

An Entropy-based Adaptive Genetic Algorithm for Learning Classification Rules

Yang LI Su

Dept. of Computer Science Georgia A&M University
lyyang@cs.tamu.edu

Widyantoro H . Bint

Dept. of Computer Science Georgia A&M University
dhw7942@cs.tamu.edu

Ioerger Beoston

Dept. of Computer Science Georgia A&M University
ioerger@cs.tamu.edu

Yen Si

Dept. of Computer Science Georgia A&M University
yen@cs.tamu.edu

Abstract

Genetic algorithm is one of the commonly used approaches on data mining. In this paper, we put forward a genetic algorithm approach for classification problems. Binary coding is adopted in which an individual in a population consists of a fixed number of rules that stand for a solution candidate. The evaluation function considers four important factors which are error rate, entropy measure, rule consistency and hole ratio, respectively. Adaptive asymmetric mutation is applied by the self-adaptation of mutation inversion probability from 1-0 (0-1). The generated rules are not disjoint but can overlap. The final conclusion for prediction is based on the voting of rules and the classifier gives all rules equal weight for their votes. Based on three databases, we compared our approach with several other traditional data mining techniques including decision trees, neural networks and naive bayes learning. The results show that our approach outperformed others on both the prediction accuracy and the standard deviation.

Keywords: genetic algorithm, adaptive asymmetric mutation, entropy, voting-based classifier

1. Introduction

Genetic algorithms have been successfully applied to a wide range of optimization problems including design, scheduling, routing, and control, etc. Data mining is also one of the important application fields of genetic algorithms. In data mining, GA can be used to either optimize parameters for other kinds of data mining algorithms or discover knowledge by itself. In this latter task the rules that GA found are usually more general because of its global search nature. In contrast, most other data mining methods are based on the rule induction paradigm, where the algorithm usually performs a kind of local search. The advantage of GA becomes more obvious when the search space of a task is large.

In this paper we put forward a genetic algorithm approach for classification problems. First we use binary coding in which an individual solution candidate consists of a fixed number of rules. In each rule,

k bits are used for the possible k values of a certain attribute. Continuous attributes are modified to threshold-based boolean attributes before coding. Rule consequent is not explicitly coded in the string, instead, the consequent of a rule is determined by the majority of training examples it matches.

Four important factors are considered in our evaluation functions. Error rate is calculated by the predicting results on the training examples. Entropy is used to measure the homogeneity of the examples that a certain rule matches. Rule consistency is a measure on the consistency of classification conclusions of a certain training example given by a set of rules. Finally, hole ratio is used to evaluate the percentage of training examples that a set of rules does not cover. We try to include related information as complete as possible in the evaluation function so that the overall performance of a rule set can be better.

An adaptive asymmetric mutation operator is applied in our reproduction step. When a bit is selected to mutate, the inversion probability from 1-0 (0-1) is not 50% as usual. The value of this probability is asymmetric and self-adaptive during the running of the program. This is made to reach the best match of a certain rule to training examples. For crossover, two-point crossover is adopted in our approach.

We used three real databases to test our approach: credit database, voting database and heart database. We compared our performance with four well-known methods from data mining, namely Induction Decision Trees (ID3) (Quinlan, 1986), ID3 with Boosting (Quinlan, 1996), Neural Networks, and Naive Bayes (Mitchell, 1997). The appropriate state-of-the-art techniques are incorporated in these non-GA methods to improve their performances. The results show that our GA approach outperformed other approaches on both the prediction accuracy and the standard deviation.

In the rest of the paper, we will first give a brief overview of related work. Our GA approach is then discussed in detail. The results of our application on three real databases and the comparison with other data mining methods are followed. Finally we make some concluding remarks.

2. Related works

Many researchers have contributed to the application of GA on data mining. In this section, we will give a brief overview on a few representative works.

In early 90's, De Jong et al. implemented a GA-based system called GABIL that continually learns and refines concept classification rules (De Jong, 1993). An individual is a variable-length string representing a set of fixed-length rules. Traditional bit inversion is used on mutation. In crossover, corresponding crossover points in the two parents must semantically match. The fitness function contains only the percent of correct examples classified by an individual rule set. They compared the performance of GABIL with that of four other traditional concept learners on a variety of target concepts.

