

A Two-Stage Gene Selection Algorithm by Combining ReliefF and mRMR

Yi Zhang

School of Computing and
Information Sciences
Florida International University
Miami, Florida 33199
Email: yzhan004@cs.fiu.edu

Chris Ding

Department of Computer Science and Engineering
University of Texas
Arlington, TX 76019
Email: chqding@uta.edu

Tao Li

School of Computing and
Information Sciences
Florida International University
Miami, Florida 33199
Email: taoli@cs.fiu.edu

Abstract—Gene expression data usually contains a large number of genes, but a small number of samples. Feature selection for gene expression data aims at finding a set of genes that best discriminate biological samples of different types. In this paper, we present a two-stage selection algorithm by combining ReliefF and mRMR: In the first stage, ReliefF is applied to find a candidate gene set; In the second stage, mRMR method is applied to directly and explicitly reduce redundancy for selecting a compact yet effective gene subset from the candidate set. We also perform comprehensive experiments to compare the mRMR-ReliefF selection algorithm with ReliefF, mRMR and other feature selection methods using two classifiers as SVM and Naive Bayes, on seven different datasets. The experimental results show that the mRMR-ReliefF gene selection algorithm is very effective.

Index Terms—Gene selection algorithms, reliefF, mRMR, mRMR-reliefF.

I. INTRODUCTION

Gene expression refers to the level of production of protein molecules defined by a gene. Monitoring of gene expression is one of the most fundamental approach in genetics and molecular biology. The standard technique for measuring gene expression is to measure the mRNA instead of proteins, because mRNA sequences hybridize with their complementary RNA or DNA sequences while this property lacks in proteins. The DNA arrays, pioneered in [4] [8], are novel technologies that are designed to measure gene expression of tens of thousands of genes in a single experiment. The ability of measuring gene expression for a very large number of genes, covering the entire genome for some small organisms, raises the issue of characterizing cells in terms of gene expression, that is, using gene expression to determine the fate and functions of the cells. The most fundamental of the characterization problem is that of identifying a set of genes and its expression patterns that either characterize a certain cell state or predict a certain cell state in the future [13].

When the expression dataset contains multiple classes, the problem of classifying samples according to their gene expression becomes much more challenging, especially when the number of classes exceeds five [18]. Moreover, the special characteristics of expression data adds more challenge to the classification problem. Expression data usually contains

a large number of genes (in thousands) and a small number of experiments (in dozens). In machine learning terminology, these datasets are usually of very high dimensions with undersized samples. In microarray data analysis, many gene selection methods have been proposed to reduce the data dimensionality [22].

Gene selection aims to find a set of genes that best discriminate biological samples of different types. The selected genes are “biomarkers”, and they form “marker panel” for analysis. Most gene selection schemes are based on binary discrimination using rank-based schemes [7], such as information gain, which reduces the entropy of the class variables given the selected attributes. In expression data, many gene groups interact closely and gene interactions are important biologically and may contribute to class distinctions. However, the majority of the rank-based schemes assume the conditional independence of the attributes given the target variable and are thus not effective for problems involving much feature interaction [15].

In this paper, we present a two-stage selection algorithm by combining ReliefF [15] and mRMR [19]. ReliefF, a general and successful attribute estimator, is able to effectively provide quality estimates of attributes in problems with dependencies between attributes. mRMR (minimal-redundancy-maximal-relevance) method selects genes that have the highest relevance with the target class and are also maximally dissimilar to each other. mRMR is computationally expensive. The integration of ReliefF and mRMR thus leads to an effective gene selection scheme. In the first stage, ReliefF is applied to find a candidate gene set. This filters out many unimportant genes and reduces the computational load for mRMR. In the second stage, mRMR method is applied to directly and explicitly reduce redundancy and select a compact yet effective gene subset from the candidate set. We perform comprehensive experiments to compare the mRMR-ReliefF selection algorithm with ReliefF, mRMR and other feature selection methods using two classifiers on seven different datasets. The experimental results show that the mRMR-ReliefF gene selection is very effective. The rest of the paper is organized as the follows: Section II discusses the related work; our mRMR-ReliefF gene selection algorithm is presented in

Section III; experimental results are presented in Section IV; and Finally Section V concludes.

II. RELATED WORK

Generally two types of feature selection methods have been studied in the literature: filter methods [12] and wrapper methods [11]. As pointed out in [24], the essential differences between the two methods are:

- (1) that a wrapper method makes use of the algorithm that will be used to build the final classifier while a filter method does not, and
- (2) that a wrapper method uses cross validation to compare the performance of the final classifier and searches for an optimal subset while a filter method uses simple statistics computed from the empirical distribution to select attribute subset.

