

Wrapper-Filter Feature Selection Algorithm Using A Memetic Framework

Zexuan Zhu, Yew-Soon Ong, and Manoranjan Dash

Abstract—This paper presents a novel hybrid wrapper and filter feature selection algorithm for classification problem using a memetic framework. It incorporates filter ranking method in the traditional genetic algorithm to improve classification performance and accelerate the search in identifying the core feature subsets. Particularly, the method adds or deletes a feature from a candidate feature subset based on the univariate feature ranking information. Our empirical study on commonly used datasets from the UCI repository and microarray datasets show that the proposed method outperforms existing methods in terms of classification accuracy, number of selected features and computational efficiency. Further, we investigate several major issues of memetic algorithm to identify a good balance between local and genetic search so as to maximize search quality and efficiency in the hybrid filter and wrapper memetic algorithm.

Index Terms—Feature Selection, Filter, Wrapper, Memetic Algorithm (MA), Genetic Algorithm (GA), Hybrid Genetic Algorithm, Relief, Gain Ratio, Chi-Square

I. INTRODUCTION

FEATURE selection has become the focus of many research areas in recent years. With the rapid advance of computer and database technologies, datasets with hundreds and thousands of variables or features are now ubiquitous in pattern recognition, data mining, and machine learning [1-4]. To process such huge datasets is a challenging task because traditional machine learning techniques usually work well only on small datasets. Feature selection addresses this problem by removing the irrelevant, redundant, or noisy data. It improves the performance of the learning algorithm, reduces the computational cost and provides better understandings of the datasets.

Feature selection algorithms may be widely categorized into two groups: filter and wrapper methods [2, 4-10]. Filter methods evaluate the goodness of the feature subset by using the intrinsic characteristic of the data. They are relatively computationally cheap, since they do not involve the induction algorithm. However, they also take the risk of selecting subsets of features which may not match the chosen induction algorithm. Wrapper methods, on the contrary, directly use the induction algorithm to evaluate the feature subsets. They generally outperform filter methods in terms of prediction

accuracy, but are generally computationally more intensive. In summary, wrapper and filter methods can complement each other, in that filter methods can search through the feature space efficiently while the former provides good accuracy. In this paper, we propose a novel wrapper-filter feature selection algorithm (WFFSA) using a memetic framework [11-16], i.e., a combination of genetic algorithm (GA) [17-20] and local search (LS). Memetic algorithms (MAs) are population-based meta-heuristic search methods inspired by Darwinian's principles of natural evolution and Dawkins' notion of a meme defined as a unit of cultural evolution that is capable of local refinements. Recent studies on MAs have revealed their successes on a wide variety of real world problems. Particularly, they not only converge to high quality solutions, but also search more efficiently than their conventional counterparts [11-16].

The goal of WFFSA is to improve classification performance and accelerate the search to identify important feature subsets. In particular, the filter method fine-tunes the population of GA solutions by adding or deleting features based on univariate feature ranking information. Hence, our focus here is on filter methods that are able to assess the goodness or ranking of the individual features. We denote such filter methods as filter ranking methods in this paper and investigate the proposed WFFSA for several filter ranking methods. Empirical study of WFFSA on several commonly used datasets from the UCI repository [21] and several microarray datasets indicates that it outperforms recent existing methods in the literature in terms of classification accuracy, selected feature size and efficiency. Further, we also investigate the balance between local and genetic search to maximize the search quality and efficiency of WFFSA.

The rest of this paper is organized as follows. Section II describes the wrapper-filter feature selection algorithm based on a memetic framework. The experimental results and discussions are presented in Section III. Finally, Section IV concludes this study.

II. WRAPPER-FILTER FEATURE SELECTION ALGORITHM—WFFSA

In this section, we introduce the proposed Wrapper-Filter Feature Selection Algorithm for classification problems which is depicted in Figure 1. In the first step, the GA population is initialized randomly with each chromosome encoding a candidate feature subset. Subsequently, on all or portion of the chromosomes, a local search or meme is applied in the spirit of Lamarckian learning [11, 14]. The mechanism to do local

Z. Zhu, Y.S.Ong and M. Dash are with the Division of Information Systems, School of Computer Engineering, Nanyang Technological University, 639798, Singapore (e-mail: zhuzexuan@pmail.ntu.edu.sg, asysong@ntu.edu.sg, asmdash@ntu.edu.sg).

Fig. 1 The Procedure of WFFSA

improvement can be to reach a local optimum or to improve the solution. Genetic operators are then used to generate the next population. This process repeats until the stopping conditions are satisfied. We describe each component in detail as follows.

A. Encoding Representation and Initialization

In the feature selection problem, a representation for candidate feature subset must be chosen and encoded as a chromosome. In most studies, a chromosome is a binary string of length equal to the total number of features so that each bit encodes a single feature (as shown in Figure 2). A bit of '1' ('0') implies the corresponding feature is selected (excluded). The length of the chromosome is denoted here as n . The maximum allowable number of bit '1' in each chromosome is denoted as m . When prior knowledge about the optimal number of features is available, we may limit m to no more than the pre-defined value; otherwise m is equal to n . At the start of the search, a population size of p is randomly initialized.

Fig. 2 Representation of chromosome as a binary bit string

B. Objective Function

The objective function is simply defined by the classification accuracy:

$$Fitness(c) = J(S_c) \quad (1)$$

where S_c denotes the corresponding selected feature subset encoded in chromosome c , and the feature selection criterion function $J(S_c)$ evaluate the significance for the given feature subset S_c , in this study, $J(S_c)$ is specified as the classification accuracy for S_c . Note that when two chromosomes are found having similar fitness, i.e., the difference between their fitness is less than a small value of ϵ , the one with a smaller number of selected features is given higher chances of surviving to the next generation.

C. Local Search Improvement Procedure

Much work on the use of domain knowledge and heuristics has resulted in highly effective search [12, 14, 15]. Taking this cue, we consider here the use of filter ranking methods as memes or local search heuristics in our WFFSA. Later, we show in Section III that using filter ranking methods as memes, MA is capable of converging to improved classification accuracy and at lower number of selected features when compared to existing methods recently proposed in the literature.