In GIL (Janikow, 1993), an individual is also a set of rules, but attributes values are encoded directly, rather than bits. GIL has special genetic operators for handling rule sets, rules, and rule conditions. The operators can perform generalization, specialization or other operations. Besides the correctness, the evaluation function of GIL also includes the complexity of a rule set. Therefore it favors correct, simple (short) rules.

Greene and Smith put forward a GA-based inductive system called COGIN (Greene, 1993). In contrast to the above two approaches, the system's current model at any point during the search is represented as a population of fixed length rules. The population size (i.e., the number of rules in the model) will vary from cycle to cycle as a function of the coverage constraint is applied. The fitness function contains the information gain of a rule R and a penalty of the number of misclassifications made by R . Entropy measure is used to calculate the information gain of rule R based on the number of

examples rule R *Matched* and *Unmatched*. However, it couldn't evaluate the entropy measure of the entire partition formed by the classification due to its encoding method.

In view of the situation that most of the data mining work emphasizes only on the predictive accuracy and comprehensibility, Noda et al. (Noda, 1999) put forward a GA approach designed to discover the interesting rules. The fitness function consists of two parts. The first one measures the degree of interestingness of the rule, while the second measures its predictive accuracy. The computation of the consequent's degree of interestingness is based on the following idea: the larger the relative frequency (in the training set) of the value being predicted by the consequent, the less interesting it is. In other words, the rarer a value of a goal attribute, the more interesting a rule predicting it is. Since the values of the goal attribute in the databases we tested are not unevenly distributed seriously, we didn't specially consider interestingness of a rule in our current implementation but mainly focus on including related factors as complete as possible to improve the predicting accuracy. However, we did consider how to treat uneven distribution of goal attribute values somehow. We will discuss this in details in the next section.

3. Our GA approach

In this section we present our GA approach for classification problem. The key idea of the algorithm is general and should be applicable for various kinds of classification problems. Some parameter values used in the algorithm might be task dependent.

3.1 Individual's encoding

Each individual in the population consists of a fixed number of rules. In other words, the individual itself is a complete solution candidate. In our current implementation, we set this fixed number as 10 which well satisfies the requirement of our testing databases. The antecedent of a certain rule in the individual is formed by a conjunction of n attributes, where n is number of attributes being mined. K bits will be used to stand for an attribute if this attribute has k possible values. Continuous attributes will be partitioned to threshold-based boolean attribute in which the threshold is a boundary (adjacent examples across this boundary differ in their target classification) that maximizes the information gain. Therefore two bits will be used for a continuous attribute. The consequent of a rule is not explicitly encoded in the string. In contrast, it is automatically given based on the proportion of positive/negative examples it matches in the training set. We will illustrate the encoding method by the following example.

Suppose our task has three attributes, and they have 4, 2, 5 possible values respectively. Then an individual in the population can be represented as following:

A1	A2	A3	A1	A2	A3	A1	A2	A3
0110	11	10110	1110	01	10011	1100	11	01110
Rule 1			Rule 2			Rule 10		

Ten rules are included in this individual. The architecture of each rule is same. We will use rule 1 to explain the meaning of encoding. In this example, the meaning of the antecedent of rule 1 is:

If (A1=value 2 OR value 3) AND (A2=value 1 OR value 2) AND (A3=value 1 OR value 3 OR value 4)

If all the bits belong to one attribute are 0s, it means that attribute can not equal to any possible value therefore this is meaningless. To avoid this, we add one step before the evaluation of the population. We will check each rule in each individual one by one, if the above case happens, we will randomly select one bit of that attribute and change it to one.

The consequent of a rule is not encoded in the string. It will be determined by the proportion situation of the training examples that rule matches. Suppose i is one of the classifications, the consequent of a rule will be i if $\frac{N_{matched_i}}{N_{matched}} > \frac{N_{training_i}}{N_{training}}$, (1)

where $N_{matched_i}$ is the number of examples whose classification is i and matched by that rule, $N_{matched}$ is the total number of examples that the rule matches; $N_{training_i}$ is the number of training examples whose classification is i , $N_{training}$ is the total number of training examples.

For example, if the distribution of the positive examples and negative examples in the training set is 42% and 58%, and among the examples of rule 1 matches positive / negative examples are half to half, then the consequent of rule 1 should be positive because $0.5 > 0.42$. Since the testing databases we use at this time don't have a very uneven distribution on the classification of training examples, in our current implementation we didn't specially consider the interestingness of rules but use this strategy to keep enough rules to match examples with minor classification. Our encoding method is not limited to two-category classification but applicable to multi target value problems.

