

8 FEATURE SUBSET SELECTION USING A GENETIC ALGORITHM

Jihoon Yang and Vasant Honavar

AI Research Group, Department of Computer Science
226 Atanasoff Hall, Iowa State University
Ames, IA 50011, U.S.A.

{yang|honavar}@cs.iastate.edu

Abstract: Pattern classification and knowledge discovery problems require selection of a subset of features to represent the patterns to be classified. This is due to the fact that the performance of the classifier and the cost of classification are sensitive to the choice of the features used to construct the classifier. Genetic algorithms offer an attractive approach to find near-optimal solutions to such optimization problems. This chapter presents an approach to feature subset selection using a genetic algorithm. Our experiments demonstrate the feasibility of this approach to feature subset selection in the automated design of neural networks for pattern classification and knowledge discovery.

8.1 INTRODUCTION

Many practical pattern classification tasks (e.g., medical diagnosis) require learning of an appropriate classification function that assigns a given input pattern (typically represented using a vector of attribute or feature values) to one of a finite set of classes. The choice of features, attributes, or measurements used to represent patterns that are presented to a classifier affect (among other things):

- The accuracy of the classification function that can be learned using an inductive learning algorithm: The features used to describe the patterns implicitly define a pattern language. If the language is not expressive enough, it would fail to capture the information that is necessary for classification and hence regardless of the learning algorithm used, the accuracy of the classification function learned would be limited by this lack of information.

- The time needed for learning a sufficiently accurate classification function: For a given representation of the classification function, the features used to describe the patterns implicitly determine the search space that needs to be explored by the learning algorithm. An abundance of irrelevant features can unnecessarily increase the size of the search space, and hence the time needed for learning a sufficiently accurate classification function.
- The number of examples needed for learning a sufficiently accurate classification function: All other things being equal, the larger the number of features used to describe the patterns in a domain of interest, the larger is the number of examples needed to learn a classification function to a desired accuracy (Langley, 1995; Mitchell, 1997).
- The cost of performing classification using the learned classification function: In many practical applications e.g., medical diagnosis, patterns are described using observable symptoms as well as results of diagnostic tests. Different diagnostic tests might have different costs as well as risks associated with them. For instance, an invasive exploratory surgery can be much more expensive and risky than say, a blood test.
- The comprehensibility of the knowledge acquired through learning: A primary task of an inductive learning algorithm is to extract *knowledge* (e.g., in the form of classification rules) from the training data. Presence of a large number of features, especially if they are irrelevant or misleading, can make the knowledge difficult to comprehend by humans. Conversely, if the learned rules are based on a small number of relevant features, they would much more concise and hence easier to understand, and use.

This presents us with a *feature subset selection problem* in automated design of pattern classifiers. The feature subset selection problem refers the task of identifying and selecting a useful subset of features to be used to represent patterns from a larger set of often mutually redundant, possibly irrelevant, features with different associated measurement costs and/or risks. An example of such a scenario which is of significant practical interest is the task of selecting a subset of clinical tests (each with different financial cost, diagnostic value, and associated risk) to be performed as part of a medical diagnosis task. Other examples of feature subset selection problem include large scale data mining applications, power system control (Zhou et al., 1997), construction of user interest profiles for text classification (Yang et al., 1998a) and sensor subset selection in the design of autonomous robots (Balakrishnan and Honavar, 1996).

8.2 RELATED WORK

A number of approaches to feature subset selection have been proposed in the literature. (See (Siedlecki and Sklansky, 1988; Doak, 1992; Langley, 1994; Dash and Liu, 1997) for surveys). These approaches involve searching for an optimal subset of features based on some criteria of interest. Feature subset selection problem can be viewed as a special case of the *feature weighting* problem. It

involves assigning a real-valued weight to each feature. The weight associated with a feature measures its relevance or significance in the classification task (Cost and Salzberg, 1993; Wettschereck et al., 1995). If we restrict the weights to be binary valued, the feature weighting problem reduces to the feature subset selection problem. The focus of this chapter is on feature subset selection.

Let $\mu(S)$ be a performance measure that is used to evaluate a feature subset S with respect to the criteria of interest (e.g., cost and accuracy of the resulting classifier). Feature subset selection problem is essentially an optimization problem which involves searching the space of possible feature subsets to identify one that is optimal or near-optimal with respect to μ . Feature subset selection algorithms can broadly be classified into three categories according to the characteristics of the search strategy employed.

8.2.1 Feature Subset Selection Using Exhaustive Search

In this approach, the candidate feature subsets are evaluated with respect to the performance measure μ and an *optimal* feature subset is found using exhaustive search. The Focus algorithm (Almuallim and Dietterich, 1994) employs the breadth-first search algorithm to find the minimal combination of features sufficient to construct a hypothesis that is consistent with the training examples. The algorithm proposed by (Sheinvald et al., 1990) uses the *minimum description length* criterion to select an optimal feature subset using exhaustive enumeration and evaluation of candidate feature subsets. Exhaustive search is computationally infeasible in practice, except in those rare instances where the total number of features is quite small.

8.2.2 Feature Subset Selection Using Heuristic Search

Since exhaustive search over all possible subsets of a feature set is not computationally feasible in practice, a number of authors have explored the use of *heuristics* for feature subset selection, often in conjunction with branch and bound search, a technique that is well-known in combinatorial optimization and artificial intelligence. *Forward selection* and *backward elimination* are the most common sequential branch and bound search algorithms used in feature subset selection (Narendra and Fukunaga, 1977; Foroutan and Sklansky, 1987). Forward selection starts with an empty feature set and adds a feature at a time, at each stage choosing the addition that most increases μ . Backward elimination starts with the entire feature set and at each step drops the feature whose absence least decreases μ . Both forward and backward selection procedures are optimal at each stage, but are unable to anticipate complex *interactions* between features that might affect the performance of the classifier. A related approach, called the *exchange strategy* starts with an initial feature subset (perhaps found by forward selection or backward elimination) and then tries to exchange a feature in the selected subset with one of the features that is outside it. We can often find a feature subset that is guaranteed to be the best for a given size of the feature subset without considering all possible subsets

using branch and bound search (Narendra and Fukunaga, 1977) if we assume that μ is monotone. That is, adding features is guaranteed to not decrease μ . It is worth pointing out that in many practical pattern classification scenarios, the monotonicity assumption is not satisfied. For example, addition of irrelevant features (e.g., social security numbers in medical records in a diagnosis task) can significantly worsen the generalization accuracy of a decision tree classifier (Mitchell, 1997). Furthermore, feature subset selection techniques that rely on the monotonicity of the performance criterion, although they appear to work reasonably well with linear classifiers, can exhibit poor performance with non-linear classifiers such as neural networks (Ripley, 1996).

