates for the units in the hidden layer. Finally, errors are propagated back to the connections stemming from the input units. The backpropagation network updates its weights incrementally until the network stabilizes.

• Simulated evolution and genetic algorithms: During the past decade there has been a growing interest in algorithms which rely on analogies to natural processes and Darwinian survival of the fittest. The emergence of massively parallel computers made these algorithms of practical interest. There are currently three main avenues of research in simulated evolution: genetic algorithms; evolution strategies; and evolutionary programming (Fogel, 1994). Each method emphasizes a different facet of natural evolution. Genetic algorithms stress chromosomal operations such as crossover and mutation (Booker, Goldberg, & Holland, 1990; Holland, 1975). Evolution strategies emphasize individual behavioral changes. Evolutionary programming stresses behavioral changes at the level of the species (Fogel, 1962, 1964). Fogel (1994) also provides an excellent review of the history and recent efforts in this area. Among these methods, genetic algorithms have been used successfully for various optimization problems in engineering and biomedical domains.

Genetic algorithms were developed based on the principle of genetics (Goldberg, 1989; Koza, 1992; Michalewicz, 1992). In such algorithms a population of individuals (potential solutions) undergoes a sequence of unary (mutation) and higher order (crossover) transformations. These individuals strive for survival: a selection (reproduction) scheme, biased toward selecting fitter individuals, produces the individuals for the next generation. After some number of generations the program converges—the best individual represents the optimum solution.

Over the past years there have been several studies which compared the performance of these techniques for different applications as well as some systems which used hybrid representations and learning techniques. We summarize some of these studies below.

Mooney et al. (1985) found that ID3 was faster than a backpropagation net, but the backpropagation net was more adaptive to noisy data sets. The performances of these two techniques were comparable, however. Weiss and Kapouleas (1989, 1991) suggested using a resampling technique, such as leave-one-out for evaluation, instead of using a hold-out testing data set. Discriminant

analysis methods, backpropagation net, and decisiontree-based inductive learning methods (ID3-like) were found to achieve comparable performance for several data sets. Fisher and McKusick (1989) found that using batch learning, backpropagation performed as well as ID3, but it was more noise-resistant. They also compared the effect of incremental learning versus batch learning. Kitano (1990) performed systematic, empirical studies on the speed of convergence of backpropagation networks and genetic algorithms. The results indicated that genetic search is, at best, equally efficient as faster variants of a backpropagation algorithm in very small scale networks, but far less efficient in larger networks. Earlier research by Montana and Davis (1989), however, showed that using some domain-specific genetic operators to train the backpropagation network, instead of using the conventional backpropagation delta learning rule, improved performance. Harp, Samad, and Guha (1989) also achieved good results by using GAs for neural network design.

Systems developed by Kitano (1990) and Harp et al. (1989) are also considered hybrid systems (genetic algorithms and neural networks), as are systems like COGIN (Green & Smith, 1991) which performed symbolic induction using genetic algorithms and SC-net (Hall & Romaniuk, 1990), which is a fuzzy connectionist expert system. Other hybrid systems developed in recent years employ symbolic and neural net characteristics. For example, Touretzky and Hinton (1988) and Gallant (1988) proposed connectionist production systems, and Derthick (1988) and Shastri (1991) developed different connectionist semantic networks.

Learning Systems in IR. The adaptive learning techniques cited have also drawn attention from researchers in information science in recent years. In particular, Doszkocs, Reggia, & Lin (1990) provided an excellent review of connectionist models for information retrieval and Lewis (1991) has briefly surveyed previous research on machine learning in information retrieval and discussed promising areas for future research at the intersection of these two fields.

• Neural networks and IR: Neural networks computing, in particular, seems to fit well with conventional retrieval models such as the vector space model (Salton, 1989) and the probabilistic model (Maron & Kuhns, 1960). Doszkocs et al. (1990) provided an excellent overview of the use of connectionist models in information retrieval. These models include several related information processing approaches, such as artificial neural networks, spreading activation models, associative networks, and parallel distributed processing. In contrast to more conventional information processing models, connectionist models are "self-processing" in that no external program operates on the network: the network literally processes itself, with "intelligent be-

havior" emerging from the local interactions that occur concurrently between the numerous network nodes through their synaptic connections. By taking a broader definition of connectionist models, these authors were able to discuss the well-known vector space model, cosine measures of similarity, and automatic clustering and thesaurus in the context of network representation. Based on the network representation, spreading activation methods such as constrained spreading activation adopted in GRANT (Cohen & Kjeldsen, 1987) and the branch-and-bound algorithm adopted in METACAT (Chen & Dhar, 1991) can be considered as variants of connectionist activation. However, only a few systems are considered classical connectionist systems that typically consist of weighted, unlabeled links and exhibit some adaptive learning capabilities.

The work of Belew is probably the earliest connectionist model adopted in IR. In AIR (Belew, 1989), he developed a three-layer neural network of authors, index terms, and documents. The system used relevance feedback from its users to change its representation of authors, index terms, and documents over time. The result was a representation of the consensual meaning of keywords and documents shared by some group of users. One of his major contributions was the use of a modified correlational learning rule. The learning process created many new connections between documents and index terms. Rose and Belew (1991) extended AIR to a hybrid connectionist and symbolic system called SCALIR which used analogical reasoning to find relevant documents for legal research. Kwok (1989) also developed a similar three-layer network of queries, index terms, and documents. A modified Hebbian learning rule was used to reformulate probabilistic information retrieval. Wilkinson and Hingston (1991, 1992) incorporated the vector space model in a neural network for document retrieval. Their network also consisted of three layers: queries, terms, and documents. They have shown that spreading activation through related terms can help improve retrieval performance.

While the above systems represent information retrieval applications in terms of their main components of documents, queries, index terms, authors, etc., other researchers used different neural networks for more specific tasks. Lin, Soergel, & Marchionini (1991) adopted a Kohonen network for information retrieval. Kohonen's feature map, which produced a two-dimensional grid representation for N-dimensional features, was applied to construct a self-organizing (unsupervised learning), visual representation of the semantic relationships between input documents. In MacLeod and Robertson (1991), a neural algorithm was used for document clustering. The algorithm compared favorably with conventional hierarchical clustering algorithms. Chen et al. (1992, 1993, in press) reported a series of experiments and system developments which generated an automatically created weighted network of keywords from large textual databases and integrated it with several existing man-made thesauri (e.g.,

- LCSH). Instead of using a three-layer design, Chen's systems developed a single-layer, interconnected, weighted/labeled network of keywords (concepts) for "concept-based" information retrieval. A blackboardbased design which supported browsing and automatic concept exploration using the Hopfield neural network's parallel relaxation method was adopted to facilitate the usage of several thesauri (Chen et al., 1993). In Chen, Basu, and Ng (in press-a), the performance of a branch-and-bound serial search algorithm was compared with that of the parallel Hopfield network activation in a hybrid neural-semantic network (one neural network and two semantic networks). Both methods achieved similar performance, but the Hopfield activation method appeared to activate concepts from different networks more evenly.
- Symbolic learning and IR: Despite the popularity of using neural networks for information retrieval, we see only limited use of symbolic learning techniques for IR. In Blosseville et al. (1992), the researchers used discriminant analysis and a simple symbolic learning technique for automatic text classification. Their symbolic learning process represented the numeric classification results in terms of IF-THEN rules. Text classification involves the task of classifying documents with respect to a set of two or more predefined classes (Lewis, 1992). A number of systems were built based on human categorization rules (a knowledge-based system approach) (Rau & Jacobs, 1991). However, a range of statistical techniques including probabilistic models, factor analysis, regression, and nearest neighbor methods have been adopted (Blosseville et al., 1992: Lewis, 1992: Masand, Gordon, & Waltz, 1992). Fuhr et al. (1990) adopted regressions methods and ID3 for their feature-based automatic indexing technique. Crawford, Fung, and their coworkers (Crawford et al., 1991; Crawford & Fung, 1992; Fung & Crawford, 1990) have developed a probabilistic induction technique called CONSTRUCTOR and have compared it with the popular CART algorithm (Breiman et al., 1984). Their experiment showed that CON-STRUCTOR's output is more interpretable than that produced by CART, but CART can be applied to more situations (e.g., real-valued training sets). Chen and She (1994) adopted ID3 and the incremental ID5R algorithm for information retrieval. Both algorithms were able to use user-supplied samples of desired documents to construct decision trees of important keywords which could represent the users' queries. For a test collection of about 1000 documents, both symbolic learning algorithms did a good job in identifying the concepts (keywords) which best represent the set of documents identified by users as relevant (positive examples) and irrelevant (negative examples). More testing, however, is underway to determine the effectiveness of example-based document retrieval using ID3 and ID5R.

Several recent works which involved using symbolic learning techniques in the related database areas were also identified, especially in relational database management systems (RDBMS). Cai, Cercone, and

Han (1991) and Han, Cai, and Cercone (1993) developed an attribute-oriented, tree-ascending method for extracting characteristic and classification rules from relational databases. The technique relied on some existing conceptual tree for identifying higher-level, abstract concepts in the attributes. Ioannidis, Saulys, and Whitsitt (1992) examined the idea of incorporating machine learning algorithms (UNIMEM and COB-WEB) into a database system for monitoring the stream of incoming queries and generating hierarchies with the most important concepts expressed in those queries. The goal is for these hierarchies to provide valuable input for dynamically modifying the physical and logical designs of a database. Also related to database design, Borgida and Williamson (1985) proposed the use of machine learning to represent exceptions in databases that are based on semantic data models. Li and McLeod (1989) used machine learning techniques to handle object flavor evolution in object-oriented databases.