3.2 Fitness function

It is very important to define a good fitness function that rewards the right kinds of individuals. We try to consider affecting factors as complete as possible to improve the results of classification. Our fitness function is defined as following:

$$\text{Fitness} = \text{Error rate} + \text{Entropy measure} + \text{Rule consistency} + \text{Hole ratio} \quad (2)$$

We will elaborate each part in the fitness function in the following.

1) Error rate

It is well known that accuracy is the most important and commonly used measure in the fitness function as the final goal of data mining is to get good prediction results. Since our objective function here is minimization, we use error rate to represent this information. It is calculated as:

$$\text{Error rate} = \text{percent of misclassified examples} \quad (3)$$

If a rule matches a certain example, the classification it gives is its consequent part. If it doesn't match, no classification is given. An individual consists of a set of rules, the final classification predicted by this rule set is based on the voting of those rules that match the example. The classifier gives all matching rules equal weight. For instance, in an individual (which has ten rules here), one rule doesn't match, six rules give positive classification and three rules give negative classification on a training example, then the final conclusion given by this individual on that training example is positive. If a tie happens (i.e., four positive classifications and four negative classifications), the final conclusion will be the majority classification in the training examples. If none of the rules in the individual matches that example, the final conclusion will also be the majority classification in the training examples. The error rate measure of this individual is the percent of misclassified examples among all training examples.

2) Entropy measure

Entropy is a commonly used measure in information theory. Originally it is used to characterize the (im)purity of an arbitrary collection of examples. In our implementation entropy is used to measure the homogeneity of the examples that a rule matches.

Given a collection S , containing the examples that a certain rule R matches, let P_i be the proportion of examples in S belonging to class i , then the entropy $Entropy(R)$ related to this rule is defined as:

$$Entropy(R) = -\sum_{i=1}^n (p_i \log_2(p_i)) \quad (4)$$

where n is the number of target classifications.

While an individual consists of a number of rules, the entropy measure of an individual is calculated by averaging the entropy of each rule:

$$Entropy(individual) = \frac{\sum_{i=1}^{N_R} Entropy(R_i)}{N_R} \quad (5)$$

where N_R is number of rules in the individual (in our current implementation it is 10).

The rationale of using entropy measure in fitness function is to prefer those rules that match less examples whose target values are different from rule's consequent. High accuracy does not implicitly guarantee the entropy measure is good because the final classification conclusion of a certain training example is based on the comprehensive results of a number of rules. It is very possible that each rule in the individual has a bad entropy measure but the whole rule set still gives the correct classification. Keeping low entropy value of an individual will be helpful to get better predicting results for untrained examples.

3) Rule consistency

As stated in the above sections, the final predicted classification of a training example is the majority classification made by rules in an individual. Let's consider the following classifications made by two individuals on an example:

Individual a: six rules $\rightarrow +$, four rules $\rightarrow -$, final classification: +

Individual b: nine rules $\rightarrow +$, one rule $\rightarrow -$, final classification: +

We will prefer the second individual since it is less ambiguous. To address this rule consistency issue, we add another measure in the fitness function. The calculation is similar to the entropy measure. Let $P_{correct}$ be the proportion of rules in one individual whose consequent is same with the target value of the training example, then

$$Rule_{consistency}(individual) = -p_{correct} \log_2 p_{correct} - (1 - p_{correct}) \log_2 (1 - p_{correct}) \quad (6)$$

We should notice that this formula will give the same rule consistency value when $p_{correct}$ and $(1 - p_{correct})$ switch each other. Therefore a penalty will be given when $p_{correct}$ is smaller than 0.5. In this case $Rule_{consistency} = 2 - Rule_{consistency}$.

The above calculation is based on the predicting results for one training example. The complete measure of rule consistency of an individual should be averaged by the number of training examples.