Five greedy hillclimbing procedures (with different sequential search methods) for obtaining good generalization with decision tree construction algorithms (ID3 and C4.5) (Mitchell, 1997) were proposed in (Caruana and Freitag, 1994). In related work, (John et al., 1994) used both forward selection and backward elimination to minimize the cross validation error of decision tree classifiers; (Kohavi, 1994) used hillclimbing and best-first search for feature subset selection for decision tree classifiers. (Koller and Sahami, 1996; Koller and Sahami, 1997) used forward selection and backward elimination to select a feature that is subsumed by the remaining features (determined by the *Markov blanket*, the set of features that render the selected feature conditionally independent of the remaining features) for constructing Naive Bayesian (Mitchell, 1997) and decision tree classifiers. The *Preset* algorithm (Modrzejewski, 1993) employs the *rough set theory* to select a feature subset by rank ordering the features to generate a minimal decision tree.

8.2.3 Feature Subset Selection Using Randomized Search

Randomized algorithms make use of randomized or probabilistic (as opposed to deterministic) steps or sampling processes. Several researchers have explored the use of such algorithms for feature subset selection. The *Relief* algorithm (Kira and Rendell, 1992) assigns weights to features (based on their estimated effectiveness for classification) using the randomly sampled instances. Features whose weights exceed a user-determined threshold are selected in designing the classifier. Several extensions of *Relief* have been introduced to handle noisy or missing features as well as multi-category classification (Kononenko, 1994). A randomized hillclimbing search for feature subset selection for nearest neighbor classifiers (Cover and Hart, 1967; Dasarathy, 1991) was proposed in (Skalak, 1994). The LVF and LVW algorithms (Liu and Setiono, 1996b; Liu and Setiono, 1996a) are randomized algorithms that generate several random feature subsets and pick the one that has the least number of *unfaithful* patterns in the space defined by the feature subset (LVF) or the one that has the lowest error using a decision tree classifier (LVW) giving preference to smaller feature subsets. (Two patterns are said *unfaithful* if they have the same feature values but different class labels). Several authors have explored the use of randomized population-based heuristic search techniques such as genetic algorithms (GA) for feature subset selection for decision tree and nearest neighbor classifiers (Siedlecki and

Sklansky, 1989; Brill et al., 1992; Richeldi and Lanzi, 1996) or rule induction systems (Vafaie and De Jong, 1993). A related approach used *lateral feedback* networks (Guo, 1992) to evaluate feature subsets (Guo and Uhrig, 1992). Feature subset selection techniques that employ genetic algorithms do not require the restrictive monotonicity assumption. They also readily lend themselves to the use of multiple selection criteria. This makes them particularly attractive in the design of pattern classifiers in many practical scenarios.

8.2.4 Filter and Wrapper Approaches to Feature Subset Selection

Feature subset selection algorithms can also be classified into two categories based on whether or not feature selection is done independently of the learning algorithm used to construct the classifier. If feature selection is performed independently of the learning algorithm, the technique is said to follow a *filter* approach. Otherwise, it is said to follow a *wrapper* approach (John et al., 1994). While the filter approach is generally computationally more efficient than the wrapper approach, its major drawback is that an optimal selection of features may not be independent of the inductive and representational biases of the learning algorithm that is used to construct the classifier. The wrapper approach on the other hand, involves the computational overhead of evaluating candidate feature subsets by executing a selected learning algorithm on the dataset represented using each feature subset under consideration. This is feasible only if the learning algorithm used to train the classifier is relatively fast. Figure 8.1 summarizes the filter and wrapper approaches. The approach to feature subset selection proposed in this chapter is an instance of the wrapper approach. It utilizes a genetic algorithm for feature subset selection. Feature subsets are evaluated by computing the generalization accuracy of (and optionally cost of features used in) the neural network classifier constructed using a computationally efficient neural network learning algorithm called DistAI (Yang et al., 1998b).

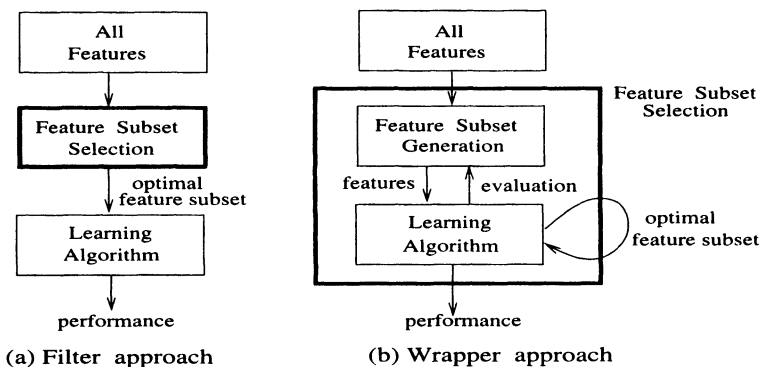


Figure 8.1 Two approaches to feature subset selection based on the incorporation of the learning algorithm.

8.3 FEATURE SELECTION USING A GENETIC ALGORITHM FOR NEURAL NETWORK PATTERN CLASSIFIERS

Feature subset selection in the context of many practical problems (e.g., diagnosis) presents an instance of a multi-criteria optimization problem. The multiple criteria to be optimized include the accuracy of classification, cost and risk associated with classification which in turn depends on the selection of features used to describe the patterns. Genetic algorithms offer a particularly attractive approach for multi-criteria optimization.

