

Feature Subset Selection Approach Based on RFE-SVM For Cancer Classification

Khin May Win, Nang Saing Moon Kham
University of Computer Studies, Yangon
winn.km05@gmail.com, moonkhamucsy@gmail.com

Abstract

The data analysis can be many different points of views from many researchers. New method is evaluated for variable subset relevance with a view to variable selection. The new criteria are derived from support vector approach based on classification problems. This search can be efficiently performed by minimizing the generalization error. Selecting a small subset of features variables not only improves the efficiency of the classification algorithms but also improve the cancer classification accuracy. The process of building classifier is divided into two components (i) selection of variables features (i.e genes) (ii) selection of classification method. This study indicates that the classification problem is more difficult than the binary one for the gene expression data sets. This new method is related to structural risk minimization and thus leads to good generalization. The proposed method is compared to some standard feature selection method with real data sets. This method is computationally efficient with better classification performance.

Keywords —support vector machines, linear kernels, variable selection, feature ranking, over fitting.

1. Introduction

Genetic information of cells is stored in DNA and all cells in an organism which have different gene expression patterns. Expression of genes can be assessed with DNA microarray or serial analysis of gene expression among several other techniques.

The development of DNA microarray

technology has been produced large amount of gene data. This technology has been applied to the field of accurate prediction and diagnosis of cancer disease. Especially accurate classification of cancer is very important issue for treatment of cancer. To precisely classify cancer we have to select genes related to cancer. We attempt to choose the cancer related genes by using feature selection and combined the classifiers to improve the performance of classification. The gene expression data usually consist of huge number of genes and the necessity of tools analyzing them to get useful information. The identification of discriminated genes is of fundamental and practical interest for biomedical field. The usual way of transforming the input data structure into a vector representation, which is suitable for training a learning algorithm.

The selection of relevant variables may also be useful to gain some insight about the concept to be learned. Other advantages of feature selection include cost reduction of data gathering and storage (in medical applications) and computational speedup. In this work, we investigate the efficiency of criteria derived from support vector machines (SVMs) for variable selection in application to classification problems. This can be seen as an extension of the SVM-RFE algorithm. In genomics expression, the data set is usually plagued with large number of variables versus the small number of records or vectors (the problem is known as the ‘curse of dimensionality’). Genes are clustered first, and usual methods used are K-means clustering and hierarchical clustering [12], Singular Value Decomposition or Principal Component Analysis, supervised clustering and fuzzy clustering methods [4, 6]. In the dual space the decision function is expressed as a linear combination of

basis functions parameterized by the supporting patterns. The supporting patterns correspond to the class are chosen automatically by the maximum margin training procedures [10]. In case of polynomial classifiers, the Perceptron representation involves an untraceable number of parameters [5]. This problem is overcome in the dual space representation, where the classification rule is weighted sum of a kernel function for each support vectors patterns.

The incremental informative content of more variables is not always significant. Among existing methods, S2N performs good combination of gene selection methods and classifiers for microarray data. However, a significant improvement is achieved from choosing the appropriate parameter to small value. It supports the performance of the SVM classifier [15]. The selection of relevant variables may be insight to enhance the generalization performance of the learning process. High order polynomial classifiers with very large training sets can therefore be handled efficiently with the feature selection method. A new feature selection method criterion function was proposed based on the Feature weighting value. For hope, the learning process may avoid redundant, noisy or unreliable information of features.

2. Related Work

From the data mining viewpoint, gene selection problem is essentially a feature selection or dimensionality reduction problem. After reviewing, the soft margin SVM classifier can perform ranking criteria derived from SVM and an associated algorithm for feature selection. Finally, relationships with other SVM-based feature selection methods are given.

2.1 Support Vector Machines

A depiction of the basic problem is started for any machine learning algorithm, namely the detection of statistically stable patterns in training data with large amount of input data. One strategy is first to train linear Support Vector Machines (SVM) on a subset of training data to create initial classifiers. For classification problem, suppose a gene expression data set, data points are into two

classes $y_i \in \{-1, 1\}$, with binary class labels $X = \{x_i | i = 1, 2, \dots, n\}$, $Y = \{y_i | i = 1, 2, \dots, n\}$.

Given this training data set, we wish to predict the class label y for a new data point x .