· Genetic algorithms and IR: Our literature search revealed several implementations of genetic algorithms in information retrieval. Gordon (1988) presented a genetic algorithms-based approach for document indexing. Competing document descriptions (keywords) are associated with a document and altered over time by using genetic mutation and crossover operators. In his design, a keyword represents a gene (a bit pattern), a document's list of keywords represents individuals (a bit string), and a collection of documents initially judged relevant by a user represents the initial population. Based on a Jaccard's score matching function (fitness measure), the initial population evolved through generations and eventually converged to an optimal (improved) population—a set of keywords which best described the documents. Gordon (1991) further adopted a similar approach to document clustering. His experiment showed that after genetically redescribing the subject description of documents, descriptions of documents found co-relevant to a set of queries will bunch together. Redescription improved the relative density of co-relevant documents by 39.74% after 20 generations and 56.61% after 40 generations. Raghavan and Agarwal (1987) have also studied the genetic algorithms in connection with document clustering. Petry et al. (1993) applied genetic programming to a weighted information retrieval system. In their research, a weighted Boolean query was modified to improve recall and precision. They found that the form of the fitness function has a significant effect upon performance. Yang and coworkers (Yang & Korfhage, 1993; Yang, Korfhage, & Rasmussen, 1993) have developed adaptive retrieval methods based on genetic algorithms and the vector space model using relevance feedback. They reported the effect of adopting genetic algorithms in large databases, the impact of genetic operators, and GA's parallel searching capability. Frieder and Siegelmann (1991) also reported a data placement strategy for parallel information retrieval systems using a genetic algorithms approach. Their results compared favorably with pseudo-optimal document allocations. In Chen and Kim (1993), a GA-NN hybrid system, called GANNET, was developed for IR. The system performed *concept optimization* for user-selected documents using genetic algorithms. It then used the optimized concepts to perform *concept exploration* in a large network of related concepts through the Hopfield net parallel relaxation procedure. A Jaccard's score was also adopted to compute the "fitness" of subject descriptions for information retrieval.

Following this overview, we present three sample implementations of neural networks, symbolic learning, and genetic algorithms, respectively, for illustration purposes. We hope that examining these implementations in the context of IR will encourage other researchers to appreciate these techniques and adopt them in their own research.

#### **Neural Networks for IR**

Neural networks provide a convenient knowledge representation for IR applications in which nodes typically represent IR objects such as keywords, authors, and citations and bidirectional links represent their weighted associations (of relevance). The learning property of backpropagation networks and the parallel search property of the Hopfield network provide effective means for identifying relevant information items in databases. Variants of the backpropagation learning in IR can be found elsewhere (Belew, 1989; Kwok, 1989). In this section, we review a Hopfield network implementation and its associated parallel search property.

# A Hopfield Network: Knowledge Representation and Procedure

The Hopfield net (Hopfield, 1982; Tank & Hopfield, 1987) was introduced as a neural net that can be used as a content-addressable memory. Knowledge and information can be stored in single-layered interconnected neurons (nodes) and weighted synapses (links) and can be retrieved based on the network's parallel relaxation method—nodes are activated in parallel and are traversed until the network reaches a stable state (convergence). It had been used for various classification tasks and global optimization (Lippmann, 1987; Simpson, 1990).

A variant of the Hopfield network for creating a network of related keywords developed by Chen (Chen & Lynch, 1992; Chen et al., 1993) used an asymmetric similarity function to produce thesauri (or knowledge bases) for different domain-specific databases. These automatic thesauri were then integrated with some existing manually created thesauri for assisting concept exploration and query refinement. A variant of the Hopfield parallel relaxation procedure for network search (Chen et al.,

1993) and concept clustering (Chen et al., in press-b) had been reported earlier.

The implementation reported below incorporated the basic Hopfield net iteration and convergence ideas. However, significant modification was also made to accommodate unique characteristics of information retrieval; for example, asymmetric link weights and the continuous SIGMOID transformation function. With the initial search terms provided by searchers and the association of keywords captured by the network, the Hopfield parallel relaxation algorithm activates neighboring terms, combines weighted links, performs a transformation function (a SIGMOID function,  $f_s$ ), and determines the outputs of newly activated nodes. The process repeats until node outputs remain unchanged with further iterations. The node outputs then represent the concepts that are strongly related to the initial search terms. A sketch of the Hopfield net activation algorithm follows:

(1) Assigning synaptic weights: For thesauri which were generated automatically using a similarity function (e.g., the COSINE function) (Everitt, 1980), the resulting links represent probabilistic, synaptic weights between any two concepts. For other external thesauri which contain only symbolic links (e.g., narrower term, synonymous term, broader term, etc.), a user-guided procedure of assigning a probabilistic weight to each symbolic link can be adopted (Chen et al., 1993).

The "training" phase of the Hopfield net is completed when the weights have been computed or assigned.  $t_{ij}$  represents the "synaptic" weight from node i to node j.

(2) Initialization with search terms: An initial set of search terms is provided by searchers, which serves as the input pattern. Each node in the network which matches the search terms is initialized (at time 0) to have a weight of 1.

$$\mu_i(0) = x_i, 0 \le i \le n-1$$

 $\mu_i(t)$  is the output of node *i* at time *t* and  $x_i$ , which has a value between 0 and 1, indicates the input pattern for node *i*.

(3) Activation, weight computation, and iteration:

$$\mu_j(t+1) = f_s \left[ \sum_{i=0}^{n-1} t_{ij} \mu_i(t) \right], 0 \le j \le n-1$$

where  $f_s$  is the continuous SIGMOID transformation function as shown below (Dalton & Deshmane, 1991; Knight, 1990)

$$f_s(net_j) = \frac{1}{1 + \exp\left[\frac{-(net_j - \theta_j)}{\theta_0}\right]}$$

where  $net_i = \sum_{i=0}^{n-1} t_{ii}\mu_i(t)$ ,  $\theta_i$  serves as a threshold or

bias and  $\theta_0$  is used to modify the shape of the SIG-MOID function.

This formula shows the *parallel relaxation* property of the Hopfield net. At each iteration, all nodes are activated at the same time. The weight computation scheme,  $net_j = \sum_{i=0}^{n-1} t_{ij}\mu_i(t)$ , is a unique characteristic of the Hopfield net algorithm. Based on parallel activation, each newly activated node derives its new weight based on the summation of the products of the weights assigned to its neighbors and their synapses.

(4) Convergence: The above process is repeated until there is no change in terms of output between two iterations, which is accomplished by checking:

$$\sum_{j=0}^{n-1} |\mu_j(t+1) - \mu_j(t)| \leq \epsilon$$

where  $\epsilon$  is the maximal allowable error (a small number). The final output represents the set of terms relevant to the starting keywords. Some default threshold values were selected for  $(\theta_{i}, \theta_{0})$ .

# A Hopfield Network Example

A sample session of the Hopfield net spreading activation is presented below. Three thesauri were incorporated in the experiment: a Public thesaurus (generated automatically from 3000 articles extracted from DIA-LOG), the ACM Computing Review Classification System (ACM CRCS), and a portion of the Library of Congress Subject Headings (LCSH) in the computing area. The links in the ACM CRCS and in the LCSH were assigned weights between 0 and 1. Several user subjects (MIS graduate students) were also asked to reviewed selected articles and create their own folders for topics of special interest to them. Notice that some keywords were folder names assigned by the users (in the format of \*.\*); for example, QUERY.OPT folder for query optimization topics; DBMS.AI folder for artificial intelligence and databases topics; and KEVIN.HOT folder for "HOT" (current) topics selected by a user, Kevin. In the example shown below, the searcher was asked to identify descriptors which were relevant to "knowledge indexed deductive search." The initial search terms were: "information retrieval," "knowledge base," "thesaurus," and "automatic indexing" (as shown in the following interaction).

TABLE 1. Sample Hopfield net iterations.

Iteration no.	Suggested terms	Activations
0	INFORMATION RETRIEVAL	1.00
	KNOWLEDGE BASE	1.00
	THESAURUS	1.00
	AUTOMATIC INDEXING	1.00
1	INDEXING	0.65
	KEVIN.HOT	0.56
	CLASSIFICATION	0.50
	EXPERT SYSTEMS	0.50
	ROSS.HOT	0.44
2	RECALL	0.50
3	INFORMATION RETRIEVAL SYSTEM EVALUATION	0.26
4	SELLING – INFORMATION STORAGE AND RETRIEVAL SYSTEMS	0.15
	•••	

Enter the number of system-suggested terms or '0' to quit >> 10

{\* The users supplied a target number of relevant terms.\*}

Given these starting terms, the Hopfield net iterated and converged after 11 iterations. The activated terms after the first four iterations and their associated levels of activation are shown in Table 1. Due to the damping effect of the parallel search property (i.e., the farther away from the initial search terms, the weaker the activation), terms activated at later iterations had lower activation values and were less relevant to the initial search terms in general. Fourteen terms were suggested after the complete Hopfield net activation. Searchers could browse the system-suggested list, select terms of interest, and then activate the Hopfield net again. The user-system interaction continued until the user decided to stop.

- {\* The system reported 14 relevant terms as shown
  below.\*}
- 1. ( ) INDEXING
- 2. ( ) SELLING INFORMATION STORAGE AND RETRIEVAL SYSTEMS
- 3. ( ) KEVIN.HOT
- 4. ( ) INFORMATION RETRIEVAL SYSTEM EVALUATION
- 5. ( ) RECALL
- 6. ( ) EXPERT SYSTEMS
- 7. ( ) CLASSIFICATION
- 8. ( ) DBMS.AI
- 9. ( ) ROSS.HOT
- 10. ( ) INFORMATION STORAGE AND RETRIEVAL SYSTEMS
- 11. ( ) INFORMATION RETRIEVAL
- 12. ( ) KNOWLEDGE BASE
- 13. ( ) THESAURUS
- 14. ( ) AUTOMATIC INDEXING

Enter numbers [1 to 14] or '0' to quit: 1, 2, 4, 5, 7, 10-14

{\* The user selected terms he deemed relevant.
 The system confirmed the selections made and
 display the source for each term. \*}

- 1. (P ) INDEXING
- 2. ( L) SELLING INFORMATION STORAGE AND RETRIEVAL SYSTEMS
- 3. (P ) INFORMATION RETRIEVAL SYSTEM EVALUATION
- 4. (P ) RECALL
- 5. (P ) CLASSIFICATION
- 6. ( L) INFORMATION STORAGE AND RETRIEVAL SYSTEMS
- 7. (P L) INFORMATION RETRIEVAL
- 8. (P ) KNOWLEDGE BASE
- 9. (P ) THESAURUS
- 10. (P L) AUTOMATIC INDEXING

Enter the number of system-suggested terms or '0' to quit >> 30

{\* The uses decide to broaden the search by requesting the Hopfield network to identify 30
new terms based on the terms he had selected. \*}

Enter number [1 to 40] or '0' to quit: 3-7, 9, 33, 35, 36, 38

• • • • • • • •

Enter numbers [1 to 67] or '0' to quite: 0
{\* The system listed his final selections. \*}

- 1. (P ) PRECISION
- 2. (P L) INFORMATION RETRIEVAL
- 3. (P ) INDEXING
- 4. (P L) AUTOMATIC INDEXING
- 5. (P ) RECALL
- 6. ( L) AUTOMATIC ABSTRACTING
- 7. ( L) AUTOMATIC CLASSIFICATION
- 8. ( L) AUTOMATIC INFORMATION RETRIEVAL
- 9. (P ) INFORMATION RETRIEVAL SYSTEM EVALUATION
- 10. (P ) THESAURUS
- 11. ( L) INFORMATION STORAGE AND RETRIEVAL SYSTEMS
- 12. (P ) KNOWLEDGE BASE
- {\* A total of 12 terms were selected. Eight terms
  were suggested by the Hopfield net algorithm. \*}

In a more structured benchmark experiment, we tested 30 sample queries using the Hopfield algorithm in an attempt to understand the general behavior of the algorithm. We tested five cases each for queries with 1 term, 2 terms, 3 terms, 4 terms, 5 terms, and 10 terms, a total of 30 cases. A few examples of the queries used, all in the computing area, were: (1 term: Natural Language Processing); (2 terms: Group Decision Support Systems, Collaboration); (3 terms: Systems Analysis and Design, Simulation and Modeling, Optimization); etc.

