



A simple and efficient algorithm for gene selection using sparse logistic regression

S. K. Shevade¹ and S. S. Keerthi^{2,*}

¹Department of Computer Science and Automation, Indian Institute of Science, Bangalore 560012, India and ²Control Division, Department of Mechanical Engineering, National University of Singapore, Singapore 117576, Republic of Singapore

Received on November 30, 2002; revised on March 12, 2003; accepted on May 29, 2003

ABSTRACT

Motivation: This paper gives a new and efficient algorithm for the sparse logistic regression problem. The proposed algorithm is based on the Gauss–Seidel method and is asymptotically convergent. It is simple and extremely easy to implement; it neither uses any sophisticated mathematical programming software nor needs any matrix operations. It can be applied to a variety of real-world problems like identifying marker genes and building a classifier in the context of cancer diagnosis using microarray data.

Results: The gene selection method suggested in this paper is demonstrated on two real-world data sets and the results were found to be consistent with the literature.

Availability: The implementation of this algorithm is available at the site <http://guppy.mpe.nus.edu.sg/~mpessk/SparseLOGREG.shtml>

Contact: mpessk@nus.edu.sg

Supplementary Information: Supplementary material is available at the site <http://guppy.mpe.nus.edu.sg/~mpessk/SparseLOGREG.shtml>

INTRODUCTION

Logistic regression is a powerful discriminative method. It also has a direct probabilistic interpretation built into its model. One of the advantages of logistic regression is that it provides the user with explicit probabilities of classification apart from the class label information. Moreover, it can be easily extended to the multi-category classification problem. In this paper the problem of logistic regression is addressed with the application of gene expression data particularly in mind.

The development of microarray technology has created a wealth of gene expression data. Typically, these data sets involve thousands of genes (features) while the number of tissue samples is in the range 10–100. The objective is to design a classifier which separates the tissue samples into pre-defined classes (e.g. tumour and normal). Different machine

learning techniques like Parzen windows, Fisher's linear discriminant, and decision trees have been applied to solve this classification problem (Brown *et al.*, 2000). Since the data dimension is very large, Support Vector Machines (SVMs) have been found to be very useful for this classification problem (Brown *et al.*, 2000; Furey *et al.*, 2000). Apart from the classification task, it is also important to remove the irrelevant genes from the data so as to simplify the inference. One of the aims of the microarray data experiments is to identify a small subset of informative genes, called marker genes, which discriminate between the tumour and the normal tissues, or between different kinds of tumour tissues. Identification of marker genes from a small set of samples becomes a difficult task and there is a need to develop an efficient classification algorithm for such an application where the number of features is much larger than the number of samples. Furey *et al.* (2000) used Fisher score to perform feature selection prior to training. This method was compared against feature selection using radius-margin bound for SVMs by Weston *et al.* (2001). Guyon *et al.* (2002) introduced an algorithm, called Recursive Feature Elimination (RFE), where features are successively eliminated during training of a sequence of SVM classifiers. Li *et al.* (2002) introduced two Bayesian classification algorithms which also incorporate automatic feature selection. In their work, the Bayesian technique of *automatic relevance determination* (ARD) was used to perform feature selection.

The main aim of this paper is to develop an efficient algorithm for the sparse logistic regression problem. The proposed algorithm is very much in the spirit of the Gauss–Seidel method (Bertsekas and Tsitsiklis, 1989) for solving unconstrained optimization problems. The algorithm is simple, *extremely easy to implement*, and can be used for feature selection as well as classifier design, especially for the microarray data set. It *does not need any sophisticated LP or QP solver*, nor does it involve any matrix operations which is a drawback of some of the previously suggested methods. Our algorithm also has the advantage of not using any extra matrix storage and therefore it can be easily used in cases where the number of features used is very high.

*To whom correspondence should be addressed.

The remaining sections of the paper contain the following: description of the problem formulation and the optimality conditions of the problem; details of the proposed algorithm; the feature selection and the classifier design procedure; computational experiments and, concluding remarks.

