Global Journals $end{transformula} {\end{transformula}} {\end{transfor$

Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

Discriminative Gene Selection Employing Linear Regression Model

Abid Hasan¹, Shaikh Jeeshan Kabeer² and Kamrul Hasan³

¹ Islamic University of Technology (IUT), Dhaka, Bangladesh

Received: 16 December 2012 Accepted: 3 January 2013 Published: 15 January 2013

7 Abstract

3

Δ

5

Microarray datasets enables the analysis of expression of thousands of genes across hundreds 8 of samples. Usually classifiers do not perform well for large number of features (genes) as is 9 the case of microarray datasets. That is why a small number of informative and discriminative 10 features are always desirable for efficient classification. Many existing feature selection 11 approaches have been proposed which attempts sample classification based on the analysis of 12 gene expression values. In this paper a linear regression based feature selection algorithm for 13 two class microarray datasets has been developed which divides the training dataset into two 14 subtypes based on the class information. Using one of the classes as the base condition, a 15 linear regression based model is developed. Using this regression model the divergence of each 16 gene across the two classes are calculated and thus genes with higher divergence values are 17 selected as important features from the second subtype of the training data. The classification 18 performance of the proposed approach is evaluated with SVM, Random Forest and AdaBoost 19 classifiers. Results show that the proposed approach provides better accuracy values compared 20 to other existing approaches i.e. ReliefF, CFS, decision tree based attribute selector and 21 attribute selection using correlation analysis. 22

23

24 Index terms—linear regression, feature selection, microarray dataset, classification.

²⁵ 1 Introduction

he explosive growth and developments of microarray applications have enabled biologists and data mining 26 engineers to study and observe thousands of gene expression data at the same time. Various attribute selection 27 methodologies have been applied in the field of microarray data and in this particular case it is termed as gene 28 selection as illustrated in [1]. Microarray data analysis has paved the way to cancer, tumor and other disease 29 classification methods that can be used for subsequent diagnosis or prognosis. The problem of microarray data 30 are many fold, firstly not all the data are relevant and often only a small portion of the data is related to the 31 purpose of interest moreover noise and inconsistent data are prominent which hampers the search for the best 32 genes for selection and classification [2]. However the major difficult aspect of microarray data is that the genes 33 34 numbering in the thousands far outweighs the number of samples number in the lower hundreds if not less. This 35 makes the task of building effective models particularly difficult and poses over fitting problems where the model 36 does not perform well for novel patterns [3]. Thus feature selection methods being developed should be efficient 37 in handling these issues. Feature selection techniques can be generally divided into two broad categories depending on how the selection 38

process interacts with classification model ??4]. The first is the filter method where the importance of a feature is determined by scoring all the features based on their inherent attribute and retaining a portion of the features with higher scores while the low scoring features are removed as shown in many works including [5] and [6].

42 Filter methods are simple, fast and they do not require consultation with the classifier however the most obvious

43 drawback is that it examines each feature individually and hence cannot harness the combined predictive power 44 of features The second feature selection methodology is the wrapper model where a classification model is built

of features The second feature selection methodology is the wrapper model where a classification model is built by using a set of training set of features whose class labels are known and then the search for the optimal subset of

features is done by repeatedly generating and evaluating possible feature using against the well known classifiers

47 [7]. As the search for the solution is built into the classification process and as it considers the combinative

48 predicting power of gene subsets the convergence time is higher the methods are usually complex.