Given a candidate chromosome c , we define X and Y as the sets of selected and excluded features encoded in c , respectively. Both X and Y are ranked using the univariate filter ranking method and with the most important feature ranked the highest. In this study, we consider three different filter ranking methods, namely, ReliefF [5], Gain Ratio [22] and Chi-Square [6]. These methods rank features based on different criteria that include Euclidean distance, information entropy and chi-square statistics respectively. We further define two basic local search operators of the WFFSA local search improvement procedure:

- i) *Add*: select a feature from Y using the linear ranking selection [23] and move it to X .
- ii) *Del*: select a feature from X using the linear ranking selection [23] and move it to Y .

The *Add* and *Del* operations are illustrated in Figure 3. Here, F5 and F4 are the highest and lowest ranked features in Y ; F3 and F6 are the highest and lowest ranked features in X . Using the *Add* and *Del* operations, F5 is the most likely feature to be moved to X while F6 is the most likely feature to be moved to Y . The two most likely resultant chromosomes after the *Add* and *Del* operations are also depicted in Figure 3.

Fig. 3 *Add* and *Del* Operations

The intensity of local search is quantified by the local search length l and interval w . Local search length defines the maximum number of *Del* and *Add* operations in each local search. Therefore, there are a total of l^2 possible combinations of *Add* and *Del* operations applied on a chromosome. Local search interval specifies the w elite chromosomes in the population that undergo local search improvement procedure in each generation. The local search improvement procedure may be applied on a chromosome until a local optimum or an improvement is reached. The locally improved chromosome is evaluated and replaces the original chromosome if it is of higher quality. Here, we further investigate three different local search strategies that are characterized by different local search intensity.

1) Improvement First Strategy

In this strategy, a random choice from the l^2 combinations of *Del* and *Add* operations is used to search on the candidate chromosome. The local search stops once an improvement is obtained either in terms of classification accuracy or a reduction in the number of selected features without deterioration in accuracy greater than ϵ . This procedure is outlined in Figure 4.

Fig. 4 The procedure of improvement first strategy

2) Greedy Strategy

In contrast to the improvement first strategy, the greedy strategy carries out all possible l^2 combinations of *Del* and *Add* operations and the best improved solution is used to replace the original chromosome in the population. The greedy strategy is outlined in Figure 5.

Fig. 5 The procedure of greedy strategy

3) Sequential Strategy

We also consider the sequential strategy described in [16]. The Hybrid genetic algorithm (HGA) was reported in [16] to generate better search performances than GA, SFS (sequential forward search), SFFS (sequential forward floating search), PTA(l,r) (plus- l and take away- r) and multi-start algorithms. There, instead of using a filter ranking method, the *Add* operation searches for the most significant feature y in Y in a

sequential manner, i.e., $y = \arg \max_{a \in Y} J(X \cup \{a\})$, and moves it to X . In the same way, the *Del* operation searches for the least significant feature x from X in a sequential manner, i.e., $x = \arg \max_{a \in X} J(X - a)$, and moves it to Y .

D. Evolutionary Operators

In the evolution process, standard GA operators such as linear ranking selection, uniform crossover and mutation operators based on elitist strategy may be applied. However, if prior knowledge on the optimum number of features is available, the number of bit '1' in each chromosome may be constrained to a maximum of m in the evolution process. Since the standard uniform crossover and mutation operators may violate this constraint, restrictive crossover and mutation are proposed here. In restrictive crossover, the crossover operations are applied only on alleles with bit '1' and in either of the two parent chromosomes. We outline the restrictive crossover in Figure 6. Based on the same principles, the restrictive mutation operator is outlined in Figure 7.

Fig. 6 The procedure of restrictive crossover

Fig. 7 The procedure of restrictive mutation

E. Computational Complexity

In this section, we analyze the computational complexity of the proposed WFFSA. The ranking of features based on the filter methods have linear time complexity in terms of feature dimensionality, they are conducted offline and the obtained rank list may be reused for each local search in WFFSA. Consequently, the computational for feature ranking is a one-time offline cost and is considered to be negligible compared to that of fitness evaluation in equation (1). Hence, we define the computational cost of a single fitness evaluation as the basic unit of computational cost in our analysis.

The computational complexity for GA can be derived as $O(pg)$, where p is the size of population and g is the number of search generations. The expected computational complexity of WFFSA with improvement first strategy is $O(pl^2wg/2)$. Here we assume that each l^2 combinations of *Add* and *Del* operation has equivalent likelihood to obtain the first improvement, therefore the average trails to obtain the first improvement is $l^2/2$. In a single search generation, $l^2w/2$ fitness function calls are incurred. The ratio of local to genetic search is thus $l^2w/2p$.

In the greedy strategy, local search evaluates all possible l^2 combinations of *Add* and *Del* operations to attain the best possible locally improved solution. The computational complexity is thus $O(pl^2wg)$. In a single search generation, l^2w fitness function calls are incurred and the ratio for local to genetic search is l^2w/p .

For the sequential strategy proposed in HGA [16], a total of $[|Y|+(|Y|-l)]l/2$ and $[|X|+(|X|-l)]l/2$ calls to the fitness function are incurred on the *Add* and *Del* operations, respectively. Since

$|X|+|Y|=n$, the total classification calls in the local search is $(n-l)lw/2$. The computational complexity of sequential strategy is thus $O(p(n-l)hwg/2)$. The ratio for local to genetic search is then $(n-l)lw/2p$.

Since we are working with classification problems where $n \gg l$, it is easy to determine that $(n-l)lw/2 \gg l^2w > l^2w/2$. Consequently, the sequential local search strategy in HGA [16] would require significantly more computations, hence incurring more time than both the improvement first and greedy strategies of the WFFSA.