TABLE 2. Results of Hopfield network testing.

Case	No. of terms	Query terms in (P, A, L)	Suggested terms in NN: (P, A, L)	No. of iterations NN	Times (seconds) NN
1	1	(1, 1, 1)	(12, 7, 7)	18	21
2	1	(1, 0, 1)	(5, 0, 16)	15	14
3	1	(1, 1, 1)	(11, 5, 11)	14	18
4	1	(0, 0, 1)	(0, 0, 20)	11	10
5	1	(1, 0, 1)	(4, 4, 19)	17	26
6	2	(2, 1, 0)	(19, 2, 3)	21	18
7	2	(2, 0, 2)	(16, 0, 8)	19	22
8	2	(2,0,0)	(20, 3, 4)	20	24
9	2	(2, 1, 1)	(11, 5, 11)	15	16
10	2	(2, 1, 2)	(11, 0, 12)	27	29
11	3	(3,0,1)	(20, 0, 18)	19	31
12	3	(1, 2, 1)	(4, 11, 8)	22	34
13	3	(2, 1, 3)	(22, 1, 8)	18	29
14	3	(1, 3, 1)	(20, 2, 2)	16	23
15	3	(1, 2, 2)	(13, 9, 3)	9	10
16	4	(2, 2, 4)	(17, 4, 4)	17	11
17	4	(3, 2, 2)	(11, 2, 13)	19	31
18	4	(2, 3, 2)	(18, 5, 6)	24	33
19	4	(1, 3, 4)	(18, 2, 5)	19	32
20	4	(1, 2, 1)	(15, 8, 3)	18	6
21	5	(1, 4, 1)	(19, 4, 6)	16	27
22	5	(4, 2, 2)	(10, 1, 12)	15	27
23	5	(3, 2, 4)	(2, 0, 18)	11	23
24	5	(5,0,1)	(19, 0, 3)	23	33
25	5	(5,0,1)	(20, 0, 1)	12	30
26	10	(8,0,3)	(11, 0, 13)	17	34
27	10	(10, 1, 3)	(13, 2, 10)	25	32
28	10	(8, 0, 4)	(16, 0, 8)	24	36
29	10	(9, 1, 5)	(19, 1, 6)	27	25
30	10	(8, 2, 3)	(20, 2, 3)	28	31
Average	5	(3.1, 1.2, 1.9)	(14.5, 2.5, 8.5)	18.8	24.5

For each query, we selected terms from different knowledge sources, "P" for the Public KB, "A" for the ACM CRCS, and "L" for the LCSH, as shown in Table 2. Some terms may have appeared in more than one knowledge source. The three knowledge sources contained about 14,000 terms and 80,000 weighted links. The results shown in Table 2 reveal the number of iterations, the computing times, and the sources of knowledge for the query terms and the system-suggested terms. The reason for investigating the source of knowledge for system-suggested terms was to show the extent to which the Hopfield algorithm branched out and utilized knowledge from various knowledge sources.

Despite the variation in the number of starting terms, the response times increased only slightly when the number of starting terms was increased. The average response time was 24.5 seconds after about an average of about 19 iterations by the Hopfield network. The reason for this was that the Hopfield net thresholds ( $\theta_0$  and  $\theta_j$ ) helped prune the search space. However, more stringent thresholds may need to be adopted to achieve reasonable real-time response for large databases.

Another important observation was that the Hopfield net appeared to invoke the different knowledge sources quite evenly. As shown in Table 2, for most queries the Hopfield net (NN) almost always produced terms from all three knowledge sources. Most terms suggested by the algorithm appeared relevant and many of them were multiple links away from the initial search terms (conventional Hypertext browsing does not traverse multiple links effectively). However, detailed user studies need to be performed to examine the usefulness of the algorithm in search, especially for large-scale applications.

#### Symbolic Learning for IR

Even though symbolic learning techniques have been adopted frequently in various database, engineering, and business domains, we see only limited use of such techniques in IR. For illustration purposes, we summarize below a symbolic learning for IR implementation based on the ID3 and ID5R algorithms (Chen & She, 1994).

ID3/ID5R: Knowledge Representation and Procedure

ID3 is a decision-tree building algorithm developed by Quinlan (1979, 1983). It adopts a divide-and-conquer strategy for object classification. Its goal is to classify mixed objects into their associated classes based the objects' attribute values. In a decision tree, one can classify a node as:

- a leaf node that contains a class name; or
- a non-leaf node (or decision node) that contains an attribute test.

Each training instance or object is represented as a list of attribute-value pairs, which constitutes a conjunctive description of that instance. The instance is labeled with the name of the class to which it belongs. Using the divide-and-conquer strategy, ID3 picks an attribute and uses it to classify a list of objects based on their values associated with this attribute. The subclasses which are created by this division procedure are then further divided by picking other attributes. This process continues until each subclass produced contains only a single type of object. To produce the simplest decision tree (a minimal tree) for classification purpose, ID3 adopts an information-theoretic approach which aims at minimizing the expected number of tests to classify an object. An entropy (a measure of uncertainty) concept is used to help decide which attribute should be selected next. In general, an attribute which can help put objects in their proper classes tends to reduce more entropy and thus should be selected as a test node.

In IR, we can assume that there exists a database (universe) of records (documents, tables, etc.). Records are described by attributes (keywords, primary keys, fields). Each record in the database then belongs to only one of two possible classes:

- the "positive" class (+): consisting of records that are desired; and
- the "negative" class (-): consisting of records that are undesired.

Different database users may desire different sets of documents due to their unique information needs, and the set of documents desired by one user often constitutes only a small portion of the entire database. Enabling the system to identify this small set of positive documents is therefore a challenging task.

In our implementation, we maintained a list of all the keywords that existed in the desired documents and used this list to decide what attributes were crucial to describing documents in the positive class. The test at each non-leaf node of the decision tree determined the presence or absence of a particular keyword: "yes" meant that the test keyword existed in a document, and "no" meant that the keyword did not exist in a document. Thus, ID3

created a binary classification tree. A sketch of the ID3 algorithm adopted follows:

(1) Compute entropy for mixed classes: Initially searchers were requested to provide a set of positive and negative documents. This set of documents served as the training examples for the ID3 algorithm. Entropy was calculated by using the following function (Quinlan, 1983):

$$entropy = -p_{pos}\log p_{pos} - p_{neg}\log p_{neg}$$

where  $p_{pos}$  and  $p_{neg}$  represented the proportions of the documents that were positive or negative, respectively.

- (2) Select the best attribute based on entropy reduction: For each untested attribute (keyword), the algorithm computed an entropy value for its use when classifying mixed documents. Each branch of the decision tree represented the existence or nonexistence of a particular keyword. The keyword that reduced the entropy most served as the next decision node in the tree. As a "greedy" algorithm, ID3 always aims at maximizing local entropy reduction and never backtracks.
- (3) Iterate until all documents are classified: Repeating steps (1) and (2), ID3 computed the entropy value of each mixed class and identified the best attribute for further classifying the class. The process was continued until each class contained either all positive or all negative documents.

Considered as an incremental version of the ID3 algorithm, ID5R, developed by Utgoff (1989), is guaranteed to build the same decision tree as ID3 for a given set of training instances (Quinlan, 1993). In ID5R, a non-leaf node contains an attribute test (same as in ID3) and a set of other non-test attributes, each with object counts for the possible values of the attribute. This additional non-test attribute and object count information at each noleaf node allows ID5R to update a decision tree without rebuilding the entire tree. During the tree rebuilding process, an old test node may be replaced by a new attribute or swapped with other positions in the tree. As in ID3, the tree-building process requires much less computation and time than other inductive learning methods, including neural networks and genetic algorithms.

To create a robust and real-time inductive learning system, a relevance feedback scheme was introduced into our implementation. Although the proposed inductive learning algorithms require users to provide examples to confirm their interests, it is inconceivable that users will be able to browse the entire database to identify such instances. An incremental, interactive feedback process, therefore, was designed to allow users to examine a few documents at a time. In essence, our ID5R algorithm was implemented such that it provided a few suggested documents based on the documents initially provided by

a user after examining a small portion of the database. When a predetermined number of desired documents had been found (say three, in our implementation), the system presented these documents to the user immediately for evaluation (as desired or undesired). This iterative system-induction and user-feedback process continued until the user decided to stop or the complete database had been traversed.

During the relevance feedback process, the newly confirmed documents, either desired or undesired, could be used by ID5R to update the decision tree it previously had constructed. It was shown that when more examples are provided by the users and when the database is more exhaustively searched, ID5R can significantly improve its classification accuracy and search performance.

# An ID3/ID5R Example

We created a small test database of 60 records. For evaluation purposes, we were able to manually select a small set of target desired documents (i.e., eight documents in the areas of information retrieval and keywording). The goal of the experiment was to present a few documents at a time to our system and see whether the system would be able to identify them after the iterative relevance feedback process. The performance of our ID5R-based system was also compared with that of the more conventional ID3 algorithm, which used only an initial set of desired documents to generate a query tree. Sample entries in the literature database are shown below, where the first column represents the document number, and the remaining columns represent different numbers of keywords (two to five) associated with the document.