Neural networks offer an attractive framework for the design of trainable pattern classifiers for real-world real-time pattern classification tasks on account of their potential for parallelism and fault and noise tolerance (Honavar, 1998a).

While genetic algorithms are generally quite effective for rapid global search of large search spaces in difficult optimization problems, neural networks offer a particularly attractive approach to finetuning promising solutions once they have been identified. Thus, it is attractive to explore combinations of global and local search techniques in the solution of difficult design or optimization problems (Mitchell, 1996). Against this background, the use of genetic algorithms for feature subset selection in the design of neural network pattern classifiers is clearly of interest.

This chapter explores GADistAl, a wrapper-based multi-criteria approach to feature subset selection using a genetic algorithm in conjunction with a relatively fast inter-pattern distance-based neural network learning algorithm called DistAl. However, the general approach can be used with any inductive learning algorithm. The interested reader is referred to (Langley, 1995; Mitchell, 1997; Honavar, 1998a) for surveys of different approaches to inductive learning.

8.3.1 Genetic Algorithms

Evolutionary algorithms include a class related to randomized, population-based heuristic search techniques which include genetic algorithms, genetic programming, evolutionary programming, and variety of related approaches (Mitchell, 1996). They are inspired by processes that are modeled after biological evolution. Central to such evolutionary systems is the idea of a population of potential solutions (individuals) that corresponds to members of a high-dimensional search space.

The individuals represent candidate solutions to the optimization problem being solved. A wide range of genetic representations (e.g., bit vectors, LISP programs, matrices, etc.) can be used to encode the individuals depending on the space of solutions that needs to be searched. In genetic algorithms the individuals are typically represented by n -bit binary vectors. The resulting search space corresponds to an n -dimensional boolean space. In the feature subset selection problem, each individual would represent a feature subset.

It is assumed that the quality of each candidate solution (or fitness of the individual in the population) can be evaluated using a fitness function. In the feature subset selection problem, the fitness function would evaluate the

selected features with respect to some criteria of interest (e.g., cost of the resulting classifier, classification accuracy of the classifier, etc.).

Evolutionary algorithms use some form of fitness-dependent probabilistic selection of individuals from the current population to produce individuals for the next generation. A variety of selection techniques have been explored in the literature. Some of the most common ones are *fitness-proportionate*, *rank-based*, and *tournament-based* selection (Mitchell, 1996). The selected individuals are subjected to the action of genetic operators to obtain new individuals that constitute the next generation. The genetic operators are usually designed to exploit the known properties of the genetic representation, the search space, and the optimization problem to be solved. Genetic operators enable the algorithm to *explore* the space of candidate solutions.

Mutation and crossover are two of the most common operators used with genetic algorithms that represent individuals as binary strings. Mutation operates on a single string and generally changes a bit at random. Thus, a string 11010 may, as a consequence of random mutation, get changed to 11110. Crossover, on the other hand, operates on two parent strings to produce two offspring. With a randomly chosen crossover position 4, the two strings 01101 and 11000 yield the offspring 01100 and 11001 as a result of crossover. Other genetic representations require the use of appropriately designed genetic operators.

The process of fitness-dependent selection and application of genetic operators to generate successive generations of individuals is repeated many times until a satisfactory solution is found (or the search fails). It can be shown that evolutionary algorithms of the sort outlined above simulate highly opportunistic and exploitative randomized search that explores high-dimensional search spaces rather effectively under certain conditions. In practice, the performance of evolutionary algorithms depends on a number of factors including: the choice of genetic representation and operators, the fitness function, the details of the fitness-dependent selection procedure, and the various user-determined parameters such as population size, probability of application of different genetic operators, etc. The specific choices made in the experiments reported in this chapter are summarized in Section 8.4.

8.3.2 Neural Networks

Neural networks are densely connected, massively parallel, shallowly serial networks of relatively simple computing elements or neurons (Ripley, 1996; Mitchell, 1997; Honavar, 1998a). Each neuron computes a relatively simple function of its inputs and transmits outputs to other neurons to which it is connected via its output links. A variety of neuron functions are used in practice. Each neuron has associated with it a set of parameters which are modifiable through learning. The most commonly used parameters are the so-called *weights*.

The computational capabilities (and hence pattern classification abilities) of a neural network depend on its architecture (connectivity), functions computed

by the individual neurons, and the setting of parameters or weights used. It is well-known that multi-layer networks of non-linear computing elements (e.g., threshold neurons) can realize any classification function.

Since the function computed by a neural network is determined by its topology as well as the computations performed by individual neurons, designing a neural network for a particular pattern classification task reduces to determination of the network architecture (number of neurons, their connectivity, etc.), the types of neurons (e.g., linear, sigmoid, threshold, etc.), as well as the parameter or weight values. This is typically accomplished through a combination of design (using a-priori knowledge or guesswork) and inductive learning (which may be used to modify, among other things, the weights, network architecture, or both) (Parekh et al., 1997; Honavar, 1998b).

8.3.3 *Genetic Algorithm Wrapper Approach to Feature Subset Selection for Neural Network Pattern Classifiers: Some Practical Considerations*

Genetic algorithms offer an attractive technique for feature subset selection for neural network pattern classifiers for several reasons, some of which were mentioned above. However, we are faced with several difficulties in using this approach in practice.

Traditional neural network learning algorithms perform an error gradient guided search for a suitable setting of weights in the weight space determined by a user-specified network architecture. This ad hoc choice of network architecture often inappropriately constrains the search for an appropriate setting of weights. For example, if the network has fewer neurons than necessary, the learning algorithm will fail to find the desired classification function. If the network has far more neurons than necessary, it can result in overfitting of the training data leading to poor generalization. In either case, it would make it difficult to evaluate the usefulness of a feature subset employed to describe (or represent) the training patterns used to train the neural network.

Gradient based learning algorithms although mathematically well-founded for unimodal search spaces, can get caught in local minima of the error function. This can complicate the evaluation of a feature subset employed to represent the training patterns used to train the neural networks. This is due to the fact that the poor performance of the classifier might be due to the failure of the learning algorithm, and not the feature subset used.