Each vector $X = \{x_i\}_{i=1}^m$ labeled by $\{y_i\}_{i=1}^m$ in the gene expression matrix may be thought of as a point in an m -dimensional expression space. In theory, a simple way to build a binary classifier is to construct a hyperplane which can separate class members. The decision function becomes:

$$f(x) = (w, \Phi(x)) + b \quad (1)$$

Unfortunately, most real-world problems involve non separable data for which there does not exist a hyperplane that successfully separates the positive from the negative examples. One solution to the inseparability problem is to map the data into a higher-dimensional space and define a separating hyperplane there. This higher dimensional space is called the feature space, as opposed to the input space occupied by the training examples.

For linearly non-separable cases, one can introduce slack variables ξ and accordingly, the discriminated function is defined by

$$y_i(w \cdot x_i + b) \geq 1 - \xi \geq 0 \quad (2)$$

measure the deviation of a data point from optimal hyper plane. SVM are designed by minimizing

$$\Phi(w, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{subject to: } y_i(w \cdot x_i + b) \geq 1 - \xi_i, \xi_i \geq 0 \quad (3)$$

Minimize over: $(w, b, \xi_1, \dots, \xi_m)$:

$$\|w\|_p^p + C \sum_{i=1}^m \xi_i \quad (4)$$

subject to: $\forall_{i=1}^m: y_i (<w, \Phi(x_i)> + b) \geq 1 - \xi_i, \xi \geq 0$

Where $\langle \cdot \rangle$ is the inner product of mapping function between two vectors. C is a user - specified constant for controlling the penalty to the violation terms denoted by each (slack variables). $C+$ and $C-$ control the penalty to the violation of positive and negative examples of each features respectively. The w and b constitute of the classifier,

$$y = \text{sign}(\langle w, \Phi(x) \rangle + b) \quad (5)$$

Furthermore, artificially separating the data in this way exposes the learning system to the risk of finding trivial solutions that over fit the data. SVMs elegantly sidestep both difficulties [4]. They avoid over fitting by choosing the maximum margin separating hyperplane from among the many that can separate the positive from negative examples in the feature space.

2.2. Correlation Based Feature Ranking Algorithms for Gene Selection

Gene selection can be viewed as a feature selection or dimensionality reduction problem. Currently, there are mainly two kinds of algorithms for gene selection: correlation-based algorithms and backwards elimination algorithms. Correlation-based feature ranking algorithms work in a forward selection. Then, some top ranked genes are selected to form the most informative gene subset [12], [13].

Some commonly used ranking matrices are:

Signal-to-Noise (S2N):

$$w_i = |\mu_i(+)-\mu_i(-)| / \sigma_i(+)+\sigma_i(-) \quad (6)$$

Fisher Criterion (FC):

$$w_i = (\mu_i(+)-\mu_i(-))^2 / \sigma_i(+)^2 + \sigma_i(-)^2 \quad (7)$$

T-Statistics (TS):

$$w_i = |\mu_i(+)-\mu_i(-)| / \sqrt{\sigma_i(+)^2/n(+) + \sigma_i(-)^2/n(-)} \quad (8)$$

At above equations, $\mu_i(+)$ and $\mu_i(-)$ are the mean values of the i th gene expression data over positive and negative samples in the training data set, respectively. $\sigma_i(+)$ and $\sigma_i(-)$ are the corresponding standard deviations. $n(+)$ and $n(-)$ denote the numbers of positive negative training samples, respectively. A larger w_i means that the i th gene is more informative for cancer classification.

3. Recursive Feature Elimination

The RFE approach operates in an iterative manner to eliminate features weighted by weak weights specified in a 2-norm SVM model. Nested subsets of features are selected in a sequential backward elimination manner, which starts with all the features and remove one feature each time. In this way, in the end, all the feature variables are ranked. At each step, the coefficients of the weight vector w of a linear SVM are used as the feature ranking criterion. The recursive elimination procedure used as follows:

- (1) Start: ranked feature $R = []$;
Selected subset $S = [1, \dots, d]$;
- (2) Repeat until all features are ranked:
 - (a) Train a linear SVM with all the training data and variables in S ;
 - (b) Compute the weight vector
 - (c) Compute the ranking scores in S : $c_i = (w_i)^2$;
 - (d) Find the feature with the smallest ranking score: $e = \arg \min_i c_i$;
 - (e) Update R : $R = R[e, R]$;
 - (f) Update S : $S = S - [e]$;
- (3) Output: Ranked feature list R

The algorithm can be generalized to remove more than one feature per step.