III. EXPERIMENTAL RESULTS AND DISCUSSIONS

In this section, we present an experimental study of WFFSA on commonly used benchmark and biological datasets. In particular, we focus on four UCI datasets having more than 15 features and four microarray datasets (ALL/AML [24], Colon [25], NCI60 [26], SRBCT [27]) having thousands of features in our study.

In the WFFSA, we employed a population size of 30 and stopping criterion of 6000 fitness function calls for UCI datasets. On microarray datasets, as the feature size are significantly larger, the population size is increased to 50 or 100 and the stopping criterion to 10000 or 20000. Further, the maximum number of selected features is unconstrained for the UCI datasets, i.e., $m=n$. On the other hand, it has been shown in the literature [8, 24-28] that microarray datasets could be learned with high accuracies with only hundreds or tens of features. In effect, we make use of this prior knowledge to constrain the maximum number of selected features m as depicted in Table 1. In effect, the restrictive crossover and mutation operations are applied on the microarray datasets. In our experimental setup, we employ crossover and mutation probabilities of $p_c=0.6$ and $p_m=0.1$, respectively. Linear ranking selection [23] with selection pressure of 1.5 is used for selection. The threshold ϵ to determine fitness similarity between two chromosomes is configured as 0.001 for UCI datasets and 0.02 for microarray datasets. The fitness of a chromosome or selected feature subset is evaluated using the 1-nearest neighbor (1NN) classifier and the leave-one-out cross validation (LOOCV). Here we use the classification accuracy estimated from LOOCV and the number of selected features as performance measures. It is worth noting that the configurations of the parameter used here have been investigated empirically for the datasets considered and are summarized in Table 1.

TABLE I
DATASETS AND PARAMETERS USED FOR EXPERIMENTS

A. Comparison of Filter, GA and WFFSA

The search performances of the filter ranking methods, GA and, WFFSA using several filter ranking methods and local search strategies on the eight datasets considered are summarized in Table 2. The comma delimited pair wise numeric values in the table represent the classification accuracy and the corresponding number of selected features. Due to the

TABLE 2
COMPARISON RESULTS OF FILTER, GA, AND WFFSA METHODS

stochastic nature of GA and WFFSA, the average results for ten independent runs are reported. The parentheses highlights the best result found across the ten runs. For each dataset, the bold-faced and bold-italic-faced representations in Table 2 emphasize the best average performance and the best solution found among all methods, respectively. Six WFFSAs obtained based on the three filter ranking methods (i.e., ReliefF—WFFSA-R or Gain Ratio—WFFSA-G or Chi-Square—WFFSA-C) and two local search strategies (i.e., Improvement First or Greedy) have been investigated. In the WFFSAs, the local search length l and interval w in improvement first and greedy strategies are configured as $l=4, w=1$ and $l=4, w=5$, respectively.

We discuss the performances for the various feature selection algorithms considered. On filter ranking methods, the t features with highest rank are chosen to induce the 1NN LOOCV classification. We increase t from 1 to m and the optimal value is determined when the best classification accuracy is obtained. The best classification accuracy and the corresponding t value for each filter ranking methods are reported in Table 2 (i.e., columns 3-5). It can be observed that WFFSAs outperform all three filter ranking methods and GA in terms of classification accuracy. The improvement in performance is more significant for the Sonar, Colon and NCI60 datasets. Moreover, WFFSAs reduce the number of selected features significantly. On the ALL/AML and SRBCT datasets, WFFSAs use less than one-third of the features required by GA and filter ranking methods to arrive at the improved classification accuracy. The best solutions found on all eight datasets shown in column 2 were attained by WFFSAs. We also study the performance of WFFSA for three different filter ranking methods since it was shown in [11] that inappropriate use of meme may result in the MA performing poorer than standard GA. The results in Table 2 indicate that WFFSA-R performs better than the other counterparts on 7 out of the 8 datasets in terms of average performance. For ReliefF filter ranking method, the feature goodness is evaluated by its ability to distinguish the near hit (nearest neighbors from the same class) and near miss (nearest neighbors from different classes). Therefore, it makes good sense that features with high score in ReliefF are more likely to help 1NN identify the correct nearest neighbor, hence generating good classification accuracy. On the other hand, Gain Ratio and Chi-Square do not appear to possess such mechanisms to compliment 1NN, hence being less effective here.

Further, we compare the results obtained by WFFSAs to the recently proposed HGA(3) [16] for the UCI datasets in Table 3. The best average performance and the best solution for each dataset are highlighted using bold typeface representation. The results indicate that WFFSAs generate the best solutions on all four datasets and give competitive results to existing methods

TABLE 3
RESULTS BY WFFSAS AND OTHER METHODS ON UCI DATASETS

in terms of average performance. Even so, it is worth noting that HGA(3) consumed more than 200000 fitness function calls to arrive at the competitive performance on these UCI datasets. In contrast, WFFSA incurs less than 6000 fitness function calls to arrive at superior or competitive performances.

To illustrate the generality and efficacy of the WFFSA framework, we consider also the use of different induction algorithms, particularly standard 1NN and Support Vector Machine with Radius Margin Bound (SVM) [29], and compare their performances with other recent studies in the literature using internal and external cross-validation schemes on the real world microarray datasets. In external cross-validation each dataset is randomly split into k stratified folds with $k-1$ folds for training and 1 fold for testing. The performance of the final selected features obtained is measured on the unseen testing data. The procedure is repeated k times. The best average performance and/or the best solution are reported in Table 4. Table 4 indicates that WFFSAs displays superior performances on most of the microarray datasets in comparison to the existing counterparts.

TABLE 4
RESULTS BY WFFSAS AND OTHER METHODS ON MICROARRAY DATASETS

B. Study on Local Search Strategies

Over the recent years, much work has shown that an appropriate balance between local and genetic search is necessary for efficient memetic search [11-15]. In this subsection, we examine the balance between local and genetic search using different local search strategies, which are characterized by different intensities of genetic and local searches. Ten independent runs of WFFSA with different combinations of strategy, local search length l and interval w are conducted on two representative UCI and microarray datasets, i.e., the Sonar and Colon datasets. The results obtained are summarized in Tables 5 and 6. The local search heuristic or filter ranking method used is ReliefF, since it is shown to give better search performance than the other counterparts.