Fortunately, constructive neural network learning algorithms (Honavar, 1998b) eliminate the need for ad hoc, and often inappropriate a-priori choices of network architectures; and can potentially discover near-minimal networks whose size is commensurate with the complexity of the classification task that is implicitly specified by the training data. Several new, provably convergent, and relatively efficient constructive learning algorithms for multi-category real as well as discrete valued pattern classification tasks have begun to appear in the literature (Parekh et al., 1997; Yang et al., 1998b). Many of these algorithms have demonstrated very good performance in terms of reduced network size, learning time, and generalization in a number of experiments with both

artificial and fairly large real-world datasets (Parekh et al., 1997; Yang et al., 1998b). However, most of them, with the exception of DistAl (Yang et al., 1998b) use time-consuming iterative training algorithms for setting the weights of the neurons.

Using genetic algorithms for feature subset selection for the design of neural network pattern classifiers involves running a genetic algorithm for several generations. In each generation, evaluation of an individual (a feature subset) requires training the corresponding neural network and computing its accuracy and cost. This evaluation has to be performed for each of the individuals in the population. Thus, it is not feasible to use computationally expensive iterative weight update algorithms for training neural network classifiers for evaluating candidate feature subsets. Against this background, DistAl offers an attractive approach to training neural networks.

8.3.4 DistAl: A Fast Algorithm for Constructing Neural Network Pattern Classifiers

DistAl (Yang et al., 1998b) is a simple and relatively fast constructive neural network learning algorithm for pattern classification. The results presented in this chapter are based experiments using neural networks constructed by DistAl. The key idea behind DistAl is to add *hyperspherical* hidden neurons one at a time based on a greedy strategy which ensures that the hidden neuron correctly classifies a maximal subset of training patterns belonging to a single class. Correctly classified examples can then be eliminated from further consideration. The process terminates when the pattern set becomes empty (that is, when the network correctly classifies the entire training set). When this happens, the training set becomes linearly separable in the transformed space defined by the hidden neurons. In fact, it is possible to set the weights on the hidden to output neuron connections without going through an iterative process. It is straightforward to show that DistAl is guaranteed to converge to 100% classification accuracy on any finite training set in time that is polynomial in the number of training patterns (Yang et al., 1998b). Experiments reported in (Yang et al., 1998b) show that DistAl, despite its simplicity, yields classifiers that compare quite favorably with those generated using more sophisticated (and substantially more computationally demanding) learning algorithms. This makes DistAl an attractive choice for experimenting with evolutionary approaches to feature subset selection for neural network pattern classifiers. Key steps in our approach are shown in Figure 8.2.

8.4 IMPLEMENTATION DETAILS

Our experiments were run using a genetic algorithm using rank-based selection strategy. The probability of selection of the highest ranked individual is p (where $0.5 < p < 1.0$ is a user-specified parameter), that of the second highest ranked individual is $p(1 - p)$, that of the third highest ranked individual is $p(1 - p)^2, \dots$, that of the last ranked individual is $1 - (\text{sum of the probabilities of$

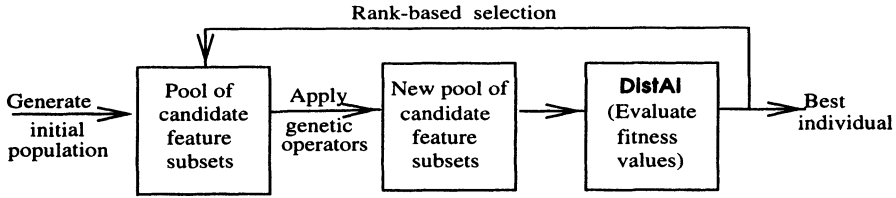


Figure 8.2 GADistAl: Feature subset selection using a genetic algorithm with DistAl.

selection of all the other individuals). The rank-based selection strategy gives a non-zero probability of selection of each individual (Mitchell, 1996). Our experiments used the following parameter settings: Population size is 50; Number of generation is 20; Probability of crossover is 0.6; Probability of mutation is 0.001; Probability of selection of the highest ranked individual is 0.6. The parameter settings were based on results of several preliminary runs. They are comparable to the typical values mentioned in the literature (Mitchell, 1996).

Each individual in the population represents a candidate solution to the feature subset selection problem. Let m be the total number of features available to choose from to represent the patterns to be classified. It is represented by a binary vector of dimension m (where m is the total number of features). If a bit is a 1, it means that the corresponding feature is selected. A value of 0 indicates that the corresponding feature is not selected. The fitness of an individual is determined by evaluating the neural network constructed by DistAl using a training set whose patterns are represented using only the selected subset of features. If an individual has n bits turned on, the corresponding neural network has n input nodes.

The fitness function has to combine two different criteria – the accuracy of the classification function realized by the neural network and the cost of performing classification. The accuracy of the classification function can be estimated by calculating the percentage of patterns in a test set that are correctly classified by the neural network in question. A number of different measures of the cost of classification suggest themselves: cost of measuring the value of a particular feature needed for classification (or the cost of performing the necessary test in a medical diagnosis application), the risk involved, etc. To keep things simple, we chose a 2-criteria fitness function defined as follows:

$$fitness(x) = accuracy(x) - \frac{cost(x)}{accuracy(x) + 1} + cost_{max} \quad (8.1)$$

where $fitness(x)$ is the fitness of the feature subset represented by x , $accuracy(x)$ is the test accuracy of the neural network classifier trained using DistAl using the feature subset represented by x , $cost(x)$ is the sum of measurement costs of feature subset represented by x , and $cost_{max}$ is an upper bound on the costs of candidate solutions. In this case, it is simply the sum of the costs associated with all of the features. This is clearly a somewhat ad hoc choice. However, it does discourage trivial solutions (e.g., a zero

cost solution with a very low accuracy) from being selected over reasonable solutions which yield high accuracy at a moderate cost. It also ensures that $\forall x \ 0 \leq fitness(x) \leq (100 + cost_{max})$. In practice, defining suitable tradeoffs between the multiple objectives has to be based on knowledge of the domain. In general, it is a non-trivial task to combine multiple optimization criteria into a single fitness function. A wide variety of approaches have been examined in the utility theory literature (Keeney and Raiffa, 1976).