In SVM-FRE, the following SVM formulation is used

$$\min_{w, b, \xi_i} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^1 \xi_i^2 \quad (9)$$

$$\text{subject to } y_i(w \cdot z_i + b) \geq 1 - \xi_i, \xi_i \geq 0 \quad (10)$$

The formulation of SVM is usually solve by the following dual problem with mapping function $\langle \phi(x_i), \phi(x_j) \rangle$, $z_i = \phi(x_i)$.i.e. kernel function.

Maximize over $(\alpha_1, \dots, \alpha_m)$:

$$J = -\frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle \phi(x_i), \phi(x_j) \rangle + \sum_{i=1}^m \alpha_i \quad (11)$$

$$\text{subject to } : \sum_{i=1}^m \alpha_i y_i = 0 \text{ and } \forall_{i=1}^m : 0 \leq \alpha_i \leq C$$

The weighting vector w is given by $\sum_{i=1}^m \alpha_i y_i \phi(x_i)$. Calculating w might be prohibitively expensive when nonlinear kernels are used. Using w_i^2 as ranking score corresponds to removing the feature whose removal change the objective function test. The approximation of the change in objective function caused by removing the i^{th} feature by expanding the objective function in Taylor series to second order.

$$\Delta J(i) = \frac{\partial J}{\partial w_i} \Delta w_i + \frac{\partial^2 J}{\partial w_i^2} (\Delta w_i)^2 \quad (12)$$

At the optimum of J , the first order term can be neglected and with $\Delta w_i = (1/2) \|\Delta w\|^2$, the equation becomes

$$\Delta J(i) = (\Delta w_i)^2 \quad (13)$$

For tasks of selecting features in the input space, it is often unnecessary to calculate the true vector of w , which corresponds to the features of mapped data. We might only need the set of weights relevant to the features of input vectors. In RFE, the weight of a feature is approximately measured by the change of the objective value J in SVM model by leaving this feature out.

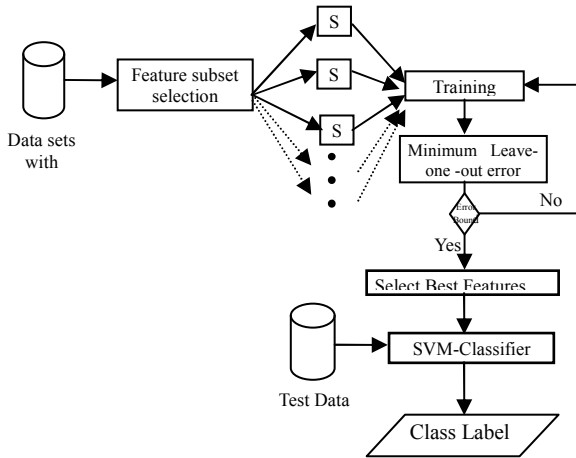


Figure 2.1 System Flow for the proposed method

4. Experimental Results

We randomly split the original dataset into a training set and a test set and keep percentages of the positive and negative samples same in the training and test sets. We summarize some basic information about the datasets, including the number of features, the sizes of the training and test sets. However, the total numbers of available samples in our mass spectrometry datasets are small. In such a case, the test error may be biased due to an “unfortunate” partition of training and test sets. Thus, instead of reporting such a test error from one division of training and test sets, merging the training set and test set and then partition the total samples again into a training set and a test set randomly by stratified sampling for 100 times; for each division, first train a linear SVM classifier on the training set (hyperparameter C is to be selected by 5-fold cross-validation on the training set) and then test it on the corresponding test set; from this 100 trials we can compute the averages of performance measures.

4.1 Leukemia Cancer Dataset

Leukemia dataset consists of 72 samples: 25 samples of acute myeloid leukemia (AML) and 47 samples of acute lymphoblastic leukemia (ALL). The source of the gene expression measurements was taken from 63 bone marrow samples and 9 peripheral blood samples. Gene expression levels in these 72 samples were measured using high density oligonucleotide microarrays of 72 samples were used as training data and the remaining were used as test data in this paper. Each sample contains expression level of 7129 genes.