The results in Tables 5 and 6 indicate that WFFSA for $l = 4$ and $w = 1$ obtains the best average accuracy on the improvement first strategy. On the Sonar dataset, WFFSA for $l = 4$ and $w = 1$ and using the improvement first strategy is found to perform significantly better than the other configurations of WFFSAs statistically, using the two-tailed paired t-test at significance level of 0.05. On the Colon dataset, the superiority of WFFSA for $l = 4, w = 1$ and an improvement first strategy is however not statistically conclusive using the same t-test.

TABLE 5
RESULTS OF DIFFERENT LOCAL SEARCH STRATEGIES, LENGTH AND INTERVAL ON SONAR DATASET

TABLE 6
RESULTS OF DIFFERENT LOCAL SEARCH STRATEGIES, LENGTH AND INTERVAL ON COLON DATASET

TABLE 7
RESULTS OF DIFFERENT LOCAL STRATERGIES ON ALL EIGHT DATASETS

Overall, we note that WFFSA with improvement first strategy generally performs better than using the greedy strategy. WFFSA obtains better results when local search is applied only to the single or five elite chromosomes instead of the entire population. In line with the observations in [12, 14], conducting local search on all chromosomes can result in unnecessary computation. Since the computational budget is often limited (i.e., based on the maximum number of fitness function calls allowed), the genetic search will be reduced proportionally due to possible redundant computation spent on local search. Without sufficient genetic search, the memetic algorithm is more likely to be trapped in the local optimum and hence may lead to poor search quality under limited computational budget.

We also compare WFFSA-R with a memetic algorithm based on sequential strategy (MA-S). In MA-S, we configure $l = 4$ and $w = 5$. The comparison results are reported in Table 7. On all the datasets, both WFFSA-R with improvement first strategy or greedy strategy, (they are labeled as WFFSA-R First and WFFSA-R Greedy in Table 7) displays significantly better performance than MA-S statistically using two-tailed paired t-test at significance level of 0.05. The superior performance of WFFSA over MA-S is more obvious on the microarray datasets where $n \gg l$. So far, the results obtained further strengthen the importance of balance tradeoff between local search and genetic search for efficient MA search. A search dominated by either genetic search (e.g., GA) or local search (MA-S) would generally not perform effectively. The WFFSA for $l = 4$ and $w = 1$ based on the improvement first strategy thus gives the most appropriate tradeoff of local and genetic search than GA, MA-S and all other configurations considered in our present study. As shown in Figure 8, the average search trends of WFFSA-R First, WFFSA-R Greedy, MA-S ($l = 4, w = 1$) and MA-S ($l = 4, w = 5$) on the Sonar and NCI datasets suggest that WFFSA-Rs search more efficiently than GA and MA-S under limited computational budget.

Fig. 8 (a)

Fig. 8 (b)

Fig. 8 Search trace (average of 10 independent runs) of GA, WFFSA-R Improvement First ($l=4, w=1$), WFFSA-R Greedy ($l=4, w=5$), MA-S ($l=4, w=1$) and MA-S ($l=4, w=5$) on (a) Sonar dataset (b) NCI dataset

IV. CONCLUSIONS

In this paper, we have proposed a novel hybrid filter and wrapper feature selection algorithm based on a memetic framework. We use filter ranking method as local search heuristic in the memetic algorithm. The experimental results presented show that the proposed method searches more efficiently and is capable of producing good classification accuracy with small number of features simultaneously. Most importantly, it outperforms GA, MA with sequential local search as well as many existing algorithms in the literature. Further, our study on various local search strategies, local search length and interval allow us to identify a suitable balance tradeoff of genetic and local search in the memetic search. This

allows us to maximize the effectiveness and efficiency of the proposed hybrid filter and wrapper feature selection algorithm for classification problem using a memetic framework.