Wrapper methods could perform better but would require much more computational costs than filter methods. Most gene selection schemes are based on binary discrimination using rank-based filter methods [7], such as t-statistics and information gain etc. The majority of the rank-based schemes assume the conditional independence of the attributes given the target variable and are thus not effective for problems involving much feature interaction [15].

In this paper, we present a mRMR-ReliefF selection algorithm by combining ReliefF and mRMR. ReliefF is able to effectively provide quality estimates of attributes in problems with dependencies between attributes and mRMR selects genes that have the highest relevance with the target class and are also maximally dissimilar to each other. The integration of ReliefF and mRMR thus leads to an effective gene selection scheme with much gene interaction.

III. MRMR-RELIEFF GENE SELECTION

Section III-A and Section III-B discuss ReliefF and mRMR respectively. Section III-C presents the mRMR-ReliefF selection algorithm.

A. ReliefF

ReliefF is a simple yet efficient procedure to estimate the quality of attributes in problems with strong dependencies between attributes [15]. In practice, ReliefF is usually applied in data pre-processing as a feature subset selection method.

The key idea of the ReliefF is to estimate the quality of genes according to how well their values distinguish between instances that are near to each other. Given a randomly selected instance Ins_m from class L , ReliefF searches for K of its nearest neighbors from the same class called nearest hits H , and also K nearest neighbors from each of the different classes, called nearest misses M . It then updates the quality estimation W_i for gene i based on their values for Ins_m , H , M . If instance Ins_m and those in H have different values on gene i , then the quality estimation W_i is decreased. On the other hand, if instance Ins_m and those in M have different values on the the gene i , then W_i is increased. The whole process is repeated n times which is set by users. The

```

Input: Gene variables and labels
Output:  $W$  for the gene rank
Set all weights  $W := 0$ ;
foreach Iteration  $n$  do
  Randomly select an instance  $Ins_m$ ;
  Find  $K$  nearest hits  $H$ ;
  foreach class  $c \neq Label_m$  do
    | from class  $c$  find  $K$  nearest misses  $M_c$ ;
  end
  foreach  $g_i$  do
    | Update  $W_i$ ;
  end
end

```

Fig. 1. The reliefF algorithm

algorithm is shown in Figure 1 and updating W_i can use Equation 1:

$$W_i = W_i - \frac{\sum_{k=1}^K D_H}{n \cdot K} + \sum_{c=1}^{C-1} P_c \cdot \frac{\sum_{k=1}^K D_{M_c}}{n \cdot K} \quad (1)$$

where n_c is the number of instances in class c , D_H (or D_{M_c}) is the sum of distance between the selected instance and each H (or M_c), P_c is the prior probability of class c .

Detailed discussions on ReliefF can be found in [15] and recently, it was shown that ReliefF is an on-line solution to a convex optimization problem, maximizing a margin-based algorithm [23].

B. mRMR

The mRMR (minimum redundancy maximum relevance) method [19] selects genes that have the highest relevance with the target class and are also minimally redundant, i.e., selects genes that are maximally dissimilar to each other. Given g_i which represents the gene i , and the class label c , their mutual information is defined in terms of their frequencies of appearances $p(g_i)$, $p(c)$, and $p(g_i, c)$ as follows.

$$I(g_i, c) = \iint p(g_i, c) \ln \frac{p(g_i, c)}{p(g_i)p(c)} dg_i dc \quad (2)$$

The Maximum-Relevance method selects the top m genes in the descent order of $I(g_i, c)$, i.e. the best m individual features correlated to the class labels.

$$\max_S \frac{1}{|S|} \sum_{g_i \in S} I(g_i, c) \quad (3)$$

Although we can choose the top individual genes using Maximum-Relevance algorithm, it has been recognized that "the m best features are not the best m features" since the correlations among those top features may also be high [5]. In order to remove the redundancy among features, a Minimum-Redundancy criteria is introduced

$$\min_S \frac{1}{|S|^2} \sum_{g_i, g_j \in S} I(g_i, g_j) \quad (4)$$

where mutual information between each pair of genes is taken into consideration. The minimum-redundancy maximum-relevance (mRMR) feature selection framework combines both optimization criteria of Eqs.(2,3).