4.2 Colon Cancer Dataset

Colon dataset consists of 62 samples of colon epithelial cells taken from colon-cancer patients. Each sample contains over 2000 gene expression levels. 40 of 62 samples are colon cancer samples and the remaining are normal samples. Each sample was taken from tumors and normal healthy parts of the colons of the same patients and measured using high density oligonucleotide array out of 62 samples were training data and the remaining were test data .

4.3 Lymphoma Cancer Dataset

Lymphoma data sets cell diffuse large cell lymphoma (B-DLCL) is a heterogeneous group of tumors, 4026 genes containing. Gene expression profiling has revealed two distinct tumor subtypes of B-DLCL: germinal center B cell-like DLCL and activated B cell-like DLCL. These data sets contains 77 tissue samples, 58 are diffuse large B-cell lymphomas (DLBLC) and remaining 19 samples are follicular lymphomas (FL).

4.4 Preprocessing

In real-world data, the representation of data often uses too many features, but only a few of them may be related to the target concept. There may be redundancy, where certain features are correlated so that is not necessary to include all of them in modeling; and interdependence, where two or more features between them convey important information that is obscure if any of them is included on its own. For the Colon and Leukemia, each sample was standardized to zero mean and unit variance across genes.

4.4.1 Data Normalization

Normalization is a scaling down transformation of the features. Within a feature there is often a large difference between the maximum and minimum values, e.g. 0.01 and 1000. When normalization is performed the valued magnitudes and scaled to appreciably low values.

z- score normalization becomes;

$$v' = \frac{v - \text{mean}_A}{\text{stand_dev}_A} \quad (15)$$

where v is the old feature weight value and v' the new weight value one, mean value for each feature and standard deviation is one. The feature selection is concerning with these facts.

4.4.2 Features' Characteristics

Feature subset selection is the identifying and removing as much relevant features as possible. This reduces the dimensionality of the data and enables learning algorithm to operate faster and more effectively. Generally features are characterized as:

- **Relevant** : These are features have an influence on the output and their role cannot be assumed by the rest.
- **Redundant**: A redundancy exists whenever a feature can take the role of another.
- **Information**: Feature X is preferred to feature Y if the information gain from feature X is greater than from feature Y.
- **Dependence**: The coefficient is a classical dependence measure and can be used to find the correlation of feature X with class (1) is higher than the correlation of feature Y with class (1), then feature X is preferred to Y.

Table 1. Accuracy Comparison with other algorithms on Leukemia Dataset

Models	8 genes	Best (<=8)	Mean (<=20)	Std (<=20)
S2N correlation	0.8264	0.8356	0.8451	0.0254
FC correlation	0.8126	0.8237	0.8634	0.0213
Default SVM-RFE	0.8041	0.9012	0.0509	0.0489
Extended SVM-RFE	0.9816	0.9006	0.0600	0.0601

Table 2. Accuracy Comparison with other algorithms on Colon Dataset

Models	8 genes	Best (<=8)	Mean (<=20)	Std (<=20)
S2N correlation	0.8648	0.8649	0.8611	0.0254
FC correlation	0.8226	0.8664	0.8611	0.0223
Default SVM-RFE	0.88871	0.9034	0.0559	0.0599
Extended SVM-RFE	0.9616	0.9677	0.0622	0.0632

Table 3 . Accuracy Comparison with other algorithms on Lymphoma Dataset

Models	8 genes	Best (<=8)	Mean (<=20)	Std (<=20)
S2N correlation	0.8268	0.8568	0.8604	0.0125
FC correlation	0.8266	0.8568	0.8604	0.0125
Default SVM-RFE	0.9056	0.9355	0.0549	0.0759
Extended SVM-RFE	0.9622	0.9624	0.0905	0.0600