010 generic, keyword, reference 013 modeling, thesaurus, terrorism 014 modeling, simulation, thesaurus, terrorism 018 keyword, thesaurus 021 ID3, AI, NN 022 file, keyword 023 hierarchy, interface, index 030 carat, AI, expert, keyword, thesaurus 031 AI, protocol, thesaurus 048 keyword, retrieval 049 cross-reference, remote use, redundancy expectations, market, maintenance, quel, interface 107 IT, computerized, MIS 149 database, query, keyword 152 sort, indexing, merge, keyword 177 country, code, keyword, ISO

Initially the user was able to identify the following documents as desired (+) or undesired (-), respectively (documents which the user had seen before):

006 thesaurus, remote use, keyword (+)

008 retrieval, interface (+)

**083** syntax checking, remote use, test, user (-)

**084** interface, protocol, standardization (–)

Providing negative documents was optional. If a user could not think of an example of a document which was undesired, the system by default automatically generated one negative document which contained no keyword identical to any that was present in the desired set. The initial positive keyword list then consisted of all keywords from desired documents; that is, thesaurus, remote use, keyword, retrieval, interface (in that order). Therefore, the set of initial training instances can be represented as:

Initial Training Instances							
у	у	n	n	(+)			
n	n	у	у	(+)			
у	n	n	n	(-)			
n	n	n	У	(-)			
	y n y n	y y n n	y y n n n y	y y n n n n y y			

If a document contained a particular keyword in the keyword list, its attribute value was labeled "y" ("yes"), otherwise the value was "n" ("no"). Based on the set of training instances, ID3 first computed the entropy value when adopting "thesaurus" (the first keyword obtained from the desired documents). It then computed the entropy values when adopting other positive keywords. The "thesaurus" keyword produced the most entropy reduction and was thus selected as the first decision node. Following the same computation, "retrieval" was selected as the next (and last) decision node. ID3 constructed the decision tree shown in Figure 1. In the figure, for example, [2, 1] means 2 instances were in the negative class and 1 instance was in the positive class. The deci-

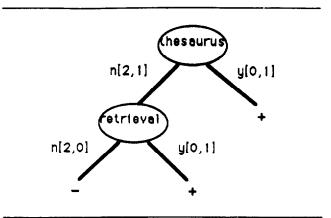


FIG. 1. Initial tree created for an IR example.

sion tree in Figure 1 can be represented as production rules: (1) IF a document has "thesaurus" as a keyword THEN it is desired (one +, the rightmost branch); (2) IF a document does not have "thesaurus" as a keyword, but has "retrieval" THEN it is also a desired document (one +, the middle branch); (3) IF a document does not have "thesaurus" or "retrieval" as a keyword THEN it is an undesired document (two-, the leftmost branch).

Based on this decision tree, the system searched the database for similar documents and identified three more documents as presented below:

013 modeling, thesaurus, terrorism (+)

014 modeling, simulation, thesaurus, terrorism (+)

018 keyword, thesaurus (+)

These documents were then presented to the user, who provided feedback as to whether or not they were desired. If the user confirmed that document 018 was desired but rejected documents 013 and 014, ID5R used the new (contradictory) evidence to update its current tree. The new training instances for ID5R were:

New Training Instances							
у	n	n	n	n	(-)		
у	n	n	n	n	(-)		
y	n	У	n	n	(+)		

The system produced a new tree as shown in Figure 2. This new tree looked different from the original one and can be summarized by the following rules: (1) IF a document has "keyword" as a keyword THEN it is desired (two +, the rightmost branch); (2) IF a document does not have "keyword" as a keyword, but has "retrieval" THEN it is also a desired document (one +, the middle branch); (3) IF a document does not have "keyword" or "retrieval" as a keyword THEN it is an undesired document (four -, the leftmost branch). The whole process was repeated until the entire database was traversed. For this particular example, the final decision tree was the same as the one shown in Figure 2.

To determine how ID5R performed during the user relevance feedback process we examined its *recall* at each point of relevance feedback and compared its performance with that of ID3. ID3 used only the initial document feedback from the users to construct a decision tree and used the tree to search the database. ID5R, on the other hand, collected new evidence during each iteration and updated its trees accordingly. The *recall* measure was defined as:

$$Recall = \frac{\text{Number of relevant records retrieved}}{\text{Total number of relevant records in database}}$$

We developed a test database of about 1000 documents from the 1992 COMPENDEX CD-ROM collection of computing literature. We then identified 10 research topics, each of which had between 5 and 20 relevant documents in the database (manually identified). The testing was conducted by comparing the recall of the ID3 algorithm and that of the ID5R incremental approach using the 10 research topics.

Detailed results of the experiment are presented in Table 3. ID5R and ID3 achieved the same levels of performance for 5 of the 10 test cases (cases 3 and 6-9). After we examined these cases carefully, we found that the initial documents presented for these cases had very precise keywords assigned to them. New instances provided during relevance feedback were consistent with the initial documents, thus ID5R did not revise its decision tree. (At each interaction, ID5R searched only a portion of the entire database. The trees constructed by ID3 remained constant because ID3 did not have any interaction with its users. However, to compare its results with those of the ID5R fairly, ID3's performance at each interaction was computed based on the same documents visited by ID5R. As more documents were examined, ID3's classification results may also have improved.)

For the other five test cases, ID5R's performance increased gradually until it reached 93.1%. ID3 had been able to reach 74.9%. These research topics tended to have more diverse keywords in the initial documents provided. ID5R appeared to benefit from incremental query tree revision based on the relevance feedback information provided by users. In all 10 cases, ID5R was able to terminate in eight interactions. The response times were often less than a second for each decision-tree building process.

In conclusion, the symbolic ID3 algorithm and its ID5R variant both were shown to be promising techniques for inductive document retrieval. By using the en-

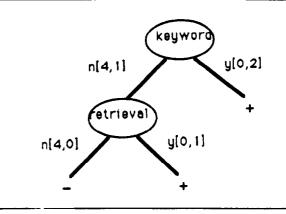


FIG. 2. Updated tree after relevance feedback.

TABLE 3. Results of ID3 and ID5R testing.

Case	Int. 1 ID3/ID5R	Int. 2 ID3/ID5R	Int. 3 ID3/ID5R	Int. 4 ID3/ID5R	Int. 5 ID3/ID5R	Int. 6 ID3/ID5R	Int. 7 ID3/ID5R	Int. 8 ID3/ID5R	Target
1	1/1	1/2	2/3	5/6	6/9				10
2	0/0	0/1	0/2	1/4	1/5	2/7	2/8	3/10	11
3	1/1	2/2	3/3	4/4	·	·	·	•	4
4	1/1	1/1	1/2	1/3	2/4	5/7			10
5	0/0	0/1	0/2	3/5					6
6	1/1	2/2	5/5	6/6					6
7	1/1	2/2	3/3	5/5					5
8	2/2	3/3	3/3	6/6	7/7				8
9	5/5	7/7	8/8	9/0	10/10	11/11			12
10	1/1	2/2	3/3	4/4	7/7	7/10			10
Avg. hits	1.3/1.3	2/2.3	2.8/3.4	4.4/5.2	5.1/6.2	5.6/7.1	5.6/7.2	5.7/7.4	8.2
Avg. recall	16.0/16.0	16.5/31.2	35.0/40.1	55.5/64.1	66.3/79.3	74.0/90.4	74.0/91.3	74.9/93.1	

tropy concept in selecting keywords, both algorithms were able to create minimal and understandable decision trees efficiently. However, ID5R's incremental learning and relevance feedback capabilities made it more robust and appealing for large-scale, real-time IR applications.

## Genetic Algorithms for IR

Often compared with the neural networks and the symbolic learning methods, the self-adaptiveness property of genetic algorithms is also extremely appealing for IR applications.

A Genetic Algorithm: Knowledge Representation and Procedure

Genetic algorithms (GAs) (Goldberg, 1989; Kohonen, 1989; Michalewicz, 1992) are problem-solving systems based on principles of evolution and heredity. A GA maintains a population of individuals,  $P(t) = x_1, \ldots$  $x_n$  at iteration t. Each individual represents a potential solution to the problem at hand and is implemented as some (possibly complex) data structure S. Each solution  $x_i$  is evaluated to give some measure of fitness. Then a new population at iteration t+1 is formed by selecting the fitter individuals. Some members of the new population undergo transformation by means of genetic operators to form new solutions. There are unary transformations  $m_i$  (mutation type), which create new individuals by a small change in a single individual and higher order transformations  $c_i$  (crossover type), which create new individuals by combining parts from several (two or more) individuals. For example, if parents are represented by a five-dimensional vector  $(a_1, a_2, a_3, a_4, a_5)$  and  $(b_1, b_2, b_3, a_4, a_5)$  $b_4$ ,  $b_5$ ), then a crossover of chromosomes after the second gene produces offspring  $(a_1, a_2, b_3, b_4, b_5)$  and  $(b_1, b_2, a_3, b_4, b_5)$  $a_4$ ,  $a_5$ ). The control parameters for genetic operators (probability of crossover and mutation) need to be carefully selected to provide better performance. The intuition behind the crossover operation is information exchange between different potential solutions. After some number of generations the program converges—the best individual hopefully represents the optimum solution. Michalewicz (1992) provided an excellent algorithmic discussion of GAs. Goldberg (1989, 1994) presented a good summary of many recent GA applications in biology, computer science, engineering, operations research, physical sciences, and social sciences.

Genetic algorithms use a vocabulary borrowed from natural genetics in that they talk about *genes* (or bits), chromosomes (individuals or bit strings), and population (of individuals). Populations evolve through generations. Our genetic algorithm was executed in the following steps:

(1) Initialize population and evaluate fitness: To initialize a population, we needed first to decide the number of genes for each individual and the total number of chromosomes (popsize) in the initial population. When adopting GAs in IR, each gene (bit) in the chromosome (bit string) represents a certain keyword or concept. The loci (locations of a certain gene) decide the existence (1, ON) or nonexistence (0, OFF) of a concept. A chromosome therefore represents a document that consists of multiple concepts. The initial population contains a set of documents which were judged relevant by a searcher through relevance feedback. The goal of a GA was to find an optimal set of documents which best matched the searcher's needs (expressed in terms of underlying keywords or concepts). An evaluation function for the fitness of each chromosome was selected based on Jaccard's score matching function as used by Gordon (1988) for document indexing. The Jaccard's score between two sets, X and Y, was computed as:

$$\#(X \cap Y)/\#(X \cup Y)$$

where #(S) indicated the cardinality of set S. The Jaccard's score is a common measure of association in information retrieval (van Rijsbergen, 1979).

(2) Reproduction (selection): Reproduction is the selection of a new population with respect to the probability distribution based on the fitness values. Fitter individuals have better chances of being selected for reproduction (Michalewicz, 1992). A roulette wheel with slots (F) sized according to the total fitness of the population was defined as follows:

$$F = \sum_{i=1}^{popsize} fitness(V_i)$$

where  $fitness(V_i)$  indicated the fitness value of chromosome  $V_i$  according to the Jaccard's score.