Studies such as [8] and [9] have shown that the biological state of individuals is defined by their gene expression 49 values. Therefore genes which have different expression profiles are more likely to properly identify biological 50 states than genes having similar expression profiles. In this paper a linear regression model is proposed where 51 one class of training dataset is considered as the base condition and generates the regression coefficients for each 52 of the genes in the base class. Using the regression coefficients of the base condition a regression representation 53 for the other class is generated and the difference in expression profiles between the genes of the base and non-54 base classes are measured. Genes with higher difference in expression profiles are given more importance and 55 scoring of genes are generated. The base class serves as domain knowledge that is used to guide the search for 56 discriminating genes in the dataset, [10] and [11] Year the best features. The detailed procedure of the proposed 57 method is provided in the following section. In the simulation and result analysis section it is seen that very high 58 59 classification accuracy rates are achieved using only a very small number of the genes and the proposed method 60 generated better results compared to other filtering approaches. The proposed approach has been applied on 6 61 microarray datasets and their effectiveness was determined by testing them in three different types of classifiers: Support Vector Machine (SVM), Random Forest and AdaBoost. 62

This paper is divided into 4 sections with section 1 giving an overview of the working domain and very brief introduction to the proposed approach. Section 2 elaborates the proposed approach in detail. Section 3 covers the simulation and result analysis part of the research where the proposed method is compared with Relief F, CFS, Chi-Squared value and Gain Ratio; it is seen that the proposed approach performs better then these existing methods. Section 4 provides the conclusion and provides scope for further research or development of this research work.

69 2 II.

70 3 Proposed Approach a) Theoretical Background

⁷¹ Linear regression is a statistical approach that can be used for predicting and forecasting. It has been traditionally ⁷² used to model relationships between a set of explanatory variables $12 \{ , ..., \}$ n A a a =

and the output variable b x. The idea is to derive a model using which the predictor or the output variable can be estimated using the explanatory variables [12]. In traditional feature selection applications the set of features are the input variables and the class labels are the output variables. Considering one feature a, the hypothesis function for this simple linear regression is $0 \ 1 \ x \ b \ x \ a = +(1)$

⁷⁷ where o x and 1

x are the parameters and x b is the predictor variable. The objective is to find the values of the parameters so that it best fits the data in the training set such that the features of unknown samples can be used for classification. j x should be chosen such that () x b a is as close to the training data (,) a b such that the following cost function is minimized? = ? = m i i i x b a b m x x F 1 2 1 0)) ((21), ((2))

82 Here) 1, 0 (x x F

is the cost function and m is the total number of samples in the training dataset. It is apparent that real world applications will require consideration of more than one feature and hence the hypothesis function will become $1 \ 1 \ 2 \ 02 \ 3 \ 3 \) \dots x n n b a x x a x a x a x a x a = + + + + + (3)$

x b a can now be written as() T b a X A x = (4)

T X is the transpose of parameter vector and A is the vector of explanatory variables. So the corresponding cost function) (X F which needs to be minimized for multiple variables, is the following:? = ? = m i i i x n b a b m x x x x F 1 2 2 1 0)) ((2 1),...,, ((5))

Gradient descent is a very popular approach that has been used in many researches including linear regression. From earlier discussion it is clear that the idea is to minimize the cost function) (X F . Gradient descent algorithm helps find the parameter value which leads to the minimum cost. The representation of equation 5 in partial derivative term is? = ? + = m i i j i i x j j a b a b m x x 1)) ((1 : ? (6)

The algorithm starts with an arbitrary value j x and keeps on changing by simultaneously updating x for n j ,...,2 , 1 , 0 =

99 until convergence for each of the j x occurs.

¹⁰⁰ 4 b) Linear Regression on Microarray Dataset

Linear regression is a statistical approach that can be used for microarray datasets provides gene expression values for different samples. Using gene expression values to find out features and hence to classify novel samples 103 is a common approach; however the application of linear regression to this task is a relatively fresh approach. In

this proposed method the gene expression values of one class of samples of a two class microarray training dataset is used as the base class. Using this portion of dataset, a model is built which acts as the domain knowledge of the dataset

106 the dataset.

107 5 Global

108) (i x b?

¹¹² represent the statistical values of expression for each gene in the non-base subtype of the training dataset.