A sequential incremental algorithm to solve the simultaneous optimizations of optimization criteria of Eqs.(2,3) is given as the following. Suppose set G represents the set of genes and we already have S_{m-1} , the feature set with $m-1$ genes. Then the task is to select the m -th feature from the set $\{G - S_{m-1}\}$. This feature is selected by maximizing the *single-variable relevance minus redundancy* function

$$\max_{g_j \in G - S_{m-1}} [I(g_i; c) - \frac{1}{m-1} \sum_{g_i \in S_{m-1}} I(g_j; g_i)] \quad (5)$$

The m -th feature can also be selected by maximizing the *single-variable relevance divided-by redundancy* function

$$\max_{g_j \in G - S_{m-1}} [I(g_i; c) / \frac{1}{m-1} \sum_{g_i \in S_{m-1}} I(g_j; g_i)] \quad (6)$$

C. mRMR-ReliefF Algorithm

As we mentioned before, ReliefF is a general and successful attribute estimator and is able to effectively provide quality estimates of attributes in problems with dependencies between attributes. However, ReliefF does not explicitly reduce the redundancy in selected genes. mRMR selects genes that have the highest relevance with the target class and are also maximally dissimilar to each other. However, mRMR is computationally expensive. For example, using the mRMR program provided in [19], we could not obtain results for several datasets with a large number of genes, e.g., ALL and GCM. The integration of ReliefF and mRMR thus leads to an effective gene selection scheme.

We can view the *quality estimation* W_i in ReliefF as maximizing the relevance score. Thus we can view the standard ReliefF algorithm as maximizing the relevance score:

$$\max_S \frac{1}{|S|} \sum_{g_i \in S} W_i \quad (7)$$

Thus our mRMR-ReliefF algorithm selection criteria becomes

$$\max_{g_j \in G - S_{m-1}} W_i - \frac{1}{m-1} \sum_{g_i \in S_{m-1}} |C(g_j, g_i)| \quad (8)$$

or

$$\max_{g_j \in G - S_{m-1}} W_i / \frac{1}{m-1} \sum_{g_i \in S_{m-1}} |C(g_j, g_i)| \quad (9)$$

where $C(g_j, g_i)$ is the Pearson correlation coefficient.

Our mRMR-ReliefF algorithm works as follows: In the first stage, ReliefF is applied to find a candidate gene set. This filters out many unimportant genes and reduces the computational load for mRMR. In the second stage, mRMR method is applied to directly and explicitly reduce redundancy and select a compact yet effective gene subset from the candidate set.

In our experiments, ReliefF is first used to choose 150 genes as the candidate set. from the all gene data. mRMR algorithm is then applied to select the final subset.

TABLE I
THE DATASET DESCRIPTION

Dataset	# Samples	# Genes	# Classes
ALL	248	12558	6
ARR	420	278	2
GCM	198	16063	14
HBC	22	3226	3
LYM	62	4026	3
MLL	72	12582	3
NCI60	60	1123	9

IV. EXPERIMENTS

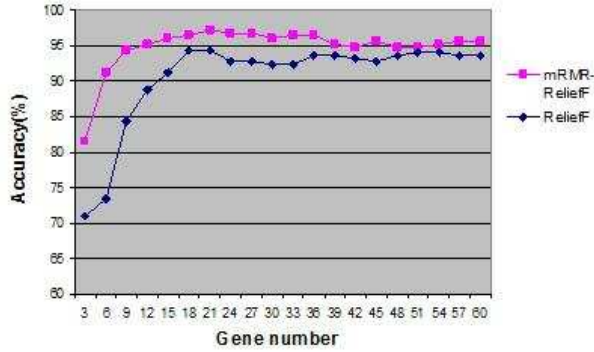
In this section, we perform comprehensive experiments to compare the mRMR-ReliefF selection algorithm with ReliefF, mRMR and other feature selection methods using two classifiers (Support Vector Machine (SVM) and Naive Bayes) on seven different datasets.

A. Datasets Description

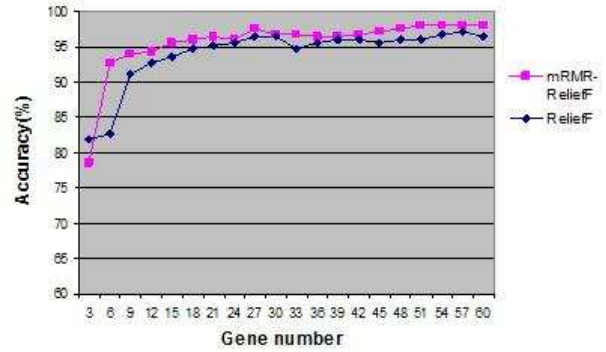
The datasets and their characteristics are summarized in Table I.