4) Hole ratio:

The last element in the fitness function is the hole ratio. It is a measure of rule's coverage on training examples. Coverage is not a problem for traditional inductive learning methods like decision trees, since the process of creating trees guarantees that all the training examples will be covered in the tree. However, this also brings a new problem that it may be sensitive to noise. GA approach does not guarantee that the generated rules will cover all the training examples. This allows flexibility and may be potentially useful for future prediction. In real implementation we still hope the coverage should reach a certain point. For instance, if a rule only matches one training example and its consequent is correct, the accuracy and entropy measure of this rule are both excellent but we do not prefer this rule because its coverage is too low.

In our fitness function the hole ratio equals to *1-coverage*, in which the latter is calculated by the union of examples that are matched and also correctly predicted by the rules in an individual. Totally misclassified examples (not classified correctly by any rule in the individual) will not be included even though they are matched by some rules. The following is the formula to calculate the hole ratio for binary classification problem (positive, negative).

$$Hole = 1 - \frac{\left| \bigcup_i P_i^+ \right| + \left| \bigcup_i N_i^- \right|}{|S|} \quad (7)$$

where P_i^+ stands for those examples whose target value is positive and classified as positive by at least one rule in the individual, N_i^- stands for those examples whose target value is negative and classified as negative by at least one rule in the individual. $|S|$ is the total number of training examples.

3.3 Adaptive asymmetric mutation

In our reproduction step, traditional bit inversion is used on mutation. However, we found many examples will not be matched if we keep number of 1's and 0's approximately equivalent in an individual (i.e., the inversion probability from 1-0 and 0-1 are both 50%). The learning process will become a majority guess if there are too many unmatched examples. Therefore, we put forward a strategy of adaptive asymmetric mutation in which the inversion probability from 1-0 (0-1) is self-adaptive during the process of run. The asymmetric mutation biases the population toward generating rules with more coverage on training examples. The self-adaptation of inversion probability makes the optimal mutation parameter be automatically adjusted.

We presented an adaptive simplex genetic algorithm before (Yang, 2000) in which the percentage of simplex operator is self-adaptive during the process of run. Similar idea is adopted here that average of fitness is used as a feedback to adjust the inversion probability. The process of self-adaptation is described as following:

1) An initial inversion probability is set (e.g., 0.5 for 1-0). Use this probability on mutation to produce a new generation. Calculate the average fitness of this generation.

2) Randomly select the direction of changing this probability (increase, decrease). Modify the probability along that direction with a small amount (0.02 in our current implementation). Use the new probability to produce the next generation and calculate the average fitness of the new generation.

3) If the fitness is better (value is smaller), continue on this direction and the amount of change is:

$$\Delta p = \max\{0.05, (1 - \text{fitness}_{\text{new}} / \text{fitness}_{\text{old}}) * 0.1\} \quad (8)$$

If the fitness is worse (value is larger), reverse the direction and the amount of change is:

$$\Delta p = \max\{0.05, (\text{fitness}_{\text{new}} / \text{fitness}_{\text{old}} - 1) * 0.05\} \quad (9)$$

Use the new probability to produce the next generation and calculate the average fitness of the new generation.

Repeat step 3 until the program ends.

4. Results and discussions

We tested our approach on three real databases. We compared our approach with four other traditional data mining techniques. This section will present the testing results.

4.1 The information of databases

1) Credit database

This database concerns credit card applications. All attribute names and values have been changed to meaningless symbols to protect confidentiality of the data. This database is interesting because there is a good mix of attributes ---- continuous, nominal with small numbers of values, and nominal with larger numbers of values. There are 15 attributes plus one target attribute. Total number of instances is 690.

2) Voting database

This database saves 1984 United States Congressional Voting Records. The data set includes votes for each of the U.S. House of Representatives Congressmen on the 16 key votes identified by the CQA. The CQA lists nine different types of votes: voted for, paired for, and announced for (these three simplified to yea), voted against, paired against, and announced against (these three simplified to nay), voted present, voted present to avoid conflict of interest, and did not vote or otherwise make a position known (these three simplified to an unknown disposition). There are 16 attributes plus one target attribute. Total number of instances is 435 (267 democrats, 168 republicans).

3) Heart database

This database concerns heart disease diagnosis. The data was provided by V.A. Medical Center, Long Beach and Cleveland Clinic Foundation: Robert Detrano, M.D., Ph.D. There are 14 attributes plus one target attribute. Total number of instances is 303.