PROBLEM FORMULATION AND OPTIMALITY

In this paper we focus on the two category classification problem. The ideas discussed here can be easily extended to the multi-category problem and those details will be addressed in a future paper. We will now describe the problem formulation using microarray tissue classification as an example. Given a microarray data set containing m tissue samples, with each tissue sample represented by the expression levels of N genes, the goal is to design a classifier which separates the tissue samples into two predefined classes. Let $\{(\tilde{x}_i, y_i)\}_{i=1}^m$ denote the training set, where \tilde{x}_i is the i th input pattern, $\tilde{x}_i \in \mathbb{R}^N$ and y_i is the corresponding target value; $y_i = 1$ means \tilde{x}_i is in class 1 and $y_i = -1$ means \tilde{x}_i is in class 2. Note that the vector \tilde{x}_i contains the expression values of N genes for the i th tissue sample and \tilde{x}_{ij} denotes the expression value of gene j for the i th tissue sample. To conveniently take into account the bias term in the regression function, let us define $x^T = (1, \tilde{x}^T)$. Let us specify the functional form of the regression function $f(x)$ as a linear model,

$$f(x_i) = \sum_{j=0}^N \alpha_j x_{ij} \quad (2.1)$$

where α denotes the weight vector. We address the first element (corresponding to the bias term) of the vectors x and α as the ‘zeroth element’. To avoid notational clutter, throughout this paper we use the index i to denote the elements of the training set ($i = 1, \dots, m$) and the index j to represent the elements of x and α ($j = 0, \dots, N$, unless specified explicitly).

Typically for the microarray data, $N \gg m$. That is, the training samples lie in a very high-dimensional space. Therefore, linear separability of two classes may not be a problem. Hence, it is appropriate to consider using a linear classifier in the input space directly without transforming it into a higher dimensional feature space. Furthermore, for such data sets it is important to note that finding a linear classifier in the input space involves estimation of a large number of parameters [in Equation (2.1)] using a very small number of training examples. It is therefore possible to derive different linear classifiers for the same problem as the problem is underdetermined. How the proposed algorithm takes care of the effect of multiple solutions is discussed in the section on feature selection and classifier design.

For diagnostic purposes, it is important to have a classifier which uses as few features (genes) as possible. Since we are

interested in finding a small subset of genes which enables us to do good classification it is appropriate to use a sparse model of the regression function, that is, in the final function $f(x)$ in (2.1), only a small number of α s corresponding to relevant features will have non-zero values. Therefore, we solve the sparse logistic regression problem, which can be formulated as the following optimization problem:

$$\begin{aligned} \min_{\alpha} \quad & \rho = \sum_i g(-y_i f(x_i)) \\ \text{s.t.} \quad & \sum_{j=1}^N |\alpha_j| \leq t \end{aligned} \quad (2.2)$$

where $t \geq 0$ is a parameter that is tuned using techniques such as cross-validation and the function g is given by:

$$g(\xi) = \log(1 + e^{\xi}) \quad (2.3)$$

which is the negative log-likelihood function associated with the probabilistic model

$$\text{Prob}(y|x) = \frac{1}{1 + e^{-y \cdot f(x)}} \quad (2.4)$$

Once the α s are determined by solving (2.2), the class of the test sample, \bar{x} , is $+1$ if $f(\bar{x}) > 0$ and -1 otherwise.

The formulation (2.2) was first suggested by Tibshirani (1996) as an extension of the ‘LASSO’ (Least Absolute Shrinkage and Selection Operator) method for the linear regression problem. The LASSO estimator for the regression problem is obtained by solving the following optimization problem:

$$\min_{\alpha} \quad \frac{1}{2} \sum_i (y_i - f(x_i))^2 \quad \text{s.t.} \quad \sum_{j=1}^N |\alpha_j| \leq t \quad (2.5)$$

In Tibshirani (1996), two iterative algorithms were suggested for the solution of (2.5). One of these algorithms treats (2.5) as a problem in $N + 1$ variables and 2^N constraints. In each iteration it solves a linear least squares problem subject to a subset of constraints E . Once a solution is found in an iteration, a constraint violated by the current solution is added to the set E and a new feasible point is found. The linear least squares problem is then resolved subject to the new set of constraints E . This procedure is repeated until the optimality conditions are satisfied. Note that, during each iteration one constraint is added to the set E and there exist only a finite (2^N) number of constraints. Therefore, this algorithm will converge in a finite number of steps. Use of active set method was also suggested for the solution of (2.5) where every time the set E contains only the set of equality constraints satisfied by the current point. The second algorithm, suggested in Tibshirani (1996), states (2.5) as a problem in $2N + 1$ variables and $2N + 1$ constraints, which can be solved using any quadratic programming solver.