Each chromosome had a certain number of slots proportional to its fitness value. The selection process was based on spinning the wheel *popsize* times—each time we selected a single chromosome for a new population. Obviously, some chromosomes were selected more than once. This is in accordance with the genetic inheritance: the best chromosomes get more copies, the average stay even, and the worst die off.

(3) Recombination (crossover and mutation): We were then ready to apply the first recombination operator, crossover, to the individuals in the new population. The probability of crossover,  $p_c$ , gave us the expected number  $p_c \times$  popsize of chromosomes which should undergo the crossover operation. For each chromosome, we generated a random number r between 0 and 1; if  $r < p_c$ , then the chromosome was selected for crossover. We then mated selected pairs of chromosomes randomly: for each pair of coupled chromosomes we generated a random number pos from the range of  $(1 \cdots m - 1)$ , where m was the total number of genes in a chromosome. The number pos indicated the position of the crossing point. The coupled chromosomes exchanged genes at the crossing point as described earlier.

The next recombination operator, mutation, was performed on a bit-by-bit basis. The probability of mutation,  $p_m$ , gave us the expected number of mutated bits  $p_m \times m \times popsize$ . Every bit in all chromosomes of the whole population had an equal chance to undergo mutation, that is, change from 0 to 1 or vice versa. For each chromosome in the crossovered population, and for each bit within the chromosome, we generated a random number r from the range of  $(0 \cdot 1)$ ; if  $r < p_m$ , we mutated the bit. Typical  $p_c$  selected ranged between 0.7 and 0.9 and  $p_m$  ranged between 0.01 and 0.03.

(4) Convergence: Following reproduction, crossover, and mutation, the new population was ready for its next generation. The rest of the evolutions were simply cyclic repetitions of the above steps until the system reached a predetermined number of generations or converged (i.e., showed no improvement in the overall fitness of the population).

## A GA Example

We present a sample session, implementation details, and some benchmark testing results below. In our system, a keyword represented a gene (bit) in GAs; a user-selected document represented a chromosome (individual); and a set of user-selected documents represented the initial population.

The keywords used in the set of user-selected documents were first identified to represent the underlying bit strings for the initial population. Each bit represented the same unique keyword throughout the complete GA process. When a keyword was present in a document, the bit was set to 1, otherwise it was 0. Each document could then be represented in terms of a sequence of 0s and 1s. The keywords of five user-selected documents are presented below. The set of unique concepts present in these sample documents is also summarized—33 keywords (genes) in total. As in the Hopfield network example, some concepts were folder names assigned by the users (in the format of \*.\*); for example, QUERY.OPT folder for query optimization topics.

We computed the fitness of each document based on its relevance to the documents in the user-selected set. Higher Jaccard's score (a value between 0 and 1) indicated stronger relevance between two documents. For document 0, we computed five different Jaccard's scores between document 0 and documents 0, 1, 2, 3, and 4, respectively (shown below). An average fitness was then computed for document 0 (0.28774). The same procedure was applied to other documents to compute their fitness. A document which included more concepts shared by other documents had a higher Jaccard's score.

Jaccard's Score of DOC0 and DOC0 = 1.000000

Jaccard's Score of DOC0 and DOC1 = 0.120000

Jaccard's Score of DOC0 and DOC2 = 0.120000

Jaccard's Score of DOC0 and DOC3 = 0.115384

Jaccard's Score of DOC0 and DOC4 = 0.083333

Average Fitness (Jaccard's Score) of Document0: 0.28774

If a user provided documents that are closely related, the average fitness for the complete document set was high. If the user-selected documents were only loosely related, their overall fitness was low. Generally, GAs did a good job optimizing a document set which was initially low in fitness. Using the previous example, the overall Jaccard's score increased over generations. The optimized population contained only one single chromosome, with an average fitness value of 0.45121. The optimized chromosome contained six relevant keywords

- DOCO DATA RETRIEVAL, DATABASE, COMPUTER NETWORKS, IMPROVEMENTS, INFORMATION RETRIEVAL, METHOD, NETWORK, MULTIPLE, QUERY, RELATION, RELATIONAL, RETRIEVAL, QUERIES, RELATIONAL DATABASES, RELATIONAL DATABASE, US, CARAT.DAT, GQP.DAT, ORUS.DAT, QUERY.OPT
- DOC1 INFORMATION, INFORMATION RETRIEVAL, INFORMATION STORAGE, INDEXING, RE-TRIEVAL, STORAGE, US, KEVIN.HOT
- DOC2 ARTIFICIAL INTELLIGENCE, INFORMATION RETRIEVAL SYSTEMS, INFORMATION RETRIEVAL, INDEXING, NATURAL LANGUAGE PROCESSING, US, DBMS, AI, GOP. DAT
- DOC3 FUZZY SET THEORY, INFORMATION RETRIEVAL SYSTEMS, INDEXING, PERFORMANCE, RETRIEVAL SYSTEMS, RETRIEVAL, QUERIES, US, KEVIN.HOT
- DOC4 INFORMATION RETRIEVAL SYSTEMS, INDEX-ING, RETRIEVAL, STAIRS, US, KEVIN.HOT

Total Set of Concepts \_

DATA RETRIEVAL, DATABASE, COMPUTER, NETWORKS, IMPROVEMENTS, INFORMATION RETRIEVAL, METHOD, NETWORK, MULTIPLE, QUERY, RELATION, RELATIONAL, RETRIEVAL, QUERIES, RELATIONAL DATABASES, RELATIONAL DATABASE, US, CARAT.DAT, GQP.DAT, ORUS.DAT, QUERY.OPT, INFORMATION, INFORMATION STORAGE, INDEXING, STORAGE, KEVIN.HOT, ARTIFICIAL INTELLIGENCE, INFORMATION RETRIEVAL SYSTEMS, NATURAL LANGUAGE PROCESSING, DBMS.AI, FUZZY SET THEORY, PERFORMANCE, RETRIEVAL SYSTEMS, STAIRS.

\_\_\_ Initial Genetic Pattern of Chromosome in Population \_\_\_

chromosome	fitness
1111111111111111111111000000000000000	[0.287744]
000010000001000100001111100000000	[0.411692]
00001000000000101000010011110000	[0.367556]
00000000001100100000010101001110	[0.427473]
00000000001000100000010101000001	[0.451212]
Average Fitness = $0.3891$	

which best described the initial set of documents. Using these "optimized" keywords, an information retrieval system could proceed to suggest relevant documents to users. The user–GA interaction continued until a search was completed or the user decided to stop.

Table 4 summarizes the results of a benchmark testing. In the testing we randomly retrieved five test cases of 1-document, 2-document, 3-document, 4-document, 5-document, and 10-document examples, respectively, from the 3000-document DIALOG-extracted database discussed earlier. There were 30 test cases in total. For each test case, an initial fitness based on the Jaccard's score was computed. For 1-document and 2-document

chromosome	fitness
00000000001000100000010101000001	[0.45121]
00000000001000100000010101000001	[0.45121]
00000000001000100000010101000001	[0.45121]
00000000001000100000010101000001	[0.45121]
00000000001000100000010101000001	[0.45121]
Average Fitness = $0.4512$	

\_\_\_\_\_ Derived Concepts from Optimized Population\_\_\_\_

RETRIEVAL, US, INDEXING, KEVIN.HOT, INFORMATION RETRIEVAL SYSTEMS, STAIRS.

test cases, their initial fitness tended to be higher due to the smaller sample size (see column 2 of Table 4). In Table 4 we also report performance measures in terms of Jaccard scores for the GA processes, the CPU times, and the average improvements in fitness.

Using the GA optimization process, our system achieved an average fitness improvement from 5.38% to 17.7%. This improvement was slightly worse than the performance improvement for indexing reported by Gordon (1988). An interesting observation was that when more initial documents were present, the initial fitness tended to be lower, which allowed the system to do a better job in improving the preciseness of the initial keywords and in identifying other relevant documents. As shown in Table 4, fitness improvement increased as a function of the number of initial documents. This finding also suggested that when initial user-supplied documents are fuzzy and not well articulated, GAs may be able to make a more significant contribution in suggesting other relevant documents. This could be quite important for complex information retrieval sessions during which searchers need help in query articulation and search refinement.

The number of documents suggested by GANNET after the first GA process was between 9 and 13, with an average of about 11 documents. The CPU times required of the GA process also was quite reasonable, with an average of 0.168 seconds. The response times were significantly better than the Hopfield net activation. In conclusion, by using reproduction and the genetic operators, GAs provided an interesting system-aided way of analyzing users' intermediate search results and suggesting other potentially relevant documents.

#### **Conclusion and Future Directions**

Information retrieval research has been advancing very quickly over the past few decades. Researchers have experimented with techniques ranging from probabilistic models and the vector space model to the knowledge-based approach and the recent machine learning tech-

TABLE 4. Results of genetic algorithms testing.

No.	Init. score	GA score	Impr. (%)	CPU (sec.)	Docs. selected
1	1.0	1.0	0.0	0.067	7
1 2	1.0	1.0	0.0 0.0	0.067 0.05	25
	1.0	1.0		0.05	23 7
3 4	1.0	1.0	0.0 0.0	0.067	9
5	1.0	1.0	0.0	0.03	5
J I doc.	1.0	1.0 1.0	0.0	0.067	10.6
i doc.	avg.	1.0	0.0	0.00	10.0
1	0.5139	0.5139	0.0	0.083	10
2	0.5833	0.5833	0.0	0.1	8
3	0.6111	0.6111	0.0	0.083	5
4	0.6486	0.6486	0.0	0.067	10
5	0.7857	0.7857	0.0	0.083	16
2 docs.	avg.	0.6285	0.0	0.08	9.8
1	0.3841	0.3984	3.72	0.023	8
2	0.4157	0.4360	4.88	0.1	5
3	0.4286	0.4611	7.1	0.1	13
4	0.5032	0.5215	3.6	0.133	5
5	0.5899	0.6349	7.6	0.083	16
3 docs.	avg.	0.4904	5.38	0.088	9.4
1	0.2898	0.3010	3.8	0.117	22
2	0.3078	0.3142	2.1	0.1	15
3	0.3194	0.3495	9.4	0.283	5
4	0.3319	0.3442	3.7	0.25	11
5	0.4409	0.5060	14.7	0.25	10
4 docs.	avg.	0.3629	6.74	0.2	12.6
1	0.3048	0.3370	10.5	0.4	12
2	0.3068	0.3267	6.4	0.15	7
3	0.3194	0.3575	11.9	0.52	5
4	0.4655	0.5671	21.8	0.3	21
5	0.6181	0.7171	16.0	0.12	21
5 docs.	avg.	0.4610	13.32	0.298	13.2
1	0.2489	0.2824	13.5	0.32	18
2	0.2038	0.2282	12.9	0.35	8
3	0.2016	0.2343	16.2	0.47	6
4	0.4997	0.6201	24.1	0.13	5
5	0.3727	0.4540	21.8	0.13	11
10 docs.	avg.	0.3638	17.7	0.28	9.6
All	avg.	0.5511	7.19	0.168	10.87

niques. At each stage, significant insights regarding how to design more useful and "intelligent" information retrieval systems have been gained.