4.2 The description of non-GA approaches

We used four well-known methods from machine learning, namely Induction Decision Trees (ID3) (Quinlan, 1986), Decision Trees with Boosting (Quinlan, 1996), Neural Networks, and Naïve Bayes (Mitchell, 1997), to compare the performance of our improved GA. Appropriate state-of-the-art techniques have been incorporated in most of the non-GA methods to improve their performance. The following is a description of the non-GA approaches we used for the performance comparison studies.

1) Induction Decision Trees

The construction of a decision tree is divided into two stages. First, creating an initial, large decision tree using a set of training set. Second, pruning the initial decision tree, if applicable, using a validation set. Given a noise-free training set, the first stage will generate a decision tree that can classify correctly all examples in the set. Except that the training set covers all instances in the domain, the initial decision tree generated will over fit the training data, which then reduce its performance in the test data. The second stage helps alleviate this problem by reducing the size of the tree. This process has an effect in generalizing the decision tree that hopefully could improve its performance in the test data.

During the construction of an initial decision tree, the selection of the best attribute is based on either the *information gain* (IG) or *gain ratio* (GR). A binary split is applied in nodes with continuous-valued attributes. The best cut-off value of a continuous-valued attribute is locally selected within each node in the tree based on the remaining training examples. The tree's node expansion stops either when the remaining training set is homogeneous (e.g., all instances have the same target attribute values) or when no attribute remains for selection. The decision on a leaf resulting from the latter case is determined by selecting the majority of the target attribute value in the remaining training set.

Decision tree pruning is a process of replacing sub-trees with leaves to reduce the size of the decision tree while retaining and hopefully increasing the accuracy of tree's classification. To obtain the best result from the induction decision tree method, we varied the use of pruning algorithm to the initial decision tree. We considered using the following decision trees pruning algorithms: *critical-value pruning* (Mingers, 1987), *minimum-error pruning* (Niblett & Bratko, 1986), *pessimistic pruning*, *cost-complexity pruning* and *reduced-error pruning* (Quinlan, 1987),.

2) Induction Decision Trees with Boosting

Decision Tree with Boosting is a method that generates a sequence of decision trees from a single training set by re-weighting and re-sampling the samples in the set (Quinlan, 1996; Freund & Schapire, 1996). Initially, all samples in the training set are equally weighted so that their sum is one. Once a decision tree has been created, the samples in the training set are re-weighted in such a way that misclassified examples will get higher weights than the ones that are easier to classify. The new samples weights are then renormalized, and next decision tree is created using higher-weighted samples in the training set. In effect, this process enforces the more difficult samples to be learned more frequently by decision trees. The trees generated are then given weights in accordance with their performance in the training examples (e.g., their accuracy in correctly classifying the training data). Given a new instance, the new instance class is selected from the maximum weighted average of the predicted class over all decision trees.

3) Neural Networks

Inspired in part by biological learning systems, neural networks approach is built from a densely interconnected set of simple units. Since this technique offers many design selections, we fixed some of them to the ones that had been well proven to be good or acceptable design choices. In particular, we use a network architecture with one hidden layer, the back-propagation learning algorithm (Rumelhart, Hinton & William, 1986), the delta-bar-delta adaptive learning rates (Jacobs, 1988), and the Nguyen-Widrow weight initialization (Nguyen & Widrow, 1990). Discrete-valued attributes fed into the networks input layer is represented as 1-of-N encoding using bipolar values to denote the presence (e.g., value 1) and the absence (e.g., value -1) of an attribute value. Continuous-valued attribute is scaled into a real value in the range [-1, 1]. We varied the design choices for batch versus incremental learning, and 1-of-N versus single network output encoding.

4) Naive Bayes

Naive Bayes classifier is a variant of the Bayesian learning that manipulates directly probabilities from observed data and uses these probabilities to make an optimal decision. This approach assumes that attributes are conditionally independent given a class. Based on this simplifying assumption, the probability of observing attributes' conjunction is the product of the probabilities for the individual attributes. Given an instance with a set of attribute-value pairs, the Naive Bayes approach will choose a class that maximizes the conditional probability of the class given the conjunction of attributes values. Although in practice the independence assumption is not entirely correct, it does not necessarily degrade the system performance (Domingos & Pazzani, 1997).