In Osborne *et al.* (2000), the problem (2.5) was treated as a convex programming problem and an efficient algorithm for computing the linear model coefficients was given using duality theory. This algorithm was stated only for the cases where the objective function is a quadratic loss function. These ideas, however, can be extended to problems with non-quadratic objective functions, for example (2.2), by optimizing the quadratic approximation of the cost function in (2.2) at any given point and repeating this procedure until convergence. Recently, Roth (2002) extended the ideas given in Osborne *et al.* (2000) to a more general class of cost functions including the one given in (2.2). The proof of convergence of this algorithm to a general class of loss functions is also given therein. In this method, the problem (2.2) is transformed to a L_1 constrained least squares problem which can then be solved using iteratively re-weighted least squares (IRLS) method on a transformed set of variables. To avoid solving large dimensional constrained least squares problems, especially for the microarray data, Roth (2002) uses a wrapper approach where a maximum violating variable is added to a small variable set, X , and the constrained least squares problem is solved with respect to the variables in the set X only. This procedure is repeated till the optimality conditions are satisfied.

All of the above methods work well and are very efficient; but they rely on mathematical programming solvers that require detailed matrix operations. The main contribution of this paper is to utilize the special structure of (2.2) and devise a simple algorithm which is extremely easy to implement; it neither uses any mathematical programming package nor needs any matrix operations. The simplicity of our algorithm is evident from the pseudocode. See the supplementary information for details.

Using optimality conditions it can be shown that there exists a $\gamma > 0$ for which, (2.2) is equivalent to the following unconstrained optimization problem:

$$\min_{\alpha} W = \gamma \sum_{j=1}^N |\alpha_j| + \sum_i g(-y_i f(x_i)) \quad (2.6)$$

This means that the family of classifiers obtained by varying t in (2.2) and the family obtained by varying γ in (2.6) are the same. The addition of the penalty term, $\sum_j |\alpha_j|$, to the original objective function can be seen as putting a Laplacian prior over the vector α . The above formulation thus promotes the choice of a sparse model, where the final α s are either large or zero. In this paper, we shall concentrate on the model in (2.6).

We solve (2.6) directly without converting it into its dual. The problem (2.6) can be solved using an appropriate unconstrained nonlinear programming technique. We choose to basically use the Gauss-Seidel method which uses coordinate-wise descent approach, mainly because the method is extremely easy to implement while also being very efficient. In this method, one variable is optimized at a time

keeping the other variables fixed, and the process is repeated till the optimality conditions are satisfied. Asymptotic convergence of this method for a more general version of the problem has been proved in (Bertsekas and Tsitsiklis, 1989; Chapter 3, Proposition 4.1). It should be noted that strict convexity of g plays an important role in this proof.

We now derive the first order optimality conditions for (2.6). Since (2.6) is a convex programming problem these conditions are both necessary and sufficient for optimality. Let us define

$$\begin{aligned} \xi_i &= -y_i f(x_i) \\ F_j &= \sum_i \frac{e^{\xi_i}}{1 + e^{\xi_i}} y_i x_{ij} \end{aligned} \quad (2.7)$$

The first order optimality conditions for the problem (2.6) can be easily derived from geometry: (i) since W is differentiable with respect to α_0 , $\partial W / \partial \alpha_0 = 0$; (ii) if $j > 0$ and $\alpha_j \neq 0$, then, since W is differentiable with respect to α_j at such a α_j , we have $\partial W / \partial \alpha_j = 0$ and (iii) if $j > 0$ and $\alpha_j = 0$, then, since W is only directionally differentiable with respect to α_j at $\alpha_j = 0$, we require the right side derivative of W with respect to α_j to be non-negative and the left side derivative to be non-positive. These conditions can be rewritten as the following algebraic conditions.

$$\begin{aligned} F_j &= 0 & \text{if } j = 0 \\ F_j &= \gamma & \text{if } \alpha_j > 0, \quad j > 0 \\ F_j &= -\gamma & \text{if } \alpha_j < 0, \quad j > 0 \\ -\gamma &\leq F_j \leq \gamma & \text{if } \alpha_j = 0, \quad j > 0 \end{aligned}$$