- **ALL:** The ALL dataset [25] is a dataset that covers six subtypes of acute lymphoblastic leukemia: BCR (15), E2A (27), Hyperdip (64), MLL (20), T (43), and TEL (79). Here the numbers in the parentheses are the numbers of samples. The dataset is available at [2].
- **ARR:** The Arrhythmia (ARR) dataset contains 420 samples and 278 features with two classes [3].
- **GCM:** The GCM dataset [20] consists of 198 human tumor samples of fifteen types. breast (12), prostate (14), lung (12), colorectal (12), lymphoma (22), bladder (11), melanoma (10), uterus (10), leukemia (10), renal (11), pancreas (11), ovary (120), mesothelioma (11), CNS (20), and MET (9). The prediction accuracy of 78% is reported in [20] using one-versus-the rest SVM with all the genes.
- **HBC:** The HBC dataset consists of 22 hereditary breast cancer samples and was first studied in [10]. The dataset has three classes and can be downloaded at [9].
- **LYM:** The Lymphoma dataset is a dataset of the three most prevalent adult lymphoid malignancies and available at [14] and it was first studied in [1].
- **MLL:** The MLL-leukemia dataset consists of three classes and can be downloaded at [16].
- **NCI60:** The NCI60 dataset was first studied in [21]. cDNA microarrays were used to examine the variation in gene expression among the 60 cell lines from the National Center Institute's anticancer drug screen. The dataset spans nine classes and can be downloaded at [9] [17].

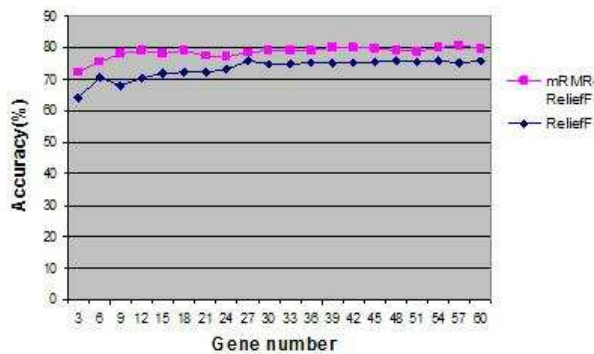
Note that in these datasets, the samples in each class is generally small, and unevenly distributed. This, together with the large number of classes, especially for NCI60, GCM, makes the classification task more complex.



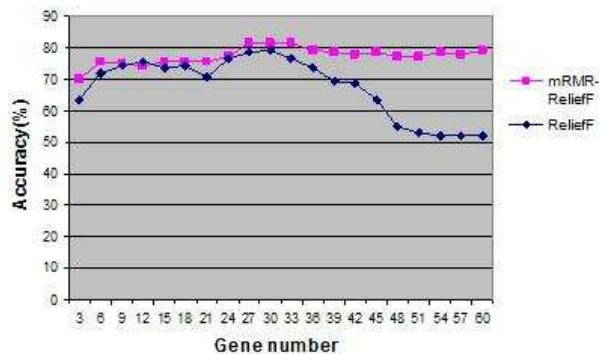
(a) Results of Naive Bayes: ALL dataset



(b) Results of SVM: ALL dataset



(c) Results of Naive Bayes: ARR dataset



(d) Results of SVM: ARR dataset

Fig. 2. Comparison of ReliefF and mRMR-ReliefF Algorithms (I)

B. Compare ReliefF, mRMR and mRMR-ReliefF algorithm

First we compare the mRMR-ReliefF algorithm with ReliefF and mRMR. We perform our comparisons using SVM and Naive Bayes classifiers on the seven datasets. Both SVM and Naive Bayes have been widely used in previous studies. Figure 2, Figure 3, and Figure 4 show the classification accuracy results as a function of the number of selected genes on the seven datasets respectively. In addition, because of mRMR is computationally expensive, using the program provided in [19], we could not obtain results for several datasets with a large number of genes, e.g., ALL and GCM. Thus in the figures, we only include the accuracy values for ReliefF and the mRMR-ReliefF algorithm and these values are all obtained via 10-fold cross validation.

Table II presents the detail of the accuracy values of applying SVM and Naive Bayes classification on the top 30 selected genes, for some unavailable results which can not be computed by mRMR, we note them as ”-”.

From the above comparative study, we observe that:

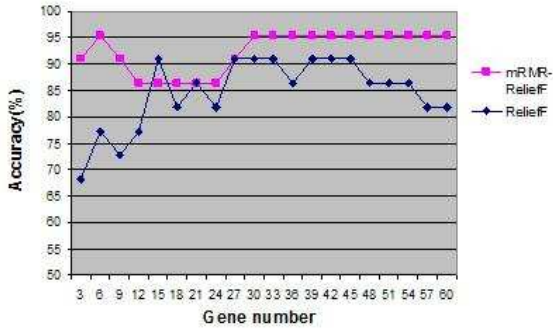
- The performance of mRMR algorithm is pulled down by its expensive computational cost, and it can not fulfill gene selection on the database with large features using

the limited memory.