REFERENCES

- [1] A. Jain and D. Zongker, "Feature Selection: Evaluation, Application, and Small Sample Performance," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 19, no. 2, pp. 153-158, 1997.
- [2] M. Dash and H. Liu, "Feature Selection for Classification," *Intelligent Data Analysis: An Int'l J.*, vol. 1, no. 3, pp. 131-156, 1997.
- [3] M. Dash and H. Liu, "Consistency-based Search in Feature Selection," *Artificial Intelligence*, vol. 151, no. 1-2, pp. 155-176, 2003.
- [4] R. Kohavi and G. H. John, "Wrapper for Feature Subset Selection," *Artificial Intelligence*, vol. 97, no. 1-2, pp. 273-324, 1997.
- [5] M. Robnic-Sikonja and I. Kononenko, "Theoretical and Empirical Analysis of ReliefF and RReliefF," *Machine Learning*, vol. 53, no. 1-2, pp. 23-69, 2003.
- [6] H. Liu and R. Setiono, "Chi2: Feature Selection and Discretization of Numeric Attributes," *In Proceedings of 7th IEEE Int'l Conference on Tools with Artificial Intelligence*, pp. 388-391, 1995.
- [7] E. Xing, M. Jordan, and R. Karp, "Feature Selection for High-Dimensional Genomic Microarray Data," *Proc. 15th Int'l Conf. Machine Learning*, pp. 601-608, 2001.
- [8] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene Selection for Cancer Classification Using Support Vector Machines," *Machine Learning*, vol. 46, no. 1-3, pp. 389-422, 2002.
- [9] K. Z. Mao, "Feature Subset Selection for Support Vector Machines Through Discriminative Function Pruning Analysis," *IEEE Transactions on System, Man and Cybernetics, Part B*, vol. 34, no. 1, pp. 60-67, 2004.
- [10] C. N. Hsu, H. J. Huang and S. Dietrich, "The ANNIGMA-Wrapper Approach to Fast Feature Selection for Neural Nets," *IEEE Transactions on System, Man and Cybernetics, Part B*, vol. 32, no. 2, pp. 207-212, 2004.
- [11] Y. S. Ong and A. J. Keane, "Meta-Lamarckian in Memetic Algorithm," *IEEE Trans. Evolutionary Computation*, vol. 8, no. 2, pp. 99-110, 2004.
- [12] Y. S. Ong, M. H. Lim, N. Zhu and K. W. Wong, "Classification of Adaptive Memetic Algorithms: A Comparative Study", *IEEE Transactions On Systems, Man and Cybernetics, Part B*, vol. 36, no. 1, pp. 141-152, 2006
- [13] H. Ishibuchi, T. Yoshida, and T. Murata, "Balance Between Genetic Search and Local Search in Memetic Algorithm for Multiobjective Permutation Flowshop Scheduling," *IEEE Trans. Evolutionary Computation*, vol. 7, no. 2, pp. 204-223, 2003.
- [14] N. Krasnogor, "Studies of the Theory and Design Space of Memetic Algorithms," *Phd Thesis, University of the West of England, Bristol*, 2002.
- [15] Q. Zhang, J. Sun, E. Tsang and J. Ford, "Hybrid Estimation of Distribution Algorithm for Global Optimization", *Engineering Computations*, vol. 21, no.1, pp 91-107, 2004.
- [16] I. S. Oh, J. S. Lee, and B. R. Moon, "Hybrid Genetic Algorithm for Feature Selection," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 26, no. 11, pp. 1424-1437, 2004.
- [17] W. Siedelecky and J. Sklansky, "On Automatic Feature Selection," *International Journal of Pattern Recognition and Artificial Intelligence*, no.2, pp. 197-220, 1988
- [18] M. L. Raymer, W. F. Punch, E. D. Goodman, L. A. Kuhn, and A. K. Jain, "Dimensionality Reduction Using Genetic Algorithms," *IEEE Trans. Evolutionary Computation*, vol. 4, no. 2, pp. 164-171, 2000.
- [19] L. Li, C. R. Weinberg, T. A. Darden, and L. G. Pedersen, "Gene Selection for Sample Classification Based on Gene Expression Data: Study of Sensitivity to Choice of Parameters of the GA/KNN Method," *Bioinformatics*, vol. 17, no. 12, pp. 1131-1142, 2001.
- [20] J. H. Yang and V. Honavar, "Feature Selection Using a Genetic Algorithm," *IEEE Intelligent Systems*, vol. 13, no. 2, pp. 44-49, 1998.
- [21] P. M. Murphy and D. W. Aha, "UCI Repository for Machine Learning Database," *Technical Report, Dept. of Information and Computer Science, Univ. of California, Irvine, Calif*, 1994.
- [22] J. R. Quinlan, "C4.5: Programs for Machine Learning," *San Mateo, Morgan Kaufman*, 1993.
- [23] J. E. Baker, "Adaptive Selection Methods for Genetic Algorithms," *In Proc. Int'l Conf. Genetic Algorithm and Their Applications*, pp. 101-111, 1985.

- [24] X. Zhou and K. Z. Mao, "LS Bound Based Gene Selection for DNA Microarray Data," *Bioinformatics*, vol. 21, no. 8, pp. 1559-1564, 2005.
- [25] U. Alon, N. Barkai, D. Notterman, K. Gish, S. Ybarra, D. Mack, and A. Levine, "Broad Patterns of Gene Expression Revealed by Clustering Analysis of Tumor and Normal Colon Tissues Probed by Oligonucleotide Arrays," *Proc. Natl. Sci. USA*, vol. 96, no. 12, pp. 6745-6750, 1999.
- [26] C. H. Ooi and P. Tan, "Genetic Algorithms Applied to Multi-Class Prediction for the Analysis of Gene Expression Data," *Bioinformatics*, vol. 19, no. 1, pp. 37-44, 2003.
- [27] J. Khan, J. S. Wei, M. Ringer, L. H. Saal, M. Ladanyi, F. Westermann, F. Berthold, M. Schwab, C. R. Antonescu, C. Peterson, and P. S. Meltzer, "Classification and Diagnostic Prediction of Cancers Using Expression Profiling and Artificial Neural Networks," *Nature Medicine*, vol. 7, no. 6, pp. 673-679, 2001.
- [28] T. Li, C. Zhang and M. Ogihara, "A Comparative Study of Feature Selection and Multiclass Classification Methods for Tissue Classification Based on Gene Expression", *Bioinformatics*, vol. 20, no. 15, pp. 2429-2437, 2004.
- [29] O. Chapelle, V. Vapnik, O. Bousquet and S. Mukherjee, "Choosing Multiple Parameters for Support Vector Machines", *Machine Learning*, vol. 46, no. 1, pp. 131-159, 2002

Procedure of WFFSA

```
1 Begin
2 Initialize: Randomly generate an initial population of feature subsets;
3 While (Stopping conditions are not satisfied)
4     Evaluate all feature subsets encoded in the population;
5     For each subset chosen to undergo the local improvement process
6         Perform local search and replace it with locally improved
           solution in the spirit of Lamarckian learning;
7     End For
8     Perform evolutionary operators based on selection, crossover, and
           mutation;
9 End While
10 End
```