113 6 c) Proposed Algorithm

In our proposed method, basic idea of linear regression has been used. We have tried to predict a potential feature 114 115 from one of the subtypes of microarray training datasets using the knowledge acquired from the other subtype 116 of the same training dataset. At first the microarray dataset is divided into two segments test and training dataset in the similar way as most supervised learning algorithm does. One of the biggest problems of microarray 117 data; redundancy has been handled by measuring the similarity in expression values of the genes in both types. 118 We have eliminated those genes having similar expression values considering their ineffectiveness as important 119 features for classification. Moreover, removing these genes gives the algorithm an efficient way of starting feature 120 selection procedure. Training samples are then divided into two subtypes: base type and non-base type, built 121 based on their class information. Next the parameter vector X for 1 S is generated using equation 6 and from 122 the parameter vector X,) (a x b? is calculated for 2 S. After the divergences and the differences are calculated, 123 genes are sorted according to difference values in the descending order. From the sorted list of genes) 100 ,..., 124 30, 20, 10 (= N N highest ranked genes are chosen and their classification accuracy is evaluated using different 125 classifiers. Section 3 shows the detailed performance evaluation of the proposed approach and its superiority 126 compared to other existing feature selection methods. 127

128 **7** III.

129 8 Materials and Methods

To find out how the proposed algorithm works, we have established the experiments using four different microarray datasets. We have compared our proposed feature selection algorithm with several other attribute selection procedures. Following sections describe a short description of microarray datasets and performance evaluations of the proposed method.

134 9 a) Datasets

The datasets are obtained from different authors. Datasets are converted into convenient way for this particular research.

The original prostate dataset was used in [13]. The dataset contains the 12,533 gene expression measurements of 102 samples. 50 of these 102 samples contain normal tissues not containing prostate tumor while 52 had prostate tumor.

Prostate cancer dataset was originally taken from dataset GSE2443 [14]. The dataset contains 12,627 gene expression values of 20 samples. Among them 10 samples contain androgen dependent tumor while other 10 contain androgen-independent tumor.

The lung cancer dataset contains two types of cancer: malignant pleural mesothelioma (MPM) and adenocarcinoma (ADCA) of lung. Among 181 tissue samples, 31 of them had MPM and 150 of them had ADCA. Each of the samples was described by 12,533 gene expression value [15].

The colon dataset was used in [16]. The dataset contains 62 samples collected from colon cancer patients. 40 tumor biopsies are from tumors and 22 normal biopsies are from healthy parts of colon of same patients. The number of genes used in this expression is around 2000.

¹⁴⁹ 10 b) Performance Evaluation

The implementation of the proposed algorithm of feature selection was done on MATLAB and the performance evaluation of the selection set of features for classification was performed on publicly available weka tool [16]. We have used 10 fold cross-validation for SVM classifier. The random forest procedure was run with 10 trees and AdaBoost classifier uses 10 iteration and weighted threshold of 100.

154 11 i. Results

155 The classification accuracy of the features selected by proposed method and its comparison with Another aspect of

156 the proposed feature selection method is; we have not used any threshold for how many features for classification

will be selected. Several different subsets of features have been used for classification and thus select the best subset based on its classifying ability. Figure 1 shows the error rate in classification by the classifier with any particular feature subset. For a particular feature selection j using a particular classifier i, the average error rate is calculated With the increase of the number of features, the error rate is decreased for most all the datasets. However, for prostate cancer dataset, although the error rate increases with first few subsets of features but at the end it too shows the same characteristics as the other microarray datasets shows.

¹⁶³ 12 ii. Discussion

We have proposed a new approach of feature selection using linear regression analysis. The algorithm works twofold. At the initial stage of the algorithm, we have eliminated redundant gene by measuring the similarity in expression values. Linear regression analysis then applied on one subtype (base type) of the training dataset to build the regression model. This model then applied on the other subtype (non-base type) of the training dataset to find out the divergence of the expression values of genes in that subtype. The more deviation shown by the gene, the more important it is considered as a feature. This way set of features selected for classification of the datasets.