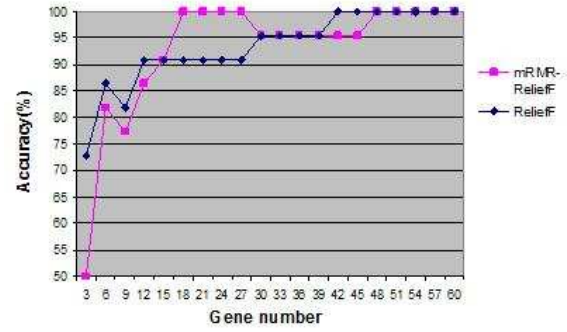
- Relief algorithm is not stable enough when only a small number of genes are selected. And when the number of selected genes is greater than 30, the variations of classification performance of both ReliefF and mRMR-ReliefF algorithms are generally small.
- The mRMR-ReliefF selection algorithm leads to significantly improved class predictions. With the same number of selected genes, the gene set obtained by the mRMR-ReliefF selection is more representative of the target class, therefore leading to better class prediction or generalization property.

C. Comparison with Other Methods

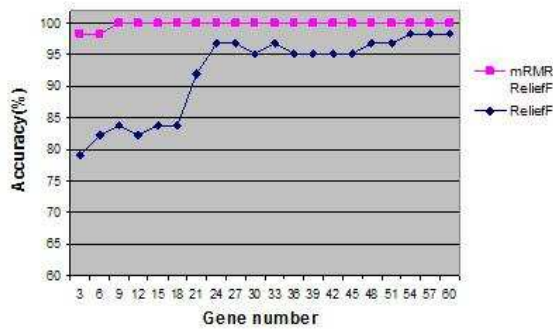
1) *Methods Description:* We also compare our mRMR-ReliefF selection algorithm with other gene selection algorithms, including Max-Relevance, Information Gain, Sum Minority, Twoing Rule, F-statistic [6], and GSNR [26]. These methods have been reported in previous work. The first four methods have been used either in machine learning (information gain) or in statistical learning theory (twoing rule and sum minority), and all of them measure the effectiveness of a



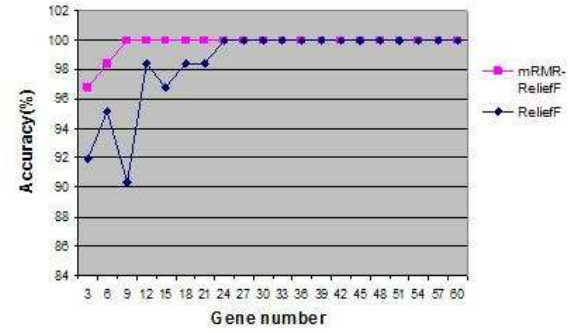
(a) Results of Naive Bayes: HBC dataset



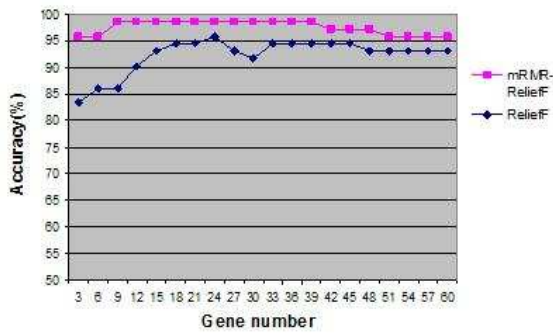
(b) Results of SVM: HBC dataset



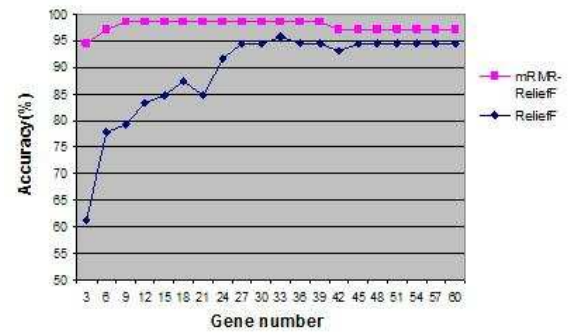
(c) Results of Naive Bayes: Lymphoma dataset