Fig. 1 The Procedure of WFFSA

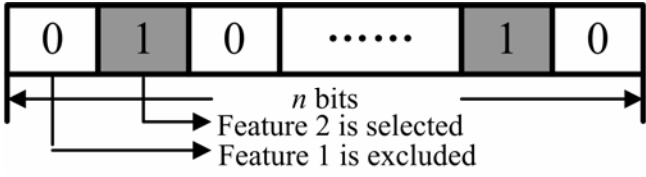


Fig. 2 Representation of chromosome as a binary bit string

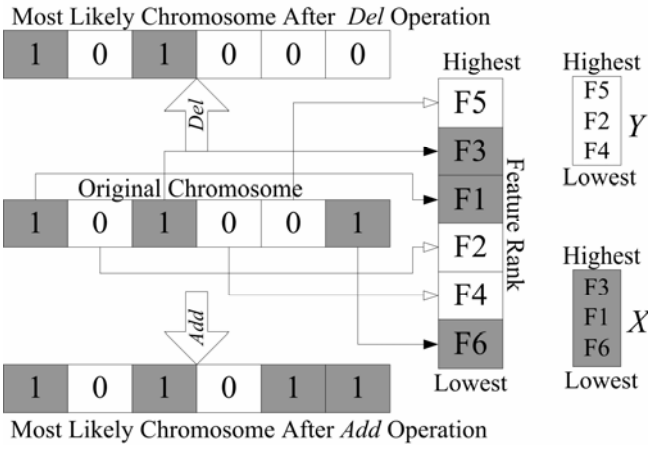


Fig. 3 Add and Del Operations

```

Procedure Improvement First Strategy
1 Begin
2   Initialize  $l$  and  $w$ ;
3   For each chromosome  $c$  among the  $w$  elitists
4     For ( $j = 1$  to  $l^2$ )
5       Generate a unique random pair  $(k,d)$ ;
6       Repeat  $k$  times of Add operation;
7       Repeat  $d$  times of Del operation;
8       Calculate fitness of improved chromosome  $c'$  using  $F(c')=J(c)$ ;
          // $|\cdot|$  denotes the cardinality of a vector
9       If ( $(F(c') > F(c))$  or ( $|F(c') - F(c)| < \epsilon$  and  $|c'| < |c|$ ))
10        Replace the genotype  $c$  with the improved  $c'$ ;
11        Break and consider the next elite chromosome;
12      End If
13    End For
14  End For
15 End

```

Fig. 4 The procedure of improvement first strategy

```

Procedure Greedy Strategy
1 Begin
2   Initialize  $l$  and  $w$ ;
3   For each chromosome  $c$  among the  $w$  elitists
4      $c_{best} = c$ ;
5     For each of the  $l^2$  combinations  $(k,d)$ 
6       Repeat  $k$  times of Add operation;
7       Repeat  $d$  times of Del operation;
8       Calculate fitness of improved chromosome  $c'$  using  $F(c')=J(c')$ ;
9       If  $((F(c') > F(c_{best}))$  or  $(|F(c') - F(c_{best})| < \epsilon$  and  $|c'| < |c_{best}|)$ 
10         $c_{best} = c'$ ; //update the best improved chromosome
11      End If
12    End For
13    Replace the genotype  $c$  with the best improved  $c_{best}$ ;
14  End For
15 End

```

Fig. 5 The procedure of greedy strategy

Procedure Restrictive Crossover

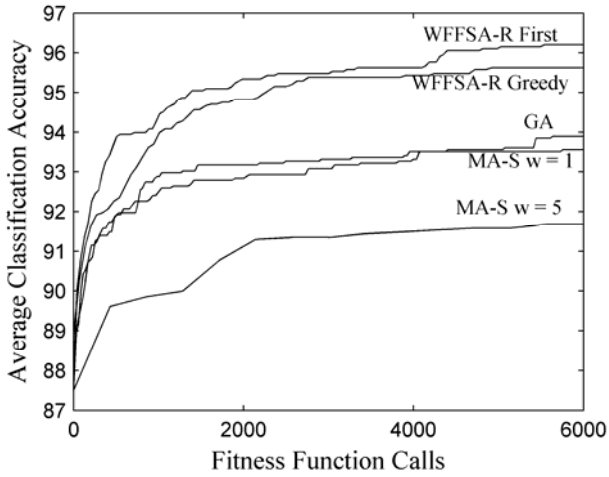
```
1 Begin
2   Randomly select two parents  $p_1$  and  $p_2$ ;
3   Randomly generate a number  $r$  within  $[0,1]$ ;
4   If ( $r < p_c$ ) //  $p_c$  denotes the crossover probability
      // ensure the number of bit '1' in the offspring dose not exceed  $m$ 
5      $k = \text{Min}(|p_1|, |p_2|)$ ; //  $|p_1|, |p_2| < m$ 
6     For ( $i = 1$  to  $k$ )
7       Locate the allele  $L_1$  of the  $i^{\text{th}}$  bit '1' in  $p_1$ ;
8       Locate the allele  $L_2$  of the  $i^{\text{th}}$  bit '1' in  $p_2$ ;
9       Crossover  $p_1$  and  $p_2$  in positions  $L_1$  and  $L_2$  with probability 0.5;
10    End For
11  End If
12 End
```

Fig. 6 The procedure of restrictive crossover

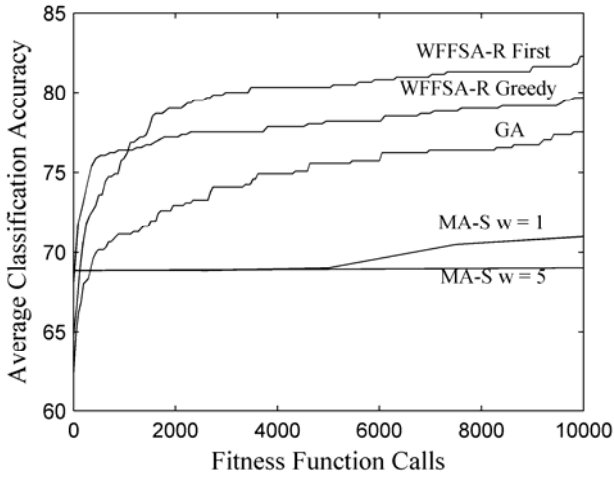
Procedure Restrictive Mutation

```
1 Begin  
   // mutating chromosome  $c$ ;  $p_m$  denotes the mutation probability  
2   For (  $i = 1$  to  $|c|$  )  
3     Locate the position  $L_1$  of the  $i^{\text{th}}$  bit '1' in  $c$ ;  
4     Randomly select a bit '0' with position  $L_0$ ;  
5     Swap positions  $L_1$  and  $L_0$  with probability  $p_m$ ;  
6   End For  
7   For (  $i = 1$  to  $m-|c|$  )  
8     Randomly flip a bit '0' with probability  $p_m$ ;  
9   End For  
10 End
```