Our main focus in this study is to classify accurately with less number of features. Table 1-4 shows the superiority of the classification accuracy by the features selected by the proposed method for different classifiers. Although, for colon dataset, classification accuracy by the features selected by ReliefF and CFS approach shows better result than the proposed method. However the result is still comparable and the number of feature selected by the proposed method is considerably fewer than the other method of attribute selection. Also Figure 1 summarizes the effect of different feature subsets for classification.

177 **13 IV.**

178 14 Conclusion

Linear regression based feature selection shows promising results in classification of microarray datasets. The proposed approach might be applied on more microarray datasets and the results obtained might be used to

improve some of the parameters of the proposed method. The results will also help to understand the performance of the proposed approach on a broader scale. The proposed approach can also be extended for multiclass

of the proposed approach on a broader scale. The proposed approach can also be extended for multiclass
 approaches to be applied in other data mining domains. In the future Incorporation of other knowledge might
 help the proposed method to enhance the performance and significance of the result.



Figure 1:



Figure 2: C

1

				Classifiers				
Attribute	Attribute selection		SVM	Random For	Random Forest		AdaBoost	
method								
		Ν	Acc $(\%)$	Ν	Acc $(\%)$	Ν	Acc $(\%)$	
ReliefF		100	80.33	150	81.97	100	88.52	
CFS		17	67.21	17	59.02	17	67.21	
Chi-Squared value		17	68.85	17	60.66	17	65.57	
GainRatio Value		1190	83.61	1190	72.13	1190	85.24	
Proposed		30	90.16	20	86.66	50	98.36	
				Classifiers				
Attribute selection			SVM	Random Forest		AdaBoost		
method								
		Ν	Acc $(\%)$	Ν	Acc $(\%)$	Ν	Acc $(\%)$	
ReliefF		200	58.33	150	66.67	100	50.00	
CFS		44	58.33	44	25.00	44	33.33	
Chi-Squared value		44	58.33	44	66.67	44	41.67	
GainRatio Value		188	58.33	188	58.33	188	25.00	
Proposed		20	91.67	20	75.00	10	83.33	

Figure 3: Table 1 :

3

			Classifiers						
Attribute selection method	SVM	Random forest		AdaBoost					
	Ν	Acc $(\%)$	Ν	Acc $(\%)$	Ν	Acc $(\%)$			
ReliefF	100	93.96	200	93.96	100	97.98			
CFS	37	89.93	37	91.27	37	93.30			
Chi-Squared value	37	89.93	37	92.62	37	94.63			
GainRatio Value	705	91.27	705	95.30	705	97.31			
Proposed	30	97.98	40	97.98	30	97.98			
Table 4 : Colon dataset classification accuracy									
			Classifiers						
Attribute selection method		SVM	Random forest		AdaBoost				
	Ν	Acc $(\%)$	Ν	Acc $(\%)$	Ν	Acc (%)			
ReliefF	200	87.88	200	84.85	200	72.73			
CFS	8	69.7	8	66.67	8	57.58			
Chi-Squared value	8	81.82	8	69.70	8	63.64			
GainRatio Value	62	81.82	62	66.67	62	69.69			
Proposed	10	84.85	20	75.76	20	78.79			

Figure 4: Table 3 :

 $\mathbf{2}$

Figure 5: Table 2 :

 ¹© 2013 Global Journals Inc. (US) Year
 ²CDiscriminative Gene Selection Employing Linear Regression Model
 ³© 2013 Global Journals Inc. (US) Year
 ⁴© 2013 Global Journals Inc. (US)

14 CONCLUSION

- [Gordon et al.], G J Gordon, R V Jensen, L L Hsiao, S R Gullans, J E Blumenstock, S R Ramaswamy, W.
- [Duval and Hao (2009)] 'Advances in meta heuristics for gene selection and classification of microarray data'.
 Beatrice Duval, Jin-Kao Hao. Briefings in Bioinformatics 2009. July 2009. 11 (1) p. .
- [Alon et al. ()] 'Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon
 tissues probed by oligonucleotide arrays'. U Alon, N Barakai, D Notterman, K Gish, S Ybarra, D Mack.
 Proceedings of National Academy of Science, (National Academy of Science) 1999. 96 p. .