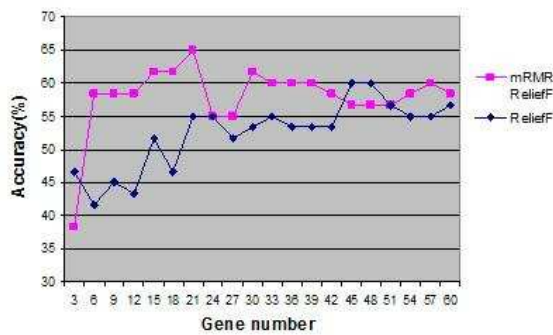
(d) Results of SVM: Lymphoma dataset



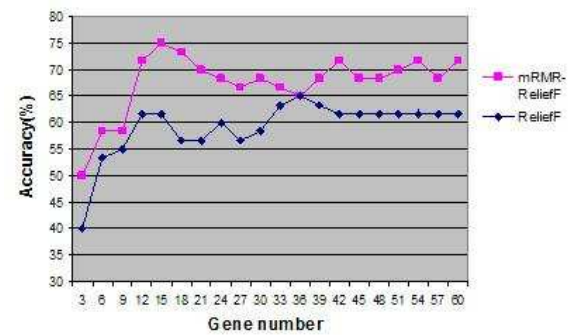
(e) Results of Naive Bayes: MLL dataset



(f) Results of SVM: MLL dataset

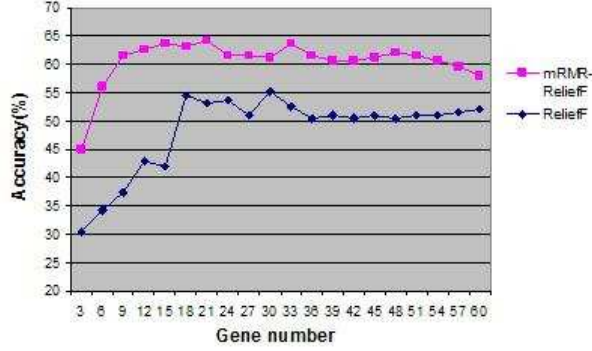


(g) Results of Naive Bayes: NCI60 dataset

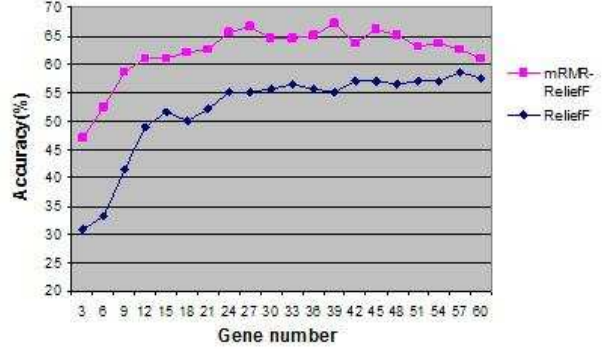


(h) Results of SVM: NCI60 dataset

Fig. 3. Comparison of ReliefF and mRMR-ReliefF Algorithms (II)



(a) Results of Naive Bayes: GCM dataset



(b) Results of SVM: GCM dataset

Fig. 4. Comparison of ReliefF and mRMR-ReliefF Algorithms (III)

TABLE II
THE COMPARISON IN RELIEFF, MRMR AND MRMR-RELIEFF ALGORITHMS (GENE NUMBER = 30)

Feature Selection Method	Classifier	ALL	ARR	LYM	HBC	NCI60	MLL	GCM
ReliefF	SVM	96.37%	79.29%	100%	95.45%	58.33%	94.44%	55.25%
	Naive Bayes	92.34%	75%	95.16%	90.91%	53.33%	91.67%	55.56%
mRMR	SVM	-	75.35%	100%	95.45%	53.33%	-	-
	Naive Bayes	-	73.21%	97.33%	87.51%	51.20%	-	-
mRMR-ReliefF	SVM	96.77%	81.43%	100%	95.45%	68.33%	98.61%	64.65%
	Naive Bayes	95.97%	79.05%	100%	95.45%	61.67%	98.61%	61.11%

feature by evaluating the strength of class prediction when the prediction is made by splitting it into two regions, the high region and the low region, by considering all possible split points [22]. More detailed descriptions on these methods can be found in [22].

F-statistic is chosen to score the relevance between the genes and the classification variable. The F-statistic of gene i in C classes has the following form [6] :

$$W_i = \frac{\sum_{c=1}^C n_c \cdot (\bar{g}_{ic} - \bar{g}_i) / (C - 1)}{\sum_{c=1}^C \{(n_c - 1) [\sum_{i=1}^{n_c} (g_{jic} - \bar{g}_{ic})^2 / n_c] / (n - C)\}} \quad (10)$$

where C is the number of classes, \bar{g}_i is the mean of gene i variables, n_c is the number of samples in class c , \bar{g}_{ic} is the mean of gene i in class c , and g_{jic} is sample j in gene i value in class c .

As to GSNR, it has been proposed and used in [26]. GSNR is a measure of the ratio between inter-group and intra-group variations. Higher GSNR values indicate higher discrimination power for the gene. The GSNR value for gene i is given by:

$$W_i = \frac{\sum_{c=1}^C |\bar{g}_{jc} - \sum_{c=1}^C \bar{g}_{jc} / C| / C}{\sum_{i=1}^C \sum_{i=1}^{n_c} |g_{jic} - \bar{g}_{ic}| / n_c} \quad (11)$$

Both F-statistic and GSNR select m genes in the descent order of W_i , and the best subset of genes is satisfied the following description:

$$\max_S \frac{1}{|S|} \sum_{g_i \in S} W_i \quad (12)$$

2) *Results Analysis:* Table III presents the classification accuracy comparison using SVM and Naive Bayes classifier when the number of selected gene is 30. From Table III, we observe that:

- Gene selection improves class prediction. Note that the accuracy of SVM using feature selection generally outperforms that without feature selection. This implies that feature selection can effectively reduce the insignificant dimensions and noise to improve classification accuracy.
- The mRMR-ReliefF algorithm is shown to achieve better performance comparing with other gene selection algorithms on almost all datasets. The experimental compar-

TABLE III
THE COMPARISONS IN SEVEN METHODS (GENE NUMBER = 30)

Feature Selection Method	Classifier	ALL	ARR	LYM	HBC	NCI60	MLL	GCM
No feature sel	SVM	91.94%	51.04%	95.16%	77.27%	63.33%	97.22%	51.52%
	Naive Bayes	85.23%	49.57%	95.04%	70.11%	45.22%	93.13%	40.33%
mRMR-ReliefF	SVM	96.77%	81.43%	100%	95.45%	68.33%	98.61%	64.65%
	Naive Bayes	95.97%	79.05%	100%	95.45%	61.67%	98.61%	61.11%
Maxrel	SVM	89.11%	74.53%	100%	72.73%	51.67%	77.78%	60.61%
	Naive Bayes	88.71%	73.49%	100%	63.64%	48.33%	80.56%	46.97%
Information Gain	SVM	97.58%	80.13%	98.39%	100%	61.67%	98.67%	46.67%
	Naive Bayes	92.74%	77.21%	93.55%	86.38%	60%	97.22%	47.47%
Sum Minority	SVM	93.95%	76.42%	98.39%	95.45%	55%	90.28%	55.05%
	Naive Bayes	91.13%	74.32%	95.16%	81.82%	46.67%	91.67%	49.49%
Twoing Rule	SVM	96.77%	79.37%	98.39%	90.91%	61.67%	97.22%	45.96%
	Naive Bayes	90.32%	72.19%	93.55%	86.36%	45%	95.83%	46.46%
F-statistic	SVM	97.17%	67.12%	96.77%	90.91%	63.33%	77.22%	39.10%
	Naive Bayes	80.27%	71.55%	98.52%	85.41%	60.15%	80.13%	39.81%
Gsnr	SVM	93.18%	77.24%	100%	95.45%	63.37%	90.25%	40.74%
	Naive Bayes	90.11%	70.43%	100%	85.65%	58.25%	87.22%	39.81%

isons demonstrate the effectiveness of the integration of ReliefF and mRMR.

- ReliefF achieves good performance on most of the data sets. Although its performance is not always as good as that of the mRMR-ReliefF algorithm. It outperforms mRMR, Maxrel, Sum Minority and partially wins information gain, twoing rule.
- Only a small number of genes are needed for classification purpose. In our experiments, the variations of the classification accuracy are small when the number of selected genes is greater than 30.

V. CONCLUSION

In this paper, we present an mRMR-ReliefF selection algorithm by combining ReliefF and mRMR. ReliefF is able to effectively provide quality estimates of attributes in problems with dependencies between attributes and mRMR method selects genes that have the highest relevance with the target class and are also maximally dissimilar to each other. The integration of ReliefF and mRMR thus leads to an effective gene selection scheme: In the first stage, ReliefF is applied to find a candidate gene set; In the second stage, mRMR is applied to select a compact yet effective gene subset from the candidate set. Comprehensive experiments are conducted to compare the mRMR-ReliefF selection algorithm with ReliefF, mRMR and other feature selection methods using two classifiers on seven different datasets. The experimental results show that the mRMR-ReliefF gene selection is very effective.

ACKNOWLEDGMENT

Tao Li is partially supported by a IBM Faculty Research Award, NSF CAREER Award IIS-0546280 and NIH/NIGMS S06 GM008205. Chris Ding is supported in part by a University of Texas STARS Award. We would like to thank Dingding Wang for assisting with the experiments on several gene selection algorithms. We are also grateful to the anonymous reviewers for their helpful comments.