Thus, if we define

$$\begin{aligned} \text{viol}_j &= |F_j| & \text{if } j = 0 \\ &= |\gamma - F_j| & \text{if } \alpha_j > 0, \quad j > 0 \\ &= |\gamma + F_j| & \text{if } \alpha_j < 0, \quad j > 0 \\ &= \psi_j & \text{if } \alpha_j = 0, \quad j > 0 \end{aligned} \quad (2.8)$$

where $\psi_j = \max(F_j - \gamma, -\gamma - F_j, 0)$, then the first order optimality conditions can be compactly written as

$$\text{viol}_j = 0 \quad \forall j \quad (2.9)$$

Since, in asymptotically convergent procedures it is hard to achieve exact optimality in finite time, it is usual to stop when the optimality conditions are satisfied up to some tolerance, τ . For our purpose, we will take it that any algorithm used to solve (2.6) will terminate successfully when

$$\text{viol}_j \leq \tau \quad \forall j. \quad (2.10)$$

This can be used as a stopping condition for (2.6). Throughout, we will refer to optimality as optimality with tolerance.

In the next section, we describe the actual details of our algorithm for solving the problem (2.6).

A SIMPLE ALGORITHM FOR SPARSE LOGISTIC REGRESSION

The algorithmic ideas discussed here are based on the Gauss–Seidel method used for solving unconstrained optimization problem. In solving (2.6) using Gauss–Seidel method, one variable α_j which violates the optimality conditions is chosen and the optimization subproblem is solved with respect to this variable α_j alone, keeping the other α s fixed. This procedure is repeated as long as there exists a variable which violates the optimality conditions. The method terminates when the optimality conditions (2.9) are satisfied. Note that the objective function in (2.6) is strictly convex and it will strictly decrease at every step when the optimization subproblem is solved with respect to one variable. Repeated application of this procedure for the problem (2.6) will make sure that the algorithm will converge asymptotically (Bertsekas and Tsitsiklis, 1989) to the solution. Of course, if the stopping condition (2.10) is used then the algorithm will terminate in finite time.

For a given α , let us define the following sets: $I_z = \{j : \alpha_j = 0, j > 0\}$; and $I_{nz} = \{0\} \cup \{j : \alpha_j \neq 0, j > 0\}$. Also, let $I = I_z \cup I_{nz}$. The key to efficiently solving (2.6) using Gauss–Seidel method is the selection of the variable α_j in each iteration with respect to which the objective function is optimized. At optimality, it is expected that the resulting model is sparse, that is, there will be only a few weights α_j with non-zero values. It is important that the algorithm spends most of its time adjusting the non-zero α s and making the subset I_{nz} self-consistent as far as optimality of the variables in this subset is concerned. Therefore, it is appropriate to find the maximum violating variable, say α_v in the set I_z and then solve (2.6) using the variables in the set $I_{nz} \cup \{v\}$ only, until the optimality conditions for these variables are satisfied. This procedure can be repeated till no violator remains in the set I_z , at which point the algorithm can be terminated. This procedure can be best explained using Algorithm 1.

Algorithm 1

Input Training Examples

Initialize α s to 0.

while Optimality Violator exists in I_z

Find the maximum violator, v , in I_z

repeat

Optimize W w.r.t. α_v

Find the maximum violator, v , in I_{nz}

until No violator exists in I_{nz}

end while

Output A set of α s for the function in (2.1)

In this algorithm, all α s are set to zero initially which implies that only the set I_z exists. The algorithm can be thought of as a two-loop approach. The type I loop runs over the variables in the set I_z to choose the maximum violator, v . In the type II loop, W is optimized with respect to α_v , thus modifying the set I_{nz} and the maximum violator in the set I_{nz} is then found. This procedure in the type II loop is then repeated until no violators are found in the set I_{nz} . The algorithm thus alternates between the type I and type II loops until no violators exist in either of the sets, I_z and I_{nz} .