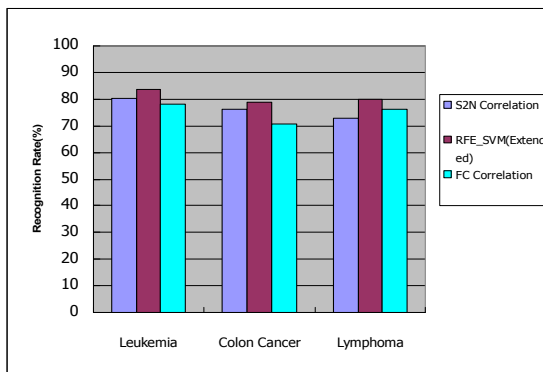


Figure.2. Performance Comparison with some correlation based methods

5. Conclusion

In this study, a new feature subset selection algorithm for classification task using SVMs was developed. The proposed method was assumed that very few features are needed to classify the given samples and smallest subset may provide more insight into the data. Looking at the performance of SVMs without SVM and with SVM in tables. The classification performance of extended SVM-RFE is much better than of other SVMs with all feature subset as input variables. In terms of dimensionality reduction, the best accuracy we get by starting from the original set and reduces the irrelevant features in each individual gene subset. The high prediction accuracy also strengths the promising application prospects of mass spectrometry patterns in the

further cancer classification.

Acknowledgements

This research was supported by SDRC (Software Development and Research Center) in University of Computer Studies, Yangon. SDRC is a Research Center designated by Myanmar Science and Engineering Foundation and Ministry of Science & Technology.

References

- [1] C.Emmanouilidis, A.Hunter, and J.MacIntyre,"A multiobjective evolutionary setting for feature selection and a commonality-based crossover operator", in Proceedings of the 2000 Congress on Evolutionary Computation (CEC00).
- [2] Cancer Program Data Set Chang, Chih-Chung and Lin,Chih-Jen." LIBSVM: a library for support Vector Machines", <http://www.csie.ntu.edu.tw/~cjlin/>, 2003 <http://www.broad.mit.edu/cgi/datasets.cgi>
- [3] Dudoit S,Fridlyand J,Speed TP: Comparison of discrimination method for the classification of tumors using gene expression data.J Amer Statist Assoc 2002,97:77-87.
- [4] Eisen, M.B and Brown ,P.O.(1999); DNA arrays for analysis of gene expression. Methods Enzymol, 303: 179-205
- [5] Furey, T.S.,Cristianini, N.,Duffy,N.,Bednarski ,D.W.,Schummer,M.and Haussler,D.(2000): Support Vector Machines Classification and validation of cancer tissue samples using microarray expression data.Bioinformatics, 16(10):906-914.
- [6] GolubTR, Sloomim DK,Tamayo P,Huard C,Gassenbeek M,Mesirov JP,Coller Fgff,Loh ML,Downing JR,Caligiuri MA,Bloomfield CD,Lander ES:Molecular classification of cancer:class discovery and class prediction by gene expression monitoring.Science 1999,286:531-537.
- [7] Harrington,C.A.,Rosenow,C., (2000); Monitoring gene expression using DNA microarrays. Opin.Microbiol.,3: 285-291.
- [8] I.Guyon, J.Weston, S.Barnhill,and V.Vapnik. Gene Selection for cancer classification using support vector machines, Machine

Learning,2000.

- [9] K.M Win and Kham N.S.M, “Minimizing Essential Set Based Feature selection for Cancer Classification”, ICCA2008, Yangon, Myanmar, Feb 14-15, 2008
- [10] Li,L., Weinberg, C.R., Darden,T.A. and Pedersen , L.G.(2001): Gene Selection for sample classification base don gene expression data: study of sensitivity to choice of parameters of GA/KNN method.Bioinformatics,17(12):1131-1142.
- [11] M.Momma and K.P.Bannett.A Pattern search method for model selection of support vector regression.In Proceeding of the SIAM International Conference on Data mining. Philadelphia,Pennsylvania,2002.SIAM
- [12] N. Cristaianini and J. Shawe-Taylor. *Introduction to Support Vector Machines*. Cambridge University Press,2000
- [13] T.Furey, N.Duffy, D.Bednarski, M. Schummer, and D.Haussler” Support Vector Machine Classification and Validation of Cancer Tissue Samples Using Microarray Expression Data,” *Bioinformatics*, vol.16, pp.906-914,2000
- [14] P.Pavlidis at al.,”Gene Functional Analysis from Heterogeneous Data,” Proc. Conf. *Research in Computational Molecular Biology (RECOMB)*, pp.249-255,2001.
- [15] Vapnik and O.Chapelle. Bounds on error expectation for support vector machines. *Neural Computation*, 12(9)2000.
- [16] Wessels Weston,J. and C.Watkins.”Support Vector Machines for Multi-Class Pattern Recognition”, Proceedings of he Seventh European Symposium On Artificial Neural Networks,1999.