8.5 EXPERIMENTS

8.5.1 Description of Datasets

The experiments reported here used a wide range of real-world datasets from the machine learning data repository at the University of California at Irvine (Murphy and Aha, 1994) as well as a carefully constructed artificial dataset (3-bit parity) to explore the feasibility of using genetic algorithms for feature subset selection for neural network classifiers. The feature subset selection using DistAI is also applied to document classification problem for journal paper abstracts and news articles.

3-bit Parity Dataset (3P). This dataset was constructed to explore the effectiveness of the genetic algorithm in selecting an appropriate subset of relevant features in the presence of redundant features so as to minimize the cost and maximize the accuracy of the resulting neural network pattern classifier. The modified training set is constructed as follows: The original features are replicated once (to introduce redundancy) thereby doubling the number of features. Then an additional set of irrelevant features are generated and are assigned random boolean values. 100 7-bit random vectors were generated and augmented with the 6-bit vectors (corresponding to the original 3 bits plus an identical set of 3 bits). Each feature in the resulting dataset is assigned a random cost between 0 and 9. The performance considering the random costs in addition to the accuracy (see Equation 8.1) was compared with that obtained by considering the accuracy alone.

Datasets from UCI Repository. In our experiments with real world datasets, our objective was to compare the neural networks built using feature subsets selected by the genetic algorithm with those that use the entire set of features. Some medical datasets include measurement costs for the features, but most of the datasets lack this information. Therefore, our experiments with the datasets from UCI repository focused on identifying a minimal subset of features that yield high accuracy neural network classifiers. Where measurement costs were available, the performance considering the cost in addition to the accuracy was compared with that obtained by considering the accuracy alone.

Document Datasets. The paper abstracts were chosen from three different sources: IEEE Expert magazine, Journal of Artificial Intelligence Research and

Neural Computation. The news articles were obtained from Reuters dataset. Each document is represented in the form of a vector of numeric weights for each of the words (terms) in the vocabulary. The weights correspond to the *term frequency-inverse document frequency* (*TFIDF*) (Salton, 1989) values for the corresponding words. The training sets for paper abstracts were generated based on the classification of the corresponding documents into two classes (interesting and not interesting) by two different individuals, resulting in two different data sets (**Abstract1** and **Abstract2**). The classifications for news articles were given based on their topics (6, 4 and 8 classes) following (Koller and Sahami, 1997), resulting in three different datasets (**Reuters1**, **Reuters2** and **Reuters3**). Since these datasets do not have measurement costs for the features, our experiments with document datasets also focused on identifying a minimal subset of features that yield high accuracy neural network classifiers.

8.5.2 Experimental Results

Two different sets of experiments were run to explore the performance of GADistAI. The first set of experiments were designed to explore the effect of feature subset selection on the performance of DistAI on a given choice of training and test sets. Each dataset was randomly partitioned into a training and test set (with 90% of the data used for training and the remaining 10% for testing). The genetic algorithm was used to select the best feature subset on the basis of this choice of training and test sets. The results were averaged over 5 independent runs of the genetic algorithm, for a given choice of training and test set. This process was repeated 10 times with 10 different choices of training and test set. The results of these experiments (which represent $5 \times 10 = 50$ runs of the genetic algorithm) are shown in Table 8.1 (GADistAI(rand)) and 8.3.

The second set of experiments explored a somewhat different, but related question. Since feature subset selection in GADistAI is guided by the fitness function, it seems reasonable to expect that the quality of fitness estimates will have some impact on the performance of DistAI. Thus, it is interesting to explore the performance of GADistAI when the fitness estimates are obtained using several training and test sets. Thus, in this set of experiments, fitness estimates used by GADistAI were obtained by averaging the observed fitness values for 10 different partitions of the data into training and test sets (i.e., 10-fold cross validation). The results are shown in Table 8.1, 8.2 and 8.4. The results in Table 8.1 represent averages over 5 independent runs (GADistAI(avg)) and the best run among the 5 runs (GADitAI(best)) of the algorithm.

Improvement in Generalization using Feature Subset Selection. To study the effect of feature subset selection on generalization, experiments were run using classification accuracy as the fitness function. The results in Table 8.1 indicate that the networks constructed using GA-selected subset of features (GADistAI (rand)) compare quite favorably with networks that use all of the features in all randomly partitioned datasets. In particular, feature subset selection resulted in substantial improvement in generalization on many of the

Table 8.1 Comparison of neural network pattern classifiers constructed by DistAI using the entire set of features with the best network constructed by GADistAI.^a