Once the maximum violating variable in a given set is chosen, some non-linear optimization technique needs to be used to solve the unconstrained optimization problem (2.6) with respect to a single variable. Note that the objective function is convex. A combination of the bisection method and Newton method was used in our algorithm. In this method, two points L and H for which the derivative of the objective function has opposite signs are chosen. This ensures that the root always lies in a bracketed interval, $[L, H]$. The Newton method is tried first, and a check is made in the algorithm to make sure the iteration obtains a solution in this interval of interest. If the Newton method takes a step outside the interval we do not accept this next point, and instead resort to the bisection method. In the bisection method, the next point is chosen to be the midpoint of the given interval. Depending upon the sign of the derivative of the objective function at this point the new interval is decided. Now, the abovementioned procedure starting from the trying of the Newton step is repeated. Since it is always ensured that the root lies in the interval of interest, the method is guaranteed to reach the solution.

It is important to note that the objective function, W , has different right-hand and left-hand derivatives with respect to α_j at $\alpha_j = 0$. Therefore, if the non-zero variable attains a value of 0 when our algorithm is executed it is necessary to see whether further progress in the objective function can be made by altering the same variable (but now in the opposite direction). This will make sure that after successful termination of the type II loop the optimality conditions for the variables in I_{nz} are satisfied.

Note from (2.7) and (2.8) that for checking the optimality conditions for each variable it is necessary to calculate ξ_i for all the examples. Since this calculation is done repeatedly, the efficient implementation of the algorithm requires that ξ is stored in the memory for all the examples. After a single variable, say α_j , is updated it becomes necessary to update $\xi_i \forall i$. This can be done by using the following simple equation:

$$\xi_i^{\text{new}} = \xi_i^{\text{old}} + (\alpha_j^{\text{old}} - \alpha_j^{\text{new}}) y_i x_{ij} \quad (3.1)$$

Let us now comment on the speed of our algorithm for (2.6). γ in (2.6) is a hyperparameter and one cannot identify any nominal value for it as a starting point. For a given value of γ and a data set containing 72 samples with 7129 features per sample, it took about 20s to solve (2.6) on a SUN UltraSparc III CPU running on 750 MHz machine. Once the

solution at a certain value of γ is found, the solution at another close-by value of γ can be efficiently found using the previous solution as the starting point. This idea is very useful for improving efficiency when γ needs to be tuned using cross-validation; see the section below.

FEATURE SELECTION AND CLASSIFIER DESIGN

Before we discuss the procedure for feature selection it is important to note that the solution of (2.6) is useful for the tasks of feature selection as well as classifier design. Since (2.6) automatically introduces sparseness into the model, the non-zero α s in the final solution of (2.6) help us in deciding the relevant features for classification. To identify the relevant features, we proceed as follows.

To design a classifier and identify the relevant features using (2.6), it is important to find the appropriate value of γ in (2.6). This can be done using techniques like k -fold cross-validation. For this purpose, the entire set of training samples is divided into k segments. For a given value of γ , the problem (2.6) is solved k times, each time using a version of the training set in which one of the segments is omitted. Each of these k classifiers is then tested on the samples from the segment which was omitted during training, and the ‘validation error’ is found by averaging the test set results over all the k classifiers. This procedure is repeated for different values of γ and the optimal value γ^* is chosen to be the one which gives minimum validation error. The final classifier is then obtained by solving (2.6) using all the training samples and setting γ to γ^* .

For a given value of γ , each of the k classifiers uses a set of genes (corresponding to non-zero α s) for classification. We assign a count of one to each gene if it is selected by a classifier at γ^* . Since there are k different classifiers designed at γ^* , a gene can get a maximum count of k meaning that it is used in all the k classifiers. To avoid any bias arising from a particular k -fold split of the training data,[†] the above procedure is repeated 100 times, each time splitting the data differently into k folds and doing the cross-validation. The counts assigned to a gene for 100 different experiments are then added to give a *relevance count* for every gene. Thus a gene can have a maximum relevance count of $100k$ since k -fold cross-validation was repeated 100 times.

One of the aims of using the microarray data is to develop a classifier using as few genes as possible. It is therefore important to rank the genes using some method and then select some of these genes in a systematic way so as to design a classifier which generalizes well on the unseen data. As the relevance count of a gene reflects its importance in some sense, it is appropriate to rank the genes based on the relevance count. Having decided the order of importance of the genes, the next

step is to choose the ‘right’ set of genes to design a classifier. For this purpose, the following systematic method can be used.