We also assume that the values of continuous-valued attributes follow Gaussian distribution. Hence, once the mean and the standard deviation of these attributes are obtained from the training examples, the probability of the corresponding attributes can be calculated from the given attribute values. To avoid zero frequency count problem that can dampen the entire probabilities calculation, we use an *m-estimate* approach in calculating probabilities of discrete-valued attributes (Mitchell, 1997).

4.3 Comparison results

For each database, *k*-fold cross-validation method is used for evaluation. In this method, a data set is divided equally into *k* disjoint subsets. *k* experiments are then performed using *k* different training-test set pairs. A training-test set pair used in each experiment is generated by using one of the *k* subsets as the test set, and using the remaining subsets as the training set. Given *k* disjoint subsets, for example, the first experiment takes the first subset for the test set, and uses the second subset through the *k*th subset for the training set. The second experiment uses subset 1 and subset 3 through subset *k* for the training

set; and takes subset 2 for the test set, and so on. All results from a particular database are averaged along with its variance over k experiment runs.

Based on their size, the credit database and voting database are partitioned into 10 disjoint subsets, the heart database is partitioned into 5 subsets. Table 1, 2, 3 show the performance comparison results of different approaches on these three databases. For decision trees with and without boosting, we present only the best experimental results after varying the data splitting methods as well as the pruning algorithms described earlier. Similarly, results from neural networks approach are selected from the ones that provide the best performance after varying the use of different network output encoding and *batch* versus *incremental* learning methods.

Table 1 The comparison results on the prediction accuracy and standard deviation (%) of credit database.

	Our GA approach	Decision trees (IG, Min-Err)	Decision trees with boosting (RG, Cost-Com, 21 trees)	Neural networks (1-of-N, batch learning)	Naive Bayes
Run 1	90.77	87.69	89.23	89.23	66.15
Run 2	89.23	84.62	86.15	86.15	78.46
Run 3	89.23	89.23	90.77	90.77	84.62
Run 4	92.31	90.77	90.77	89.23	81.54
Run 5	86.15	81.54	81.54	84.62	75.38
Run 6	89.23	87.69	87.69	87.69	80.00
Run 7	84.62	81.54	84.62	84.62	73.85
Run 8	87.69	86.15	87.69	86.15	83.08
Run 9	90.77	86.15	89.23	87.69	76.92
Run 10	86.76	88.24	91.18	86.76	75.00
Average	88.68	86.36	87.89	87.29	77.50
Standard deviation	2.37	3.06	3.08	2.03	5.36

In above table, the decision tree is generated using information gain for data splitting and minimum-error pruning. Decision trees with boosting generates 21 different decision trees,

each is constructed by using gain ratio for data splitting and cost-complexity pruning algorithm. Batch learning and 1-of-N output encoding are used in neural networks.

Table 2 The comparison results on the prediction accuracy and standard deviation (%) of voting database.

	Our GA approach	Decision trees (IG, Red-Err)	Decision trees with boosting (IG, Red-Err, 3 trees)	Neural networks (1-of-N, batch learning)	Naive Bayes
Run 1	95.35	95.35	95.35	93.02	93.02
Run 2	97.67	93.02	97.67	97.67	90.70
Run 3	97.67	97.67	97.67	97.67	93.02
Run 4	95.35	95.35	97.67	97.67	88.37
Run 5	95.35	95.35	95.35	93.02	88.37

Run 6	97.67	97.67	97.67	97.67	88.37
Run 7	95.35	93.02	93.02	95.35	93.02
Run 8	97.67	95.35	95.35	95.35	90.70
Run 9	100.00	100.00	97.67	95.35	90.70
Run 10	95.83	93.75	97.92	95.03	85.42
Average	96.79	95.65	96.54	95.86	90.17
Standard deviation	1.59	2.23	1.67	1.46	2.53

In table 2, the best results from both decision trees with and without boosting are obtained from using information gain for data splitting and reduced-error pruning algorithm. Only three decision trees are needed in the decision trees with boosting.

Table 3 The comparison results on the prediction accuracy and standard deviation (%) of heart database.

	Our GA approach	Decision trees (IG, Red-Err)	Decision trees with boosting (IG, Red-Err, 21 trees)	Neural networks (1-of-N, incr learning)	Naive Bayes
Run 1	89.83	81.36	88.14	81.36	89.83
Run 2	83.05	77.97	84.75	84.75	79.66
Run 3	81.36	76.27	76.27	74.58	83.05
Run 4	88.14	76.27	83.05	89.83	83.05
Run 5	83.33	66.67	71.67	81.67	80.00
Average	85.14	75.71	80.77	82.44	83.12
Standard deviation	3.64	5.46	5.54	5.66	4.08

In table 3, the best results from neural networks are obtained from applying incremental learning and 1-of-N network output encoding.