<i>Dataset</i>	DistAI ^c		GADistAI (rand) ^d		GADistAI (avg) ^e		GADistAI (best) ^f	
	<i>Dim</i>	<i>Acc</i>	<i>Dim</i>	<i>Acc</i>	<i>Dim</i>	<i>Acc</i>	<i>Dim</i>	<i>Acc</i>
3P	13	79.0	6.6	100	4.8	100	4	100
Annealing	38	96.6	21.0	99.5	20.0	98.8	18	99.5
Audiology	69	66.0	36.4	83.5	37.2	72.6	39	76.5
Bridges	11	63.0	5.6	81.6	4.9	56.9	5	67.0
Cancer	9	97.8	5.4	99.3	6.0	98.0	8	98.6
CRX	15	87.7	8.0	91.5	7.4	87.7	6	88.0
Flag	28	65.8	14.0	78.1	14.2	63.9	18	70.0
Glass	9	70.5	5.5	80.8	4.4	69.3	5	71.0
Heart	13	86.7	7.2	93.9	7.6	85.5	7	85.9
HeartCle^b	13	85.3	7.3	92.9	8.4	86.9	9	87.7
HeartHun^b	13	85.9	7.0	93.0	7.4	85.4	8	87.2
HeartLB^b	13	80.0	7.1	91.0	7.6	79.8	6	83.0
HeartSwi^b	13	94.2	6.6	98.3	7.4	95.3	8	96.7
Hepatitis	19	84.7	9.2	97.1	10.2	85.2	10	88.7
Horse	22	86.0	11.1	92.6	9.6	83.2	5	85.0
Ionosphere	34	94.3	17.3	98.6	16.6	94.5	13	96.0
Pima	8	76.3	3.8	79.5	4.0	73.1	2	76.8
Promoters	57	88.0	28.8	100	30.6	89.8	31	92.0
Sonar	60	83.0	30.7	97.2	32.2	84.0	28	85.5
Soybean	35	81.0	19.4	92.8	21.0	83.1	19	84.3
Vehicle	18	65.4	9.1	68.8	9.4	50.1	11	59.4
Votes	16	96.1	8.9	98.8	8.2	97.0	7	97.9
Vowel	10	69.8	6.5	78.4	6.8	70.2	6	71.5
Wine	13	97.1	6.7	99.4	8.2	96.7	7	97.1
Zoo	16	96.0	9.3	100	8.8	96.8	9	99.0
Abstract1	790	89.0	393.7	97.6	402.2	89.2	387	91.0
Abstract2	790	84.0	393.8	94.4	389.8	84.0	382	85.0
Reuters1	1568	91.6	786.1	94.9	766.0	90.2	750	91.5
Reuters2	435	88.5	218.3	97.5	222.4	90.3	195	91.5
Reuters3	1440	96.4	715.4	98.7	721.0	96.2	712	96.9

^a *Dim* is the number of features and *Acc* is the generalization accuracy.^b {Cleveland,Hungarian,Long Beach,Switzerland} heart disease.^c The standard deviation of *Acc* is less than 12.2.^d The standard deviation of *Dim* is less than 3.7 for most of the datasets and less than 20.3 for document datasets. The standard deviation of *Acc* is less than 5.0 for almost all datasets.^e The standard deviation of *Dim* is less than 3.0 for most of the datasets and less than 16.6 for document datasets. The standard deviation of *Acc* is less than 3.0 for almost all datasets.^f The standard deviation of *Acc* is less than 13.8.

datasets. (For example, 100% accuracy was obtained with **3P**, **Promoters**,

and **Zoo** datasets). Also, the number of features selected is significantly smaller than the total number of features present in the original data representation.

Table 8.2 Comparison between various approaches for feature subset selection.⁹

<i>Dataset</i>	non-GA		ADHOC		GADistAI	
	<i>Dim</i>	<i>Acc</i>	<i>Dim</i>	<i>Acc</i>	<i>Dim</i>	<i>Acc</i>
Annealing	-	-	8	95.0	18	99.5
Cancer	4	74.7	-	-	8	98.6
CRX	6	85.0	7	85.1	6	88.0
Glass	4	62.5	4	70.5	5	71.0
Heart	3	79.2	5	80.8	7	85.9
Hepatitis	4	84.6	-	-	10	88.7
Horse	4	85.3	-	-	5	85.0
Pima	-	-	3	73.2	2	76.8
Sonar	-	-	16	76.0	28	85.5
Vehicle	-	-	7	69.6	11	59.4
Votes	4	97.0	5	95.7	7	97.9
Reuters1	40	94.1	-	-	750	91.5
Reuters2	40	90.0	-	-	195	91.5
Reuters3	80	98.6	-	-	712	96.9

⁹ A '-' indicates that the result is not reported in the corresponding reference.

The results also indicate that the networks constructed using GA-selected subset of features by average fitness values (GADistAI (avg)) compare favorably with networks that use all of the features in most of the datasets. Clearly, GADistAI outperformed DistAI (with all features) in the parity problem: it successfully selected relevant features that resulted in 100% accuracy. For the remaining datasets, the improvement in generalization ranged from modest in some cases to marginal in others. The best individual generated by GADistAI (GADistAI (best)) outperformed DistAI in almost all datasets. Again, the number of features selected is significantly smaller than the total number of features present in the original data representation in all of the datasets.

Table 8.2 compares the results of GADistAI with the results of other GA-based (ADHOC) (Richeldi and Lanzi, 1996) and the best among several non GA-based approaches (non-GA) that are available in the literature (Liu and Setiono, 1996a; Liu and Setiono, 1996b; Kohavi, 1994; Koller and Sahami, 1996; Koller and Sahami, 1997). The results indicate that GADistAI gave higher generalization accuracy than the other techniques or comparable accuracy in almost all cases (except **Vehicle** dataset) although it occasionally selected more features. GADistAI produced feature subsets with larger number of features than the approach in (Koller and Sahami, 1996; Koller and Sahami, 1997) for **Reuters** datasets. This can be explained by that the former found the feature subsets using a genetic algorithm for datasets with relatively large number of features while the latter set up the number of features to select a-priori.

It should be noted that it is not generally feasible to do a completely fair and thorough comparison between different approaches without the complete knowledge of the parameters and the set up used in the experiments.

Minimizing Cost and Maximizing Accuracy using Feature Subset Selection. The selection was based on both the generalization accuracy and the measurement cost of features. **3P**, **HeartCle**, **Hepatitis** and **Pima** datasets were used for the experiment (with the random costs in the 3-bit parity problem). The results shown in Table 8.3 and 8.4 are obtained by the fitness estimates over randomly partitioned datasets and the average fitness estimates over datasets arranged by 10-fold cross validation, respectively.

In Table 8.3, the fitness function that combined both accuracy and cost outperformed that based on accuracy alone in every respect: the number of features used, generalization accuracy, and cost. This is not surprising because the former tries to minimize cost (while maximizing the accuracy), which reduces the number of features, while the latter emphasizes only on the accuracy.

Table 8.4 also shows the fitness function that combined both accuracy and cost outperforms that based on accuracy alone in all datasets except **HeartCle**. The generalization accuracy was higher and the cost was also higher with the fitness function that is based on accuracy alone in **HeartCle** dataset. This explains how the fitness function (Equation 8.1) works in **GADistAl** and verifies the rationale behind it. Also, note that some of the runs resulted in feature subsets which did not necessarily have minimum cost. This suggests the possibility of improving the results by the use of a more principled choice of a fitness function that combines accuracy and cost.