Let S denote the set of features being analysed. To examine the performance of these features, the k -fold cross-validation experiment [on (2.6)] is repeated 100 times using only the features in the set S . To avoid any bias towards a particular value of γ , the validation error at every value of γ is averaged over those 100 experiments. The minimum of these validation error values over different choices of γ then gives the average validation error for the set S . Ideally, one would like to identify the set of few features for which the average validation error is minimum. For this purpose, the top ranked genes (ranked according to relevance count) are added to the set S one-by-one, and the average validation error is calculated. The procedure is repeated as long as adding the next top ranked gene to the set S does not reduce the average validation error significantly. This final set S of genes is then used along with all the training examples to design a classifier for diagnostic purposes.

Note that this final set of genes is derived from all the training samples available. Typically for the microarray data set, the training set size is very small and there are no test samples available. Therefore, one needs to use the entire training data for classifier design. The average validation error mentioned earlier uses the entire training set for gene selection and so it is not an ‘unbiased estimate’. To examine how good the feature selection procedure is, it becomes necessary to give an unbiased estimate of this procedure. This estimate can be obtained using techniques of external cross-validation suggested in (Ambroise and McLachlan, 2002). For this purpose, the training data is split into P different segments ($P = 10$ is a good choice). The gene selection procedure described above is executed P times, each time using the training samples derived by omitting one of the P segments and testing the final classifier (obtained with the reduced set of genes) on the samples from the segment which was omitted during training. This procedure gives an unbiased estimate of the performance of the gene selection method described earlier. We will call this estimate as the average test error. As one might expect, the average test error is usually higher than the average validation error.

NUMERICAL EXPERIMENTS

In this section we report the test performance of the proposed gene selection algorithm on some publicly available data sets: colon cancer and breast cancer. As the problem formulation (2.6) and the one used in (Roth, 2002) are the same,[‡] the results of different algorithms for the same formulation are expected to be close.

[†] This issue is particularly important for microarray data since the number of examples is usually small.

[‡] The algorithm for the solution of (2.6) and the actual method of selecting relevant genes used by Roth (2002), however, are quite different from our method described in the previous section.

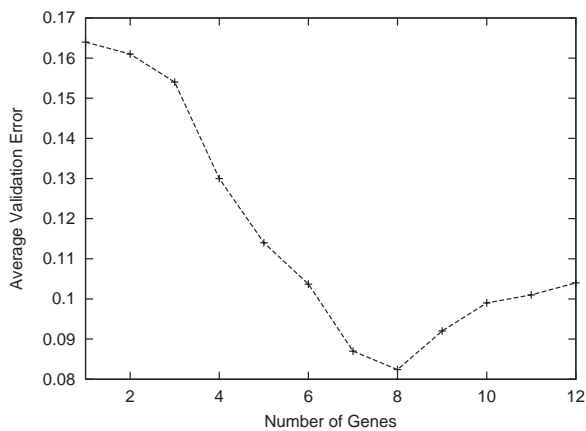


Fig. 1. Variation of validation error on the colon cancer data set, shown as the number of genes used in the feature set. The average validation error reached the minimum when eight genes, listed in Table 1, were used.

For the colon cancer data set (Alon *et al.*, 1999) the task is to distinguish tumour from normal tissue using microarray data with 2000 features per sample. The original data consisted of 22 normal and 40 cancer tissues. This dataset is available at <http://lara.enm.bris.ac.uk/colin>. We also evaluated our algorithm on the breast cancer data set (West *et al.*, 2001). The original data set consisted of 49 tumour samples. The number of genes representing each sample is 7129. The aim is to classify these tumour samples into estrogen receptor-positive (ER+) and estrogen receptor-negative (ER−). West *et al.* (2001) observed that in five tumour sample cases, the classification results using immunohistochemistry and protein immunoblotting assay conflicted. Therefore, we decided to exclude these five samples from our training set. This data set is available at the following site: http://mgm.duke.edu/genome/dna_micro/work/.

Each of these data sets was pre-processed using the following procedure. The microarray data set was arranged as matrix with m rows and N columns. Each row of this matrix was standardized to have mean zero and unit variance. Finally, each column of the consequent matrix was standardized to have mean zero and unit variance. This transformed data was then used for all the experiments.

The results of the gene selection algorithm on these different data sets are reported in Figures 1 and 2 and Tables 1 and 2. As discussed in the previous section, the relevance count of every gene was obtained by repeating the k -fold cross-validation procedure 100 times, where k was set to 3. Tables 1 and 2 denote the important genes used in the classifier design for each of the data sets along with their relevance counts.