From the above results we can see that our GA approach outperformed other approaches on both the average prediction accuracy and the standard deviation. The advantage of our GA approach becomes more obvious on heart database, which is most difficult to learn among the three. During the process of running, we also found that the training accuracy and testing accuracy of GA approach are basically in the same level, while the training accuracy is often much higher than testing accuracy for other approaches. This proves that GA approach is less sensitive to noise and might be more effective for future prediction.

5. Conclusions and future work

In this paper we put forward a genetic algorithm approach for classification problem. An individual in a population is a complete solution candidate that consists of a fixed number of rules. Rule consequent is not explicitly encoded in the string but determined by the match situation on training examples of the rule. To consider the performance affecting factors as complete as possible, four elements are included in the fitness function which are predicting error rate, entropy measure, rule consistency and hole ratio,

respectively. Adaptive asymmetric mutation and two-point crossover are adopted in reproduction step. The inversion probability of 1-0 (0-1) in mutation is self-adaptive by the feedback of average fitness during the run. The generated classifier after evolution is voting-based. Rules are not disjoint but allowed to overlap. Classifier gives all rules the equal weight for their votes. We tested our algorithm on three real databases and compared the results with four other traditional data mining approaches. It is shown that our approach outperformed other approaches on both prediction accuracy and the standard deviation.

Further testing on various databases is in progress to test the robustness of our algorithm. Splitting continuous attribute into multiple intervals rather than just two intervals based on a single threshold is also considered to improve the performance.

Bibliography

- De Jong, K. A., Spears, W. M., and Gordon, D. F. (1993) Using genetic algorithms for concept learning. *Machine Learning*, 13, 161-188.
- Domingos, P. and Pazzani, M. (1997) On the Optimality of the Simple Bayesian Classifier under Zero-One Loss. *Machine Learning*, 29, 103-130.
- Freund, Yoav, and Schapire, R. E. (1996) Experiments with a new boosting algorithm. In *Machine Learning: Proceedings of the Thirteenth International Conference*, pp. 148-156.
- Greene, D., P. and Smith, S. F. (1993) Competition-based induction of decision models from examples. *Machine Learning*, 13, 229-257.
- Jacobs, R.A. (1988) Increased Rates of Convergence Through Learning Rate Adaptation. *Neural Networks*, 1(4): 295-307.
- Janikow, C. Z. (1993) A knowledge-intensive genetic algorithm for supervised learning. *Machine Learning*, 13, 189-228.
- Mingers, J. (1987) Expert Systems – Rule Induction with Statistical Data. *Journal of the Operational Research Society*, 38, 39-47.
- Mitchell, Tom. (1997) *Machine Learning*. New York: McGraw-Hill.
- Niblett, T. (1986) Constructing Decision Trees in Noisy Domains. In I. Bratko and N. Lavrac (Eds). *Progress in Machine Learning*. England: Sigma Press.
- Noda, E., Freitas, A. A. and Lopes, H. S. (1999) Discovering interesting prediction rules with a genetic algorithm. In *Proceedings of 1999 Congress on Evolutionary Computation (CEC' 99)*, pp. 1322-1329.
- Nguyen, D. and Widrow, B. (1990) Improving the Learning Speed of Two-Layer Networks by Choosing Initial Values of the Adaptive Weights. *International Joint Conference on Neural Networks*, San Diego, CA, III:21-26.
- Quinlan, J.R. (1986) Induction of Decision Trees. *Machine Learning*, 1, 81-106.
- Quinlan, J.R. (1987) Simplifying Decision Trees. *International Journal of Man-Machine Studies*, 27, 221-234.
- Quinlan, J. R. (1996) Bagging, Boosting, and C4.5. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pp. 725-730.
- Rumelhart, D. E., Hinton, G.E., and William, R. J. (1986) Learning Representations by Back-Propagation Error. *Nature*, 323:533-536.
- Yang, L. and Yen, J. (2000) An adaptive simplex genetic algorithm. In *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2000)*, July 2000, pp. 379.