REFERENCES

- [1] A. A. Alizadeh, M. B. Eisen, R. E. David, C. Ma, I. S. Lossos, A. R. osenwald, H. C. Boldrick, H. Sabet, T. Tran, X. Yu, J. I. Powell, L. Yang, G. E. Martu, T. Moore, J. Hudson, L. Lu, D. B. Lewis, R. Tibshirani, G. Sherlock, W. C. Chan, T. C. Greiner, D. D. Weisenburger, G. P. Armitage, R. Warnke, R. Levy, W. Wilson, M. R. Grever, J. C. Byrd, D. Botsten, P. O. Brown, and L. M. Staudt. Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling. *Nature*, 403:503–511, 2000.
- [2] <http://www.stjude.com/research/data/ALL1/>.
- [3] <http://www.ics.uci.edu/ml/learn/MLSummary.html>.
- [4] M. Chee, R. Yang, E. Hubbell, A. Berno, X. Huang, D. Stern, J. Winkler, D. Lockhart, M. Morris, and S. Fodor. Accessing genetic information with high density DNA arrays. *Science*, 274:610–614, 1996.
- [5] T. Cover. The best two independent measurements are not the two best. *IEEE Trans. Systems, and Cybernetics*, 4:116–117, 1974.
- [6] C. Ding, H. Peng. Minimum redundancy feature selection from microarray gene expression data. *CSB 03*, 2003.
- [7] S. Dudoit, J. Fridlyand, and T. P. Speed. Comparison of discrimination methods for the classification of tumors using gene expression data. *Journal of the American Statistical Association*, 97(457):77–87, 2002.
- [8] S. Fodor, J. Read, M. Pirrung, L. Stryer, A. Lu, and D. Solas. Light-directed, spatially addressable parallel chemical synthesis. *Science*, 251:767–783, 1991.
- [9] <http://www.columbia.edu/~xy56/project.htm>.
- [10] I. Hedenfalk, D. Duggan, Y. C. Y. M. Radmacher, M. Bittner, R. Simon, P. Meltzer, B. Gusterson, M. Esteller, O. P. Kallioniemi, B. W. B. A. Borg, and J. Trent. Gene-expression profiles in hereditary breast cancer. *The New England Journal of Medicine*, 344(8):539–548, 2001.
- [11] R. Kohavi and G. H. John. Wrappers for feature subset selection. *Artificial Intelligence*, 97(1-2):273–324, 1997.
- [12] P. Langley. Selection of relevant features in machine learning. In *AAAI Fall Symposium on Relevance*, pages 140–144, 1994.
- [13] 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*, 20(15):2429–2437, 2004.
- [14] <http://genome-www.stanford.edu/lymphoma>.
- [15] R. S. Marko and K. Igor. Theoretical and empirical analysis of relief and reliefF. *Machine Learning Journal*, 53:23–69, 2003.
- [16] <http://research.dfci.harvard.edu/korsmeyer/MLL.htm>.
- [17] <http://genome-www.stanford.edu/nci60/>.
- [18] C. Ooi and P. Tan. Genetic algorithms applied to multi-class prediction for the analysis of gene expression data. *Bioinformatics*, 19(1):37–44, 2003.
- [19] H. Peng, F. Long, and C. Ding. Feature selection based on mutual information: Criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 27, 2005.

- [20] S. Ramaswamy, P. Tamayo, R. Rifkin, S. Mukherjee, C.-H. Yeang, M. Angelo, C. Ladd, M. Reich, E. Latulippe, J. P. Mesirov, T. Poggio, W. Gerald, M. Loda, E. S. Lander, and T. R. Golub. Multiclass cancer diagnosis using tumor gene expression signatures. *98(26):15149–15154*, 2001.
- [21] D. T. Ross, U. Scherf, M. B. Eisen, C. M. Perou, C. Rees, P. Spellmand, V. Iyer, S. S. Jeffrey, M. Van de Rijn, M. Waltham, A. Pergamenschikov, J. C. F. Lee, D. Lashkari, D. Shalon, T. G. Myers, J. N. Weinstein, D. Botstein, and M. P. O. Brown. Systematic variation in gene expression patterns in human cancer cell lines. *Nature Genetics*, 24:227–235, 2000.
- [22] Y. Su, T. M. Murali, V. Pavlovic, and S. Kasif. Rankgene: Identification of diagnostic genes based on expression data. *Bioinformatics*, 2003. The program can be downloaded from <http://genomics10.bu.edu/yangsu/rankgene/>.
- [23] Y. Sun and J. Li. Iterative RELIEF for feature weighting. In *Proceedings of the 23rd International Conference on Machine Learning*, 2006.
- [24] E. P. Xing, M. I. Jordan, and R. M. Karp. Feature selection for high-dimensional genomic microarray data. In *Proc. 18th International Conf. on Machine Learning*, pages 601–608. Morgan Kaufmann, San Francisco, CA, 2001.
- [25] E.-J. Yeoh, M. E. Ross, S. A. Shurtleff, W. K. Williams, D. Patel, R. Mahrouz, F. G. Behm, S. C. Raimondi, M. V. Relling, A. Patel, C. Cheng, D. Campana, D. Wilkins, X. Zhou, J. . Li, H. Liu, C.-H. Pui, W. E. Evans, C. Naeve, L. Wong, and J. R. Downing. Classification, subtype discovery, and prediction of outcome in pediatric lymphoblastic leukemia by gene expression profiling. *Cancer Cell*, 1(2):133–143, 2002.
- [26] G. Zheng. Statistical analysis of biomedical data with emphasis on data integration. *Dissertation in Florida International University*, 2006.