For the colon cancer data set, the top three ranked genes are same as those obtained in (Li *et al.*, 2002), although the gene selection procedure adopted there was different than the one we used. The main reason for this is that the classification

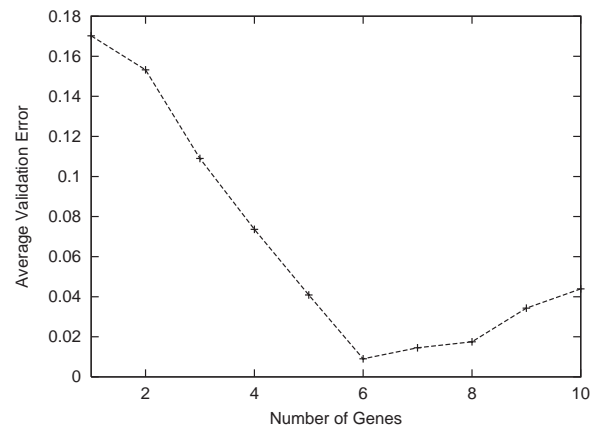


Fig. 2. Variation of average validation error on the breast cancer data set, shown as the number of genes used in the feature set. The average validation error reached the minimum when six genes were used. These genes are listed in Table 2.

Table 1. Selected top eight genes with their relevance counts for the colon cancer data set

Sr. no.	Gene annotation	Relevance count
1	<i>Homo sapiens</i> mRNA for GCAP-II/uroguanylin precursor	229
2	COLLAGEN ALPHA 2(XI) CHAIN (<i>H.sapiens</i>)	165
3	MYOSIN HEAVY CHAIN, NONMUSCLE (<i>Gallus gallus</i>)	151
4	PLACENTAL FOLATE TRANSPORTER (<i>H.sapiens</i>)	145
5	ATP SYNTHASE COUPLING FACTOR 6, MITOCHONDRIAL PRECURSOR (HUMAN)	124
6	GELSOLIN PRECURSOR, PLASMA (HUMAN)	117
7	Human mRNA for ORF, complete cds	107
8	60S RIBOSOMAL PROTEIN L24 (<i>Arabidopsis thaliana</i>)	100

Table 2. Selected top six genes with their relevance counts for the breast cancer data set

Sr. no.	Gene annotation	Relevance count
1	Human splicing factor SRp40-3 mRNA, complete cds	171
2	<i>H.sapiens</i> mRNA for cathepsin C	151
3	Human mRNA for KIAA0182 gene, partial cds	133
4	Y box binding protein-1 (YB-1) mRNA	130
5	Human transforming growth factor-beta 3 (TGF-beta3) mRNA, complete cds	127
6	Human microtubule-associated protein tau mRNA, complete cds	115

hypothesis need not be unique as the samples in gene expression data typically lie in a high-dimensional space. Moreover, if the two genes are co-regulated then it is possible that one of them might get selected during the gene selection process.

Table 3. Average test error for different data sets

Data set	Average test error
Colon cancer	0.177
Breast cancer	0.181

Some comments about the genes selected for the breast cancer data set (shown in Table 2) are worthy of mention. Cathepsins have been shown to be good markers in cancer (Kos and Schweiger, 2002). Recently, it has been shown that YB-1, also called Nuclease-sensitive element-binding protein 1 (NSEP1), is elevated in breast cancer; see the following site: <http://www.wfubmc.edu/pathology/faculty/berquin.htm>. Moreover, microtubule associated protein is the cell cycle control gene, and therefore is associated with cancer. TGF- β has been shown to be a critical cytokine in the development and progression of many epithelial malignancies, including breast cancer (Arrick and Derynck, 1996).

The generalization behaviour of our algorithm was examined using the technique of external cross-validation described at the end of the previous section. The average test error for each of the data sets is reported in Table 3. We also tested our algorithm on some other publicly available data sets. See the supplementary information for more details.