Table 8.3 Comparison of neural network pattern classifiers constructed by **GADistAl** with different fitness evaluations for randomly partitioned datasets.

<i>Dataset</i>	Accuracy only			Accuracy & Cost		
	<i>Dim</i>	<i>Acc</i>	<i>Cost</i>	<i>Dim</i>	<i>Acc</i>	<i>Cost</i>
3P	6.6	100	46.1	4.3	100	26.7
HeartCle	7.3	92.9	335.7	6.1	93.0	261.5
Hepatitis	9.2	97.1	22.8	8.3	97.3	19.0
Pima	3.8	79.5	28.5	3.1	79.5	22.8

8.6 SUMMARY AND DISCUSSION

An approach to feature subset selection using a genetic algorithm for neural network pattern classifiers is proposed in this chapter. A fast inter-pattern distance-based constructive neural network algorithm, **DistAl**, is employed to evaluate the fitness (in terms of the generalization accuracy) of candidate feature subsets in the genetic algorithm. The results presented in this chapter indicate that genetic algorithms offer an attractive approach to solving the

Table 8.4 Comparison of neural network pattern classifiers constructed by GADistAI with different fitness evaluations for datasets arranged by 10-fold cross validation.

<i>Dataset</i>	Accuracy only			Accuracy & Cost		
	<i>Dim</i>	<i>Acc</i>	<i>Cost</i>	<i>Dim</i>	<i>Acc</i>	<i>Cost</i>
3P	4.8	100	35.6	3.8	100	25.4
HeartCle	8.4	86.9	390.5	7.2	85.7	317.8
Hepatitis	10.2	85.2	23.4	10.0	85.3	23.2
Pima	4.0	73.1	29.3	4.2	76.1	20.8

feature subset selection problem in inductive learning of pattern classifiers in general, and neural network pattern classifiers in particular.

The GA-based approach to feature subset selection does not rely on monotonicity assumptions that are used in traditional approaches to feature selection which often limits their applicability to real-world classification and knowledge acquisition tasks. It also offers a natural approach to feature subset selection by taking into account, the distribution of available data. This is due to the fact that feature selection is driven by estimated fitness values, which if based on multiple partitions of the dataset into training and test data, provide a robust measure of performance of the feature subset. This is not generally the case with many of the greedy stepwise algorithms that select features based on a single partition of the data into training and test sets. Consequently, the feature subsets selected by such algorithms are likely to perform rather poorly on other random partitions of the data into training and test sets.

The approach to feature subset selection is able to naturally incorporate multiple criteria (e.g., accuracy, cost) into the feature selection process. This finds applications in cost-sensitive design of classifiers for tasks such as medical diagnosis, computer vision, among others. Another interesting application is automated data mining and knowledge discovery from datasets with an abundance of irrelevant or redundant features. In such cases, identifying a relevant subset that adequately captures the regularities in the data can be particularly useful, particularly in scientific knowledge discovery tasks. Techniques similar to the one discussed in this chapter have been successfully used recently to select feature subsets for pattern classification tasks that arise in power system security assessment (Zhou et al., 1997), sensor subsets in the design of behavior and control structures for autonomous mobile robots (Balakrishnan and Honavar, 1996).

Additional experiments with GADistAI in scientific knowledge discovery tasks in bioinformatics (e.g., discovery of protein structure–function relationships, carcinogenicity prediction, gene sequence identification) are currently in progress. Some directions for future research include: Extension of feature subset selection by incorporating *feature construction* and *genetic programming* (Koza, 1992); Extensive experimental (and wherever feasible, theoretical) comparison of the performance of the proposed approach with that of conventional

methods for feature subset selection; More principled design of multi-objective fitness functions for feature subset selection using domain knowledge as well as mathematically well-founded tools of multi-attribute utility theory (Keeney and Raiffa, 1976).

Acknowledgments

This research was partially supported by National Science Foundation Grant IRI-9409580 and John Deere Foundation Grant to Vasant Honavar. The authors wish to thank Mehran Sahami for providing **Reuters** document datasets. The authors are grateful to Dr. Pazzani of the Department of Information and Computer Science at the University of California at Irvine for managing the repository of machine learning datasets and making it available to us. An earlier version of this chapter appears in IEEE Expert. ©1998 IEEE.

References

- Almuallim, H. and Dietterich, T. (1994). Learning boolean concepts in the presence of many irrelevant features. *Artificial Intelligence*, 69(1-2):279–305.
- Balakrishnan, K. and Honavar, V. (1996). On sensor evolution in robotics. In Koza, Goldberg, Fogel, and Riolo, editors, *Proceedings of the 1996 Genetic Programming Conference – GP-96*, pages 455–460. MIT Press, Cambridge, MA.
- Brill, F., Brown, D., and Martin, W. (1992). Fast genetic selection of features for neural network classifiers. *IEEE Transactions on Neural Networks*, 3(2):324–328.
- Caruana, R. and Freitag, D. (1994). Greedy attribute selection. In *Proceedings of the Eleventh International Conference on Machine Learning*, pages 28–36, New Brunswick, NJ. Morgan Kaufmann.
- Cost, S. and Salzberg, S. (1993). A weighted nearest neighbor algorithm for learning with symbolic features. *Machine Learning*, 10(1):57–78.
- Cover, T. and Hart, P. (1967). Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13:21–27.
- Dasarathy, B. (1991). *Nearest Neighbor (NN) Norms: NN Pattern Classification Techniques*. IEEE Computer Society Press, Los Alamitos, CA.
- Dash, M. and Liu, H. (1997). Feature selection for classification. *Intelligent Data Analysis*, 1(3).
- Doak, J. (1992). An evaluation of feature selection methods and their application to computer security. Technical Report CSE-92-18, Department of Computer Science, University of California, Davis, CA.
- Foroutan, I. and Sklansky, J. (1987). Feature selection for automatic classification of non-gaussian data. *IEEE Transactions on Systems, Man and Cybernetics*, 17:187–198.
- Guo, Z. (1992). *Nuclear Power Plant Fault Diagnostics and Thermal Performance Studies Using Neural Networks and Genetic Algorithms*. PhD thesis, University of Tennessee, Knoxville, TN.