Fig. 7 The procedure of restrictive mutation



(a)



(b)

Fig. 8 Search trace (average of 10 independent runs) of GA, WFFSA-R Improvement First ($l=4, w=1$), WFFSA-R Greedy ($l=4, w=5$), MA-S ($l=4, w=1$) and MA-S ($l=4, w=5$) on (a) Sonar dataset (b) NCI dataset

TABLE 1
DATASETS AND PARAMETERS USED FOR EXPERIMENTS

Dataset	Num of Features, n	Num of Instances	Num of Classes	Population Size	p_c, p_m	Maximum Number of Selected Features, m	Feature Subset Evaluation	Stopping Criterion *
Vehicle	18	846	4	30	0.6, 0.1	18	1NN, LOOCV	6000
WDBC	30	569	2	30	0.6, 0.1	30	1NN, LOOCV	6000
Ionosphere	34	351	2	30	0.6, 0.1	34	1NN, LOOCV	6000
Sonar	60	208	2	30	0.6, 0.1	60	1NN, LOOCV	6000
ALL/AML	1000	72	2	50	0.6, 0.1	Constrained to 50	1NN, LOOCV	10000
Colon	1000	62	2	50	0.6, 0.1	Constrained to 50	1NN, LOOCV	10000
NCI60	1000	60	9	50	0.6, 0.1	Constrained to 50	1NN, LOOCV	10000
SRBCT	2308	83	4	100	0.6, 0.1	Constrained to 50	1NN, LOOCV	20000

* Here the stopping criterion is the maximum number of fitness function calls

TABLE 2
COMPARISON RESULTS OF FILTER, GA, AND WFFSA METHODS

Dataset	Best Found	ReliefF	Gain Ratio	Chi Square	GA	WFFSA-R*		WFFSA-G*		WFFSA-C*	
						First**	Greedy**	First	Greedy	First	Greedy
Vehicle (846×18)	75.06, 12	72.81, 12	70.33, 16	72.22, 15	74.88, 10.2 (75.06, 12)	75.00, 11.3 (75.06, 12)	75.06, 12 (75.06, 12)	74.98, 11.5 (75.06, 12)	75.04, 11.7 (75.06, 12)	75.06, 12 (75.06, 12)	75.04, 11.7 (75.06, 12)
WDBC (569×30)	98.24, 12	96.49, 4	96.31, 29	96.31, 29	97.82, 14.7 (98.23, 19)	97.96, 13 (98.24, 14)	98.14, 13.9 (98.24, 14)	97.98, 11.8 (98.24, 14)	97.93, 13.2 (98.07, 11)	97.88, 12 (98.24, 15)	98.00, 12.7 (98.24, 12)
Ionosphere (351×34)	96.01, 8	91.74, 7	90.88, 13	92.31, 7	94.13, 12.5 (94.87, 13)	95.00, 7.5 (95.73, 8)	94.96, 9.7 (95.44, 8)	95.13, 8.9 (96.01, 8)	95.19, 8.9 (95.44, 7)	95.01, 8.9 (95.44, 7)	94.73, 10.4 (95.16, 9)
Sonar (208×60)	97.6, 22	88.46, 15	88.94, 31	89.90, 30	93.89, 27.7 (95.19, 25)	96.30, 24 (97.12, 19)	95.63, 24.6 (96.63, 24)	94.95, 25.2 (97.6, 22)	95.22, 24 (96.63, 24)	95.82, 24.6 (96.63, 24)	95.24, 26.4 (96.15, 26)
ALL/AML (72×1000)	100, 2	98.61, 16	98.61, 17	98.61, 18	100, 26.5 (100, 20)	100, 5.1 (100, 2)	100, 8.3 (100, 4)	100, 5.5 (100, 3)	100, 6.3 (100, 4)	100, 6.1 (100, 2)	100, 7.9 (100, 3)
Colon (62×1000)	100, 14	83.87, 17	90.32, 24	87.10, 17	94.52, 37 (95.16, 31)	97.9, 10.9 (100, 14)	97.26, 16 (100, 20)	96.29, 11.9 (98.39, 11)	97.9, 14.9 (100, 17)	95.97, 10 (98.39, 15)	96.94, 17.8 (100, 16)
NCI60 (60×1000)	85.25, 14	65.57, 16	59.02, 31	65.57, 28	77.54, 38.6 (80.33, 40)	82.30, 22.7 (85.25, 19)	79.67, 23.8 (81.97, 24)	82.13, 21.5 (85.25, 18)	81.31, 27.5 (85.25, 33)	81.15, 21.8 (85.25, 14)	81.48, 26.5 (85.25, 22)
SRBCT (83×2308)	100, 7	100, 38	100, 20	100, 43	99.64, 43.3 (100, 39)	100, 15.1 (100, 7)	100, 18.7 (100, 11)	99.76, 17.1 (100, 8)	100, 21.8 (100, 16)	99.76, 14.9 (100, 9)	99.88, 17.9 (100, 8)

The value delimited by comma in each grid shows the classification accuracy and the corresponding number of selected features respectively. The values in the parentheses are the best results obtained. Bold typefaces emphasize the best average performance in each row. Bold italic typefaces emphasize the best solution found among all the methods. * Here WFFSA-R, WFFSA-G, and WFFSA-C denote WFFSA with local search heuristic of ReliefF, Gain Ratio, and Chi-Square, respectively. ** First represents the improvement first strategy with length $l = 4$ and interval $w = 1$. Greedy represents greedy strategy with $l = 4$ and $w = 5$.