In this article, we presented an extensive review of IR research that was based mainly on machine learning techniques. Connectionist modeling and learning, in particular, has attracted considerable attention due to its strong resemblance to some existing IR models and techniques. Symbolic machine learning and genetic algorithms, two popular candidates for adaptive learning in other applications, on the other hand, have been used only rarely. However, these newer techniques have been found to exhibit promising inductive learning capabilities for selected IR applications.

For researchers who are interested in examining these techniques, this study has discussed an algorithmic approach and knowledge representations appropriate for IR. We feel that the proper selection of knowledge representation and the adaptation of machine learning algorithms in the IR context are essential to the successful use of such techniques. For example, in IR a keyword could represent a node in the Hopfield net, a single bit in a genetic algorithm, or a decision node in ID3 and ID5R. Similarly, the *parallel relaxation* search of the Hopfield net, the *entropy reduction* scheme in ID3, and the *Darwinian selection* of genetic algorithms all need to be carefully studied and modified in the unique IR context.

Despite some initially successful application of se-

lected machine learning techniques for IR, there are numerous research directions that need to be pursued before we can develop a robust solution to "intelligent" information retrieval. We briefly review several important research directions below:

• Limitations of learning techniques for IR: The performance of the inductive learning techniques relies strongly on the examples provided (as in any other statistical and classification techniques) (Weiss & Kulikowski, 1991). In IR, these examples may include userprovided queries and documents collected during relevance feedback. The importance of sample size has been stressed heavily, even in the probabilistic models (Fuhr & Buckley, 1991; Fuhr & Pfeifer, 1994). In reality, user-provided relevance feedback information may be limited in quantity and noisy (i.e., contradictory or incorrect), which may have adverse effects for the IR or indexing tasks. Some learning techniques such as the neural networks approach have documented noise-resistant capability, but empirical evidence and research need to be performed to verify this characteristic in the context of IR and indexing. In our preliminary investigation, all three machine learning algorithms performed satisfactorily for small document samples, but the effect of the sample size needs to be examined more carefully.

For large-scale real-life applications, neural networks and, to some extent, genetic algorithms, may suffer from requiring extensive computation time and lack of interpretable results. Symbolic learning, on the other hand, efficiently produces simple production rules or decision-tree representations. The effects of the representations on the cognition of searchers in the real-life retrieval environments (e.g., users' acceptance of the analytical results provided by an intelligent system) remain to be determined.

Applicability to the full-text retrieval environment: In addition to extensive IR research conducted in probabilistic models, knowledge-based systems, and machine learning, significant efforts have also been made by many commercial companies in pursuit of more effective and "intelligent" information retrieval systems. In an attempt to understand the potential role of machine learning in commercial full-text retrieval systems, we examined several major full-text retrieval software packages on the market, including: BRS/SEARCH, BASIS/Plus, PixTex, and Topic.

Most full-text retrieval software has been designed to handle large volumes of text by indexing every word (and its position). This allows users to perform proximity search, morphological search (using prefix, suffix, or wildcards), and thesaurus search. BRS/SEARCH and BASIS/plus are typical of this type of software. PixTex and Topic, on the other hand, are

among the most advanced full-text retrieval systems and feature "content-based IR" and "learning" capabilities. PixTex calls its indexing process "learning." The system automatically extracts patterns from binary data (texts or images) and associates (or "learns") the storage location of the data based on neural network technology (the exact form and algorithm are not clear due to the lack of publications on and the proprietary nature of the product). By automatically storing visual scene or textual contents in terms of Huffman codes, the system can then retrieve other similar scene objects or texts during IR. Verity's Topic claims to use fuzzy logic in its design of "conceptual searching" for "intelligent" document retrieval systems. It allows users to create and reuse hierarchical, weighted query trees (thus becoming part of the corporate memory), which produce rank-ordered documents. It also appears to have some "similarity search" capability (e.g., "find me all documents like this one"). However, like PixTex, no algorithmic detail can be obtained. Despite the lack of implementation detail, we believe that with the extensive indexing capabilities provided by such full-text software, a simple user relevance feedback component and inductive machine learning algorithms, similar to the ones discussed in this research, could be incorporated to help identify what users want, based on the concepts (keywords) learned from the sample documents. As more researchers and practitioners recognize the need for concept-based and "intelligent" IR, application of machine learning algorithms presents unique challenges and opportunities.

We believe this research has shed light on the feasibility and usefulness of the newer, AI-based machine learning algorithms for IR. However, more extensive and systematic studies of various system parameters and for large-scale, real-life applications are needed. We hope by incorporating into IR inductive learning capabilities, which are complementary to the prevailing full-text, keyword-based, probabilistic, or knowledge-based techniques, we will be able to advance the design of adaptive and "intelligent" information retrieval systems.

# **Acknowledgments**

This project was supported mainly by NSF Grant #IRI-9211418, 1992–1994 (NSF/CISE, Division of Information, Robotics, and Intelligent Systems).

#### References

Appelt, D. (1985, August). The role of user modelling in language generation and communication planning. In *User Modelling Panel, Proceedings of the Ninth International Joint Conference on Artificial Intelligence*, (pp. 1298–1302). Los Altos, CA: Morgan Kaufmann Publishers. Inc.

Belew, R. K. (1989, June). Adaptive information retrieval. In Proceedings of the Twelfth Annual International ACM/SIGIR Conference on

<sup>&</sup>lt;sup>1</sup> Vended by BRS Software Products, McLean, VA, USA.

<sup>&</sup>lt;sup>2</sup> Vended by Information Dimensions Inc., Dublin, OH, USA.

<sup>&</sup>lt;sup>3</sup> Vended by Excalibur Technologies Corp., McLean, VA, USA.

<sup>&</sup>lt;sup>4</sup> Vended by Verity, Inc., Mountain View, CA, USA.

- Research and Development in Information Retrieval (pp. 11–20). NY, NY: ACM Press.
- Blair, D. C., & Maron, M. E. (1985). An evaluation of retrieval effectiveness for a full-text document-retrieval system. *Communications of the ACM*, 28, 289–299.
- Blosseville, M. J., Hebrail, G., Monteil, M. G., & Penot, N. (1992, June). Automatic document classification: Natural language processing, statistical analysis, and expert system techniques used together. In Proceedings of the Fifteenth Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval (pp. 51-57). NY, NY: ACM Press.
- Booker, L. B., Goldberg, D. E., & Holland, J. H. (1990). Classifier systems and genetic algorithms. In J. G. Carbonell (Ed.), *Machine learning, paradigms and methods* (pp. 235–282). Cambridge, MA: The MIT Press.
- Bookstein, A., & Swanson, D. R. (1975). Probabilistic models for automatic indexing. *Journal of the American Society for Information Science*, 26, 45–50.
- Borgida, A., & Williamson, K. E. (1985, August). Accommodating exceptions in a database, and refining the schema by learning from them. In *Proceedings of the 11th International VLDB Conference* (pp. 72–81). Saratoga, NY: VLDB Endowment.
- Brajnik, G., Guida, G., & Tasso, C. (1988). IR-NLI II: Applying manmachine interaction and artificial intelligence concepts to information retrieval. In *Proceedings of the Eleventh Annual International* ACM/SIGIR Conference on Research and Development in Information Retrieval (pp. 387–399). NY, NY: ACM Press.
- Brauen, T. L. (1971). Document vector modification. In G. Salton (Ed.), The Smart retrieval system—experiments in automatic document processing (pp. 456–484). Englewood Cliffs, NJ: Prentice-Hall.
- Breiman, L., Friedman, J., Olshen, R., & Stone, C. (1984). Classification and regression tree. Monterey, CA: Wadsworth.
- Buckland, M. K., & Florian, D. (1991). Expertise, task complexity, and artificial intelligence: A conceptual framework. *Journal of the American Society for Information Science*, 42, 635–643.
- Cai, Y., Cercone, N., & Han, J. (1991). Attribute-oriented induction in relational databases. In G. Piatetsky-Shapiro & W. J. Frawley (Eds.), Knowledge discovery in databases (pp. 213–228). Cambridge, MA: The MIT Press.
- Carbonell, J. G., Michalski, R. S., & Mitchell, T. M. (1983). An overview of machine learning. In R. S. Michalski, J. G. Carbonell, & T. M. Mitchell (Eds.), *Machine Learning, An Artificial Intelligence Approach* (pp. 3-23). Palo Alto, CA: Tioga.
- Chen, H., Basu, K., & Ng, T. (in press-a). An algorithmic approach to concept exploration in a large knowledge network (automatic thesaurus consultation): Symbolic branch-and-bound vs. connectionist Hopfield net activation. Journal of the American Society for Information Science.
- Chen, H., Buntin, P., She, L., Sutjahjo, S., Sommer, C., & Neely, D. (in press). Expert prediction, symbolic learning, and neural networks: An experiment on greyhound racing. *IEEE Expert*.
- Chen, H., & Dhar, V. (1987, July). Reducing indeterminism in consultation: a cognitive model of user/librarian interaction. In *Proceedings of the 6th National Conference on Artificial Intelligence (AAAI-87)* (pp. 285–289). Los Altos, CA: Morgan Kaufmann Publishers, Inc.
- Chen, H., & Dhar, V. (1990). User misconceptions of online information retrieval systems. *International Journal of Man–Machine Studies*, 32, 673–692.
- Chen, H., & Dhar, V. (1991). Cognitive process as a basis for intelligent retrieval systems design. *Information Processing and Management*, 27, 405–432.
- Chen, H., Hsu, P., Orwig, R., Hoopes, L., & Nunamaker, J. F. (1994b). Automatic concept classification of text from electronic meetings. Communications of the ACM, 37, 56-73.
- Chen, H., & Kim, J. (1993). GANNET: Information retrieval using genetics algorithms and neural networks. (Working Paper, CMI-WPS).