CONCLUSION

We have presented a simple and efficient training algorithm for feature selection and classification using sparse logistic regression. The algorithm is robust in the sense that on the data sets we have tried there was not even a single case of failure. It is useful, especially for problems where the number of features is much higher than the number of training set examples. The algorithm, when applied to gene microarray data sets, gives results comparable with other methods. It can be easily extended to the case where the aim is to identify the genes which discriminate between different kinds of tumour tissues (multi-category classification problem). The proposed algorithm can also be used for solving other problem formulations where strictly convex cost functions are used, e.g. the LASSO problem formulation where quadratic loss function is used. In fact, in the case where the quadratic loss function is used, the one-dimensional optimization problem becomes very simple since its solution can be reached in one Newton step. We are currently investigating the extension of the proposed algorithm to sparse versions of kernelized formulations like LS-SVM

(Suykens and Vandewalle, 1999) and kLOGREG (Roth, 2001, <http://www.informatik.uni-bonn.de/~roth/dagm01.pdf>).

ACKNOWLEDGEMENTS

The work of the first author was partially supported by Genome Institute of Singapore. The authors thank Phil Long for his valuable comments and suggestions on the draft of this paper.

REFERENCES

- Alon, U., Barkai, N., Notterman, D., Gish, K., Ybarra, S., Mack, D. and Levine, A. (1999) Broad patterns of gene expression revealed by clustering analysis of tumour and normal colon cancer tissues probed by oligonucleotide arrays. *Cell Biol.*, **96**, 6745–6750.
- Ambrose, C. and McLachlan, G.J. (2002) Selection bias in gene extraction on the basis of microarray gene expression data. *Proc. Natl Acad. Sci. USA*, **99**, 668–674.
- Arrick, A.B. and Derynck, R. (1996) The biological role of transforming growth factor beta in cancer development. In Waxman, J. (ed.), *Molecular Endocrinology of Cancer*, Cambridge University Press, New York, USA, pp. 51–78.
- Bertsekas, D.P. and Tsitsiklis, J.N. (1989) *Parallel and Distributed Computation: Numerical Methods*. Prentice Hall, Englewood Cliffs, NJ, USA.
- Brown, M., Grundy, W., Lin, D., Cristianini, N., Sugnet, C., Furey, T., Ares, Jr., M. and Haussler, D. (2000) Knowledge based analysis of microarray gene expression data using support vector machines. *Proc. Natl Acad. Sci. USA*, **97**, 262–267.
- Furey, T., Cristianini, N., Duffy, N., Bednarski, D., Schummer, M. and Haussler, D. (2000) Support vector machine classification and validation of cancer tissue samples using microarray expression data. *Bioinformatics*, **16**, 906–914.
- Guyon, I., Weston, J., Barnhill, S. and Vapnik, V. (2002) Gene selection for cancer classification using support vector machines. *Mach. Learning*, **46**, 389–422.
- Kos, J. and Schweiger, A. (2002) Cathepsins and cystatins in extracellular fluids—useful biological markers in cancer. *Radiol. Oncol.*, **36**, 176–179.
- Li, Yi, Campbell, C. and Tipping, M. (2002) Bayesian automatic relevance determination algorithms for classifying gene expression data. *Bioinformatics*, **18**, 1332–1339.
- Osborne, M., Presnell, B. and Turlach, B. (2000) On the LASSO and its dual. *J. Comput. Graphical Stat.*, **9**, 319–337.
- Roth, V. (2001) Probabilistic discriminative kernel classifiers for multi-class problems. In Radig, B. and Florczyk, S. (eds), *Pattern Recognition—DAGM'01*, Springer, pp. 246–253.
- Roth, V. (2002) The generalized LASSO: a wrapper approach to gene selection for microarray data. Tech. Rep. IAI-TR-2002-8, University of Bonn, Computer Science III, Bonn, Germany.
- Suykens, J. and Vandewalle, J. (1999) Least squares support vector machine classifiers. *Neural Process. Lett.*, **9**, 293–300.

- Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B*, **58**, 267–288.
- West, M., Blanchette, C., Dressman, H., Huang, E., Ishida, S., Spang, R., Zuzan, H., Olson, Jr., A.J., Marks, J.R. and Nevins, J.R. (2001) Predicting the clinical status of human breast cancer by using gene expression profiles. *Proc. Natl Acad. Sci. USA*, **98**, 11462–11467.
- Weston, J., Mukherjee, S., Chapelle, O., Pontil, M., Poggio, T. and Vapnik, V. (2001) Feature selection for SVMs. *Adv. Neural Inform. Process. Syst.* **13**, 668–674.