- Guo, Z. and Uhrig, R. (1992). Using genetic algorithms to select inputs for neural networks. In *Proceedings of COGANN'92*, pages 223–234.
- Honavar, V. (1998a). Machine learning: Principles and applications. In Webster, J., editor, *Encyclopedia of Electrical and Electronics Engineering*. Wiley, New York. To appear.
- Honavar, V. (1998b). Structural learning. In Webster, J., editor, *Encyclopedia of Electrical and Electronics Engineering*. Wiley, New York. To appear.
- John, G., Kohavi, R., and Pfleger, K. (1994). Irrelevant features and the subset selection problem. In *Proceedings of the Eleventh International Conference on Machine Learning*, pages 121–129, New Brunswick, NJ. Morgan Kaufmann.
- Keeney, R. and Raiffa, H. (1976). *Decisions with Multiple Objectives: Preferences and Value Tradeoffs*. Wiley, New York.
- Kira, K. and Rendell, L. (1992). A practical approach to feature selection. In *Proceedings of the Ninth International Conference on Machine Learning*, pages 249–256. Morgan Kaufmann.
- Kohavi, R. (1994). Feature subset selection as search with probabilistic estimates. In *AAAI Fall Symposium on Relevance*.
- Koller, D. and Sahami, M. (1996). Toward optimal feature selection. In *Machine Learning: Proceedings of the Thirteenth International Conference*. Morgan Kaufmann.
- Koller, D. and Sahami, M. (1997). Hierarchically classifying documents using very few words. In *International Conference on Machine Learning*, pages 170–178.
- Kononenko, I. (1994). Estimating attributes: Analysis and extension of relief. In *Proceedings of European Conference on Machine Learning*, pages 171–182.
- Koza, J. (1992). *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, Cambridge, MA.
- Langley, P. (1994). Selection of relevant features in machine learning. In *Proceedings of the AAAI Fall Symposium on Relevance*, pages 1–5, New Orleans, LA. AAAI Press.
- Langley, P. (1995). *Elements of Machine Learning*. Morgan Kaufmann, Palo Alto, CA.
- Liu, H. and Setiono, R. (1996a). Feature selection and classification - a probabilistic wrapper approach. In *Proceedings of the Ninth International Conference on Industrial and Engineering Applications of AI and ES*.
- Liu, H. and Setiono, R. (1996b). A probabilistic approach to feature selection - a filter solution. In *Proceedings of the Thirteenth International Conference on Machine Learning*. Morgan Kaufmann.
- Mitchell, M. (1996). *An Introduction to Genetic Algorithms*. MIT Press, Cambridge, MA.
- Mitchell, T. (1997). *Machine Learning*. McGraw Hill, New York.
- Modrzejewski, M. (1993). Feature selection using rough sets theory. In *Proceedings of the European Conference on Machine Learning*, pages 213–226. Springer.

- Murphy, P. and Aha, D. (1994). Repository of machine learning databases. Department of Information and Computer Science, University of California, Irvine, CA.
- Narendra, P. and Fukunaga, K. (1977). A branch and bound algorithm for feature subset selection. *IEEE Transactions on Computers*, 26:917–922.
- Parekh, R., Yang, J., and Honavar, V. (1997). Constructive neural network learning algorithms for multi-category real-valued pattern classification. Technical Report ISU-CS-TR97-06, Department of Computer Science, Iowa State University.
- Richeldi, M. and Lanzi, P. (1996). Performing effective feature selection by investigating the deep structure of the data. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, pages 379–383. AAAI Press.
- Ripley, B. (1996). *Pattern Recognition and Neural Networks*. Cambridge University Press, New York.
- Salton, G. (1989). *Automatic Text Processing: the Transformation, Analysis, and Retrieval of Information by Computer*. Addison-Wesley, Reading, Massachusetts.
- Sheinvald, J., Dom, B., and Niblack, W. (1990). A modelling approach to feature selection. In *Proceedings of the Tenth International Conference on Pattern Recognition*, pages 535–539.
- Siedlecki, W. and Sklansky, J. (1988). On automatic feature selection. *International Journal of Pattern Recognition*, 2:197–220.
- Siedlecki, W. and Sklansky, J. (1989). A note on genetic algorithms for large-scale feature selection. *IEEE Transactions on Computers*, 10:335–347.
- Skalak, D. (1994). Prototype and feature selection by sampling and random mutation hill-climbing algorithms. In *Proceedings of the Eleventh International Conference on Machine Learning*, pages 293–301, New Brunswick, NJ. Morgan Kaufmann.
- Vafaie, H. and De Jong, K. (1993). Robust feature selection algorithms. In *Proceedings of the IEEE International Conference on Tools with Artificial Intelligence*, pages 356–363.
- Wettschereck, D., Aha, D., and Mohri, T. (1995). A review and empirical evaluation of feature weighting methods for a class of lazy learning algorithms. Technical Report AIC95-012, Naval Research Laboratory, Navy Center for Applied Research in Artificial Intelligence, Washington, D.C.
- Yang, J., Pai, P., Honavar, V., and Miller, L. (1998a). Mobile intelligent agents for document classification and retrieval: A machine learning approach. In *14th European Meeting on Cybernetics and Systems Research. Symposium on Agent Theory to Agent Implementation*, Vienna, Austria.
- Yang, J., Parekh, R., and Honavar, V. (1998b). DistAl: An inter-pattern distance-based constructive learning algorithm. In *Proceedings of the International Joint Conference on Neural Networks*, Anchorage, Alaska. To appear.

Zhou, G., McCalley, J., and Honavar, V. (1997). Power system security margin prediction using radial basis function networks. In *Proceedings of the 29th Annual North American Power Symposium*, Laramie, Wyoming.