[Hall ()] Correlation-based feature selection for machine learning, M Hall . 1999. Department of Computer
 Science, Waikato University (PhD Thesis)

- [Witten and Frank ()] Data Mining: Practical machine learning tools and techniques, I H Witten , E Frank .
 2005. Morgan Kaufmann Publisher.
- [Yu and Liu ()] 'Efficient feature selection via analysis of relevance and redundancy'. L Yu , H Liu . J. Mach.
 Learn. Res 2004. 2004. 5 p. .
- Inza et al. (2004)] 'Filter versus wrapper gene selection approaches in DNA microarray domains'. Iñaki Inza ,
 Pedro Larrañaga , Rosa Blanco , Antonio J Cerrolaza . Artificial Intelligence in Medicine 2004. June 2004.
 31 (2) p. .
- [Singh et al. ()] 'Gene expression correlated of clinical prostate cancer behavior'. D Singh , P G Febbo , K Ross
 , D G Jackson , J Manola , C Ladd , P Tamayo , A A Renshaw , A V D'amico , J P Richie , E S Lander ,
 M Loda , P W Kantoff , T R Golub , W R Seller . *Cancer cell* 2002. 1 (2) p. .
- [Inza et al. ()] 'Gene Selection by Sequential Search Wrapper Approaches in Microarray Cancer Class Prediction'.
 Iñaki Inza , Basilio Sierra , Rosa Blanco , Pedro Larrañaga . Journal of Intelligent & Fuzzy Systems 2002. 12
 (1) p. .
- [Yu et al. ()] 'Incorporating Prior Domain Knowledge into a Kernel Based Feature Selection Algorithm'. Ting
 Yu , Simeon J Simoff , Donald Stokes . Lecture Notes in Computer Science 2007. 2007. 4426 p. .
- 207 [Kohavi and John (1997)] Kohavi, G H John. Wrappers for Feature Subset Selection, 1997. Dec. 1997. 97 p. .
- 208 [Yan and Su ()] Linear Regression Analysis, X Yan , X Su . 2009. World Scientific.
- [Best et al. ()] 'Molecular alternations in primary prostate cancer after androgen ablation therapy'. C J Best , J
 W Gillespie , Y Yi , G V Chandramouli . *Clin cancer Res* 2005. 1 (11) p. .
- [Golub et al. (1999)] 'Molecular classification of cancer: class discovery and class prediction by gene expression
 monitoring'. T R Golub , D K Slonim , P Tamayo , C Huard , M Gaasenbeek , J P Mesirov , H Coller , M L
 Loh , J R Downing , M A Caligiuri , C D Bloomfield , E S Lander . Science 1999. Oct 15. 286 (5439) p. .
- [Barzilay and Brailovsky1999 (1999)] 'On domain knowledge and feature selection using a support vector
 machine'. Ofir Barzilay, V L Brailovsky1999. Pattern Recognition Letters May 1999. 20 (5) p.
- 216 [Sørlie et al. ()] Repeated Observation of breast tumor subtypes in independent gene expression data sets, Therese
- Sørlie , Robert Tibshirani , Joel Parker , Trevor Hastie , J S Marron , Andrew Nobel , Shibing Deng , Hilde
- Johnsen, Robert Pesich, Stephanie Geisler, Janos Demeter, Charles M Perou, Per E Lønning, Patrick O
 Brown, Anne-Lise Børresen-Dale, David Botstein. 2003. PNAS. 100 p. .
- [Richards et al. ()] 'Translation of microarray data into clinically relevant cancer diagnostic tests using gene expression ratios in lung cancer and meothelioma'. G Richards , D J Suqarbaker , R Bueno . *Cancer Research*
- 222 2002. 62 p. .