TABLE 3
RESULTS BY WFFSAS AND OTHER METHODS ON UCI DATASETS

Dataset	WFFSA		Results obtained from the literature	
Vehicle	75.06, 12 <i>(75.06, 12)</i>	WFFSA+INN LOOCV	73.52, 7 <i>(73.52, 7)</i>	HGA(3)+INN LOOCV [14]
WDBC	98.14, 13.9 <i>(98.24, 12)</i>	WFFSA+INN LOOCV	93.85, 24 <i>(93.85, 24)</i>	HGA(3)+INN LOOCV [14]
Ionosphere	95.19, 8.9 <i>(96.01, 8)</i>	WFFSA+INN LOOCV	95.56, 7 <i>(95.73, 7)</i>	HGA(3)+INN LOOCV [14]
Sonar	96.30, 24 <i>(97.6, 22)</i>	WFFSA+INN LOOCV	96.34, 24 <i>(97.12, 24)</i>	HGA(3)+INN LOOCV [14]

Bold and bold italic typefaces represent best average performance and best solution found among the methods, respectively.

TABLE 4
RESULTS BY WFFSAS AND OTHER METHODS ON MICROARRAY DATASETS

Dataset	WFFSA		Results obtained from the literature	
ALL/AML	97.34, 30.5 <i>(98.57, 30.8)</i>	WFFSA+SVM 10-fold CV	97.34, 31 <i>(97.68, 50)</i>	SVM+LS bound .632+ bootstrap [23]
	100, 5.1 <i>(100, 2)</i>	WFFSA+1NN LOOCV	<i>(100, 8)</i>	SVM-RFE Holdout [8]
Colon	86.01, 31.0 <i>(87.38, 31.2)</i>	WFFSA+SVM 10-fold CV	84.77, 31 <i>(84.95, 46)</i>	SVM+LS bound .632+ bootstrap [23]
	97.9, 10.9 <i>(100, 14)</i>	WFFSA+1NN LOOCV	<i>(100, 16)</i>	SVM-RFE Holdout [8]
SRBCT	98.53, 90.5 <i>(100, 89.6)</i>	WFFSA+SVM 10-fold CV	<i>(~95.00, 150)</i>	SVM+Max Minority 4-fold CV [27]
	100, 15.1 <i>(100, 7)</i>	WFFSA+1NN LOOCV	<i>(100, 96)</i>	ANN+PCA 3-fold CV [26]
NCI60	66.43, 485.3 <i>(72.19, 485)</i>	WFFSA+SVM 4-fold CV	<i>(66.66, 150)</i>	SVM+Sum Minority 4-fold CV [27]
	82.30, 22.7 <i>(85.25, 14)</i>	WFFSA+1NN LOOCV	<i>(85.25, 13)</i>	GA+MLHD LOOCV [25]

Bold and bold italic typefaces represent best average performance and best solution found among the methods, respectively. Non-shaded and shaded tables represent using Internal or External cross validation, respectively. External cross-validation is repeated 10 times for each dataset.

TABLE 5
RESULTS OF DIFFERENT LOCAL SEARCH STRATEGIES, LENGTH AND INTERVAL
ON SONAR DATASET

	Improvement First Strategy			Greedy Strategy		
	$w = 1$	$w = 5$	$w = P$	$w = 1$	$w = 5$	$w = P$
$l = 2$	95.77, 26.3	94.95, 26.7	95.10, 26.1	94.86, 28.3	95.14, 25.3	94.52, 28.3
$l = 4$	96.30, 24	95.67, 24.7	95.38, 22.9	94.8, 23.2	95.63, 24.6	94.13, 26
$l = 8$	95.53, 21.4	95.24, 24.4	95.43, 25.5	94.91, 20.5	95.10, 25	94.04, 28.1

TABLE 6
RESULTS OF DIFFERENT LOCAL SEARCH STRATEGIES, LENGTH AND INTERVAL
ON COLON DATASET

	Improvement First Strategy			Greedy Strategy		
	$w = 1$	$w = 5$	$w = P$	$w = 1$	$w = 5$	$w = P$
$l = 2$	97.26, 17.9	97.10, 16.9	95.97, 22.2	96.13, 15	96.29, 14.7	95.65, 20.8
$l = 4$	97.9, 11.5	97.42, 12.9	97.10, 14.6	96.45, 12.6	97.26, 16	95.65, 19.9
$l = 8$	96.77, 10.5	97.26, 10.4	97.10, 13.1	96.45, 10.9	96.94, 14.1	95.32, 15.8

TABLE 7
RESULTS OF DIFFERENT LOCAL STRATEGIES ON ALL EIGHT DATASETS

Dataset	WFFSA-R First $w = 1, l = 4$	WFFSA-R Greedy $w = 5, l = 4$	MA-S* $w = 1, l = 4$	MA-S* $w = 5, l = 4$
Vehicle (846×18)	75.00, 11.25 (75.06, 12)	75.06, 12 (75.06, 12)	74.84, 10.5 (75.06, 12)	74.60, 10.6 (75.06, 12)
WDBC (569×30)	97.96, 13 (98.24, 14)	98.14, 13.9 (98.24, 14)	97.80, 14.7 (98.07, 14)	97.47, 13.7 (97.72, 12)
Ionosphere (351×34)	95.00, 7.5 (95.73, 8)	94.96, 9.7 (95.44, 8)	94.27, 11 (95.44, 10)	93.48, 12.8 (94.59, 12)
Sonar (208×60)	96.30, 24 (97.12, 19)	95.63, 24.6 (96.63, 24)	93.65, 26.5 (95.67, 25)	91.73, 28.4 (93.75, 32)
ALL/AML (72×1000)	100, 5.1 (100, 2)	100, 8.3 (100, 4)	100, 27.2 (100, 16)	100, 44 (100, 44)
Colon (62×1000)	97.9, 11.5 (100, 14)	97.26, 16 (100, 20)	91.94, 31.4 (95.16, 19)	82.26, 18 (82.26, 18)
NCI60 (60×1000)	82.30, 22.7 (85.25, 19)	79.67, 23.8 (81.97, 24)	74.36, 36 (77.05, 36)	68.85, 34 (68.85, 34)
SRBCT (83×2308)	100, 15.1 (100, 7)	100, 18.7 (100, 11)	93.13, 39 (96.39, 33)	87.83, 35.5 (91.57, 23)

* MA-S denotes memetic algorithm based on sequential local search strategy.