- Center for Management of Information, College of Business and Public Administration, University of Arizona.
- Chen, H., & Lynch, K. J. (1992). Automatic construction of networks of concepts characterizing document databases. *IEEE Transactions* on Systems, Man and Cybernetics, 22, 885–902.
- Chen, H., Lynch, K. J., Basu, K., & Ng, T. (1993). Generating, integrating, and activating thesauri for concept-based document retrieval. *IEEE Expert (special series on Artificial Intelligence in Text-Based Information Systems)*, 8, 25-34.
- Chen, H., & She, L. (1994, January). Inductive query by examples (IQBE): A machine learning approach. In *Proceedings of the 27th Annual Hawaii International Conference on System Sciences (HICSS-27), Information Sharing and Knowledge Discovery Track.* Los Alamitos, CA: IEEE Computer Society Press.
- Chiaramella, Y., & Defude, B. (1987). A prototype of an intelligent system for information retrieval: IOTA. *Information Processing and Management*, 23, 285–303.
- Cohen, P. R., & Kjeldsen, R. (1987). Information retrieval by constrained spreading activation in semantic networks. *Information Processing and Management*, 23, 255–268.
- Crawford, S. L., Fung, R., Appelbaum, L. A., & Tong, R. M. (1991). Classification trees for information retrieval. In *Proceedings of the 8th International Workshop on Machine Learning* (pp. 245–249). San Mateo, CA: Morgan Kaufmann.
- Crawford, S. L., & Fung, R. M. (1992). An analysis of two probabilistic model induction techniques. *Statistics and Computing*, 2, 83–90.
- Croft, W. B., & Thompson, R. H. (1987). I<sup>3</sup>R: A new approach to the design of document retrieval systems. Journal of the American Society for Information Science, 38, 389-404.
- Dalton, J., & Deshmane, A. (1991). Artificial neural networks. *IEEE Potentials*, 10, 33–36.
- Daniels, P. J. (1986). The user modelling function of an intelligent interface for document retrieval systems. In B. C. Brookes (Ed.), *Intelligent information systems for the information society*. Amsterdam: Fisevier
- Derthick, M. (1988). Mundane reasoning by parallel constraint satisfaction. Ph.D. thesis, Carnegie Mellon University, Pittsburgh, PA.
- Doszkocs, T. E., Reggia, J., & Lin, X. (1990). Connectionist models and information retrieval. Annual Review of Information Science and Technology (ARIST), 25, 209-260.
- Everitt, B. (1980). Cluster analysis (2nd ed.). London: Heinemann.
- Fisher, D. H., & McKusick, K. B. (1989, August). An empirical comparison of ID3 and backpropagation. In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence (IJCAI-89)* (pp. 788–793). San Mateo, CA: Morgan Kaufmann Publishers, Inc.
- Fogel, D. B. (1994). An introduction to simulated evolutionary optimization. *IEEE Transactions on Neural Networks*, 5, 3–14.
- Fogel, L. J. (1962). Autonomous automata. *Industrial Research*, 4. (pp. 14–19)
- Fogel, L. J. (1964). On the organization of intellect. Doctoral dissertation, UCLA, Los Angeles, CA.
- Fox, E. A. (1987). Development of the CODER system: A testbed for artificial intelligence methods in information retrieval. *Information Processing and Management*, 23, 341–366.
- Frawley, W. J., Pietetsky-Shapiro, G., & Matheus, C. J. (1991). Knowledge discovery in databases: An overview. In G. Piatetsky-Shapiro & W. J. Frawley (Eds.), *Knowledge discovery in databases* (pp. 1–30). Cambridge, MA: The MIT Press.
- Freund, J. E. (1971). *Mathematical statistics*. Englewood Cliffs, NJ: Prentice-Hall.
- Frieder, O., & Siegelmann, H. T. (1991, October). On the allocation of documents in multiprocessor information retrieval systems. In Proceedings of the Fourteenth Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval (pp. 230–239). NY, NY: ACM Press.
- Fuhr, N., & Buckley, C. (1991). A probabilistic learning approach for

- document indexing. ACM Transactions on Information Systems, 9, 223-248.
- Fuhr, N., Hartmann, S., Knorz, G., Lustig, G., Schwantner, M., & Tzeras, K. (1990, July-August). AIR/X—a rule-based multistage indexing system for large subject fields. In *Proceedings of the Eighth National Conference on Artificial Intelligence (AAAI-90)* (pp. 789–795). Boston, MA.
- Fuhr, N., & Pfeifer, U. (1994). Probabilistic information retrieval as a combination of abstraction, inductive learning, and probabilistic assumptions. ACM Transactions on Information Systems, 12, 92– 115.
- Fung, R., & Crawford, S. L. (1990, July-August). Constructor: A system for the induction of probabilistic models. In *Proceedings of the 8th National Conference on Artificial Intelligence (AAAI-90)* (pp. 762–769). Boston, MA.
- Gallant, S. I. (1988). Connectionist expert system. Communications of the ACM, 31, 152–169.
- Goldberg, D. E. (1989). Genetic algorithms in search, optimization, and machine learning. Reading, MA: Addison-Wesley.
- Goldberg, D. E. (1994). Genetic and evolutionary algorithms come of age. *Communications of the ACM*, 37, 113–119.
- Gordon, M. (1988). Probabilistic and genetic algorithms for document retrieval. *Communications of the ACM*, 31. 1208–1218.
- Gordon, M. D. (1991). User-based document clustering by redescribing subject descriptions with a genetic algorithm. *Journal of the Ameri*can Society for Information Science, 42, 311–322.
- Greene, D. P., & Smith, S. F. (1992). COGIN: Symbolic induction with genetic algorithms. In *Proceedings of the Tenth National Conference* on Artificial Intelligence (AAAI-92) (pp. 111–116). Cambridge, MA: The MIT Press.
- Hall, L. O., & Romaniuk, S. G. (1990, July–August). A hybrid connectionist, symbolic learning system. In *Proceedings of the Eighth National Conference on Artificial Intelligence (AAAI-90)* (pp. 783–788). Cambridge, MA: The MIT Press.
- Han, J., Cai, Y., & Cercone, N. (1993). Data-driven discovery of quantitative rules in relational databases. *IEEE Transactions on Knowledge and Data Engineering*, 5, 29–40.
- Harp, S., Samad, T., & Guha, A. (1989). Towards the genetic synthesis of neural networks. In *Proceedings of the Third International Con*ference on Genetic Algorithms. San Mateo, CA: Morgan Kaufmann.
- Hayes-Roth, F., & Jacobstein, N. (1994). The state of knowledge-based systems. *Communications of the ACM*, 37, 27–39.
- Holland, J. H. (1975). Adaptation in natural and artificial systems. Ann Arbor, MI: University of Michigan Press.
- Hopfield, J. J. (1982). Neural network and physical systems with collective computational abilities. *Proceedings of the National Academy of Science*, USA, 78(8) (pp. 2554–2558).
- Humphreys, B. L., & Lindberg, D. A. (1989, November). Building the unified medical language system. In *Proceedings of the Thirteenth Annual Symposium on Computer Applications in Medical Care*. Washington, DC: IEEE Computer Society Press.
- Ide, E. (1971). New experiments in relevance feedback. In G. Salton (Ed.), The Smart retrieval system—experiments in automatic document processing (pp. 337–354). Englewood Cliffs, NJ: Prentice-Hall.
- Ide, E., & Salton, G. (1971). Interactive search strategies and dynamic file organization in information retrieval. In G. Salton (Ed.), The Smart retrieval system—experiments in automatic document processing (pp. 373-393). Englewood Cliffs, NJ: Prentice-Hall.
- Ioannidis, Y. E., Saulys, T., & Whitsitt, A. J. (1992). Conceptual learning in database design. ACM Transactions on Information Systems, 10, 265–293.
- Kitano, H. (1990, July-August). Empirical studies on the speed of convergence of neural network training using genetic algorithms. In *Proceedings of the Eighth National Conference on Artificial Intelligence (AAAI-90)* (pp. 789–795). Cambridge, MA: The MIT Press.
- Knight, K. (1990). Connectionist ideas and algorithms. Communications of the ACM, 33, 59-74.

- Kohonen, T. (1989). Self-organization and associative memory (3rd ed.). Berlin: Springer-Verlag.
- Koza, J. R. (1992). Genetic programming: On the programming of computers by means of natural selection. Cambridge, MA: The MIT Press.
- Kwok, K. L. (1989, June). A neural network for probabilistic information retrieval. In *Proceedings of the Twelfth Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval* (pp. 21–30). NY, NY: ACM Press.
- Lebowitz, M. (1987). Concept learning in a rich input domain: Generalization-based memory. In J. G. Carbonell, R. S. Michalski, & T. M. Mitchell (Eds.), Machine learning, an artificial intelligence approach (vol. II) (pp. 193–214, 463–482). Los Altos, CA: Morgan Kaufmann.
- Lewis, D. D. (1991). Learning in intelligent information retrieval. In Proceedings of the 8th International workshop on machine learning (pp. 235–239). Los Altos, CA: Morgan Kaufmann.
- Lewis, D. D. (1992, June). An evaluation of phrasal and clustered representations on a text categorization task. In *Proceedings of the Fifteenth Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval* (pp. 37–50). NY, NY: ACM Press.
- Li, Q., & McLeod, D. (1989). Object flavor evolution through learning in an object-oriented database system. In L. Kerschberg (Ed.), Expert Database Systems, Proceedings from the Second International Conference (pp. 469–495). Menlo Park, CA: Benjamin/Cummings.
- Lin, X., Soergel, D., & Marchionini, G. (1991, October). A self-organizing semantic map for information retrieval. In *Proceedings of the Fourteenth Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval* (pp. 262–269). Chicago, IL.
- Lindberg, D. A., & Humphreys, B. L. (1990, November). The UMLS knowledge sources: Tools for building better user interface. In Proceedings of the Fourteenth Annual Symposium on Computer Applications in Medical Care. Los Alamitos, CA: Institute of Electrical and Electronics Engineers.
- Lippmann, R. P. (1987). An introduction to computing with neural networks. IEEE Acoustics Speech and Signal Processing Magazine, 4 4-27
- MacLeod, K. J., & Robertson, W. (1991). A neural algorithm for document clustering. *Information Processing & Management*, 27, 337–346.
- Maron, M. E., & Kuhns, J. L. (1960). On relevance, probabilistic indexing and information retrieval. *Journal of the ACM*, 7, 216–243.
- Martin, B. K., & Rada, R. (1987). Building a relational data base for a physician document index. *Medical Informatics*, 12, 187–201.
- Masand, B., Gordon, L., & Waltz, D. (1992, June). Classifying news stories using memory-based reasoning. In Proceedings of the Fifteenth Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval (pp. 59–65). Copenhagen, Denmark.
- McCray, A. T., & Hole, W. T. (1990). The scope and structure of the first version of the UMLS semantic network. In *Proceedings of the* Fourteenth Annual Symposium on Computer Applications in Medical Care. Los Alamitos, CA: Institute of Electrical and Electronics Engineers.
- Michalewicz, Z. (1992). Genetic algorithms + data structures = evolution programs. Berlin: Springer-Verlag.
- Michalski, R. S. (1983). A theory and methodology of inductive learning. In R. S. Michalski, J. G. Carbonell, and T. M. Mitchell (Eds.), Machine learning, an artificial intelligence approach (pp. 83–134). Palo Alto, CA: Tioga.
- Mitchell, T. M. (1982). Generalization as search. *Artificial Intelligence*, 18, 203–226.
- Monarch, I., & Carbonell, J. G. (1987). CoalSORT: A knowledge-based interface. *IEEE Expert*, 39–53.
- Montana, D. J., & Davis, L. (1989, August). Training feedforward neural networks using genetic algorithms. In *Proceedings of the Eleventh*

- International Joint Conference on Artificial Intelligence (IJCAI-89) (pp. 762-767). San Mateo, CA: Morgan Kaufmann Publishers, Inc.
- Montgomery, D. D. (1976). Design and analysis of experiments. New York; Wiley.
- Mooney, R., Shavlik, J., Towell, G., & Gove, A. (1989, August). An experimental comparison of symbolic and connectionist learning algorithms. In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence (IJCAI-89)* (pp. 775–780). San Mateo, CA: Morgan Kaufmann Publishers.
- Parsaye, K., Chignell, M., Khoshafian, S., & Wong, H. (1989). *Intelligent databases*. New York: Wiley.
- Petry, F., Buckles, B., Prabhu, D., & Kraft, D. (1993). Fuzzy information retrieval using genetic algorithms and relevance feedback. In *Proceedings of the ASIS Annual Meeting* (pp. 122-125) Medford, NJ: ASIS.
- Piatetsky-Shapiro, G. (1989). Workshop on knowledge discovery in real databases. In *International Joint Conference of Artificial Intelligence*. San Matco, CA: Morgan Kaufmann Publishers.
- Pollitt, S. (1987). Cansearch: An expert systems approach to document retrieval. *Information Processing and Management*, 23, 119–138.
- Quinlan, J. R. (1979). Discovering rules by induction from large collections of examples. In D. Michie (Ed.), Expert systems in the microelectronic age (pp. 168-201). Edinburgh: Edinburgh University Press.
- Quinlan, J. R. (1983). Learning efficient classification procedures and their application to chess end games. In R. S. Michalski, J. G. Carbonell, & T. M. Mitchell (Eds.), *Machine learning, an artificial intelli*gence approach (pp. 463–482). Palo Alto, CA: Tioga.
- Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1, 81–106.
- Quinlan, J. R. (1993). C4.5: Programs for machine learning. Los Altos, CA: Morgan Kaufmann.
- Rada, R., Mili, H., Bicknell, E., & Blettner, M. (1989). Development and application of a metric on semantic nets. *IEEE Transactions on Systems, Man, and Cybernetics*, 19, 17–30.
- Raghavan, V. V., & Agarwal, B. (1987, July). Optimal determination of user-oriented clusters: An application for the reproductive plan. In Proceedings of the Second International Conference on Genetic Algorithms and their Applications (pp. 241–246). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Rau, L. F., & Jacobs, P. S. (1991, October). Creating segmented databases from free text for text retrieval. In *Proceedings of the Four*teenth Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval (pp. 337–346). NY, NY: ACM Press.
- Rich, E. (1979, August). Building and exploiting user models. In *International Joint Conference of Artificial Intelligence* (pp. 720–722). Tokyo, Japan.
- Rich, E. (1979b) User modeling via stereotypes. Cognitive Science, 3, 329–354.
- Rich, E. (1983). Users are individuals: Individualizing user models. *International Journal of Man–Machine Studies*, 18, 199–214.
- Roberston, S. E., & Sparck Jones, K. (1976). Relevance weighting of search terms. *Journal of the American Society for Information Sci*ence, 27, 129–146.
- Rocchio, J. J. (1971). Relevance feedback in information retrieval. In G. Salton (Ed.), The Smark retrieval system—experiments in automatic document processing (pp. 313-323). Englewood Cliffs, NJ: Prentice-Hall.
- Rose, D. E., & Belew, R. K. (1991). A connectionist and symbolic hybrid for improving legal research. *International Journal of Man–Machine Studies*, 35, 1–33.
- Rumelhart, D. E., Hinton, G. E., & McClelland, J. L. (1986). A general framework for parallel distributed processing. In D. E. Rumelhart, & J. L. McClelland, & the PDP Research Group (Eds.), *Parallel distributed processing* (pp. 45–76). Cambridge, MA: The MIT Press.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning

- internal representations by error propagation. In D. E. Rumelhart, J. L. McClelland, & the PDP Research Group (Eds.), *Parallel distributed processing* (pp. 318–362). Cambridge, MA: The MIT Press.
- Rumelhart, D. E., Widrow, B., & Lehr, M. A. (1994). The basic ideas in neural networks. *Communications of the ACM*, 37, 87–92.
- Salton, G. (1989). Automatic text processing. Reading, MA: Addison-Wesley.
- Shastri, L. (1991). Why semantic networks? In J. F. Sowa (Ed.), Principles of semantic networks: Explorations in the representation of knowledge (pp. 109-136). San Mateo, CA: Morgan Kaufmann.
- Simon, H. (1991). Artificial intelligence: Where has it been, and where is it going? *IEEE Transactions on Knowledge and Data Engineering*, 3, 128–136.
- Simpson, P. K. (1990). Artificial neural sytems: Foundations, paradigms, applications, and implementations. New York: McGraw-Hill.
- Sleeman, D. (1985). UMFE: A user modeling front-end subsystem. *International Journal of Man–Machine Studies*, 23, 71–88.
- Smith, P. J., Shute, S. J., Galdes, D., & Chignell, M. H. (1989). Knowledge-based search tactics for an intelligent intermediary system. ACM Transactions on Information Systems, 7, 246–270.
- Sparck Jones, K. (1991). The role of artificial intelligence in information retrieval. *Journal of the American Society for Information Science*, 42, 558–565.
- Stepp, R. E., & Michalski, R. S. (1987). Conceptual clustering: Inventing goal-oriented classifications of structured objects. In J. G. Carbonell et al. (Eds.), *Machine learning, an artificial intelligence approach* (Vol. II) (pp. 472–498, 463–482). Los Altos, CA: Morgan Kaufmann.
- Swartout, W. (1985, August). Explanation and the role of the user model: how much will it help? In *User Modelling Panel, Proceedings* of the Ninth International Joint Conference on Artificial Intelligence (pp. 1298–1302). Los Altos, CA: Morgan Kaufmann Publishers, Inc.
- Tank, D. W., & Hopfield, J. J. (1987). Collective computation in neuronlike circuits. Scientific American, 257, 104–114.
- Touretzky, D., & Hinton, G. E. (1988). A distributed connectionist production system. *Cognitive Science*, 12, 423-466.
- Turtle, H., & Croft, W. B. (1990, September). Inference networks for document retrieval. In *Proceedings of the 13th Annual International* ACM/SIGIR Conference on Research and Development in Information Retrieval (pp. 1–24). Brussels, Belgium. NY, NY: ACM Press.
- Turtle, H., & Croft, W. B. (1991). Evaluation of an inference network-based retrieval model. ACM Transactions on Information Systems, 9, 187–222.
- Tzeras, K., & Hartmann, S. (1993, June-July). Automatic indexing based on Bayesian inference networks. In Proceedings of the 16th Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval (pp. 22-34). NY, NY: ACM Press.
- Utgoff, P. E. (1989). Incremental induction of decision trees. *Machine Learning*, 4, 161–186.
- van Rijsbergen, C. J. (1979). *Information retrieval* (2nd ed.). London: Butterworths.
- Vickery, A., & Brooks, H. M. (1987). PLEXUS—the expert system for referral. *Information Processing and Management*, 23, 99–117.
- Weiss, S. M., & Kapouleas, I. (1989, August). An empirical comparison of pattern recognition, neural nets, and machine learning classification methods. In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence (IJCAI-89)* (pp. 781–787). San Mateo, CA: Morgan Kaufmann Publishers, Inc.
- Weiss, S. M., & Kulikowski, C. A. (1991). Computer systems that learn: Classification and prediction methods from statistics, neural networks, machine learning, and expert systems. San Mateo, CA: Morgan Kaufmann.
- Widrow, B., Rumelhart, D. E., & Lehr, M. A. (1994). Neural networks: Applications in industry, business, and science. *Communications of the ACM*, 37, 93–105.

- Wilkinson, R., & Hingston, P. (1991, October). Using the Cosine measure in neural network for document retrieval. In *Proceedings of the Fourteenth Annual International ACM/SIGIR Conference on Research and Development in Information Retrieval* (pp. 202–210). Chicago, IL.
- Wilkinson, R., Hingston, P., & Osborn, T. (1992). Incorporating the vector space model in a neural network used for document retrieval. *Library Hi Tech*, 10, 69–75.
- Yang, J., & Korshage, R. R. (1993, April). Effects of query term weights modification in document retrieval: A study based on a genetic algorithm. In *Proceedings of the Second Annual Symposium on Docu-*
- ment Analysis and Information Retrieval (pp. 271–285). Las Vegas, NV: University of Nevada.
- Yang, J., Korfhage, R. R., & Rasmussen, E. (1993, November). Query improvement in information retrieval using genetic algorithms: A report on the experiments of the TREC Project. In *Text Retrieval Conference (TREC-1)* (pp. 31–58). Washington, DC: NIST.
- Yu, C. T., & Salton, G. (1976). Precision weighting: An effective automatic indexing method. *Journal of the ACM*, 23, 76–88.
- Zissos, A. Y., & Witten, I. H. (1985). User modeling for a computer coach: A case study. *International Journal of Man—Machine Studies*, 23, 729–750.