

Gene expression

Penalized Cox regression analysis in the high-dimensional and low-sample size settings, with applications to microarray gene expression data

Jiang Gui¹ and Hongzhe Li^{2,*}¹Department of Statistics and ²Rowe Program in Human Genetics, University of California, Davis, CA 95616, USA

Received on November 17, 2004; revised on March 4, 2005; accepted on March 30, 2005

Advance Access publication April 6, 2005

ABSTRACT

Motivation: An important application of microarray technology is to relate gene expression profiles to various clinical phenotypes of patients. Success has been demonstrated in molecular classification of cancer in which the gene expression data serve as predictors and different types of cancer serve as a categorical outcome variable. However, there has been less research in linking gene expression profiles to the censored survival data such as patients' overall survival time or time to cancer relapse. It would be desirable to have models with good prediction accuracy and parsimony property.

Results: We propose to use the L_1 penalized estimation for the Cox model to select genes that are relevant to patients' survival and to build a predictive model for future prediction. The computational difficulty associated with the estimation in the high-dimensional and low-sample size settings can be efficiently solved by using the recently developed least-angle regression (LARS) method. Our simulation studies and application to real datasets on predicting survival after chemotherapy for patients with diffuse large B-cell lymphoma demonstrate that the proposed procedure, which we call the LARS–Cox procedure, can be used for identifying important genes that are related to time to death due to cancer and for building a parsimonious model for predicting the survival of future patients. The LARS–Cox regression gives better predictive performance than the L_2 penalized regression and a few other dimension-reduction based methods.

Conclusions: We conclude that the proposed LARS–Cox procedure can be very useful in identifying genes relevant to survival phenotypes and in building a parsimonious predictive model that can be used for classifying future patients into clinically relevant high- and low-risk groups based on the gene expression profile and survival times of previous patients.

Supplementary information: <http://dna.ucdavis.edu/~hli/LARSCox-Appendix.pdf>

Contact: hli@ucdavis.edu

INTRODUCTION

DNA microarray technology permits simultaneous measurements of expression levels for thousands of genes, which offers the possibility of a powerful, genome-wide approach to the genetic basis of different types of tumors. The genome-wide expression profiles can be used for molecular classification of cancers, for studying varying levels of drug responses in the area of pharmacogenomics

and for predicting different patients' clinical outcomes. The problem of cancer class prediction using gene expression data, which can be formulated as predicting binary or multi-category outcomes, has been studied extensively and has been demonstrated great promise in recent years (Alon *et al.*, 1999; Golub *et al.*, 1999; Alizadeh *et al.*, 2000; Garber *et al.*, 2001; Sorlie *et al.*, 2001). However, there has been less development in relating gene expression profiles to other phenotypes, such as quantitative continuous phenotypes or censored survival phenotypes such as time to cancer recurrence or time to death. Due to the large variability in time to a certain clinical event such as cancer recurrence among cancer patients, studying possibly censored survival phenotypes can be more informative than treating the phenotypes as binary or categorical variables.

The Cox regression model (Cox, 1972) is the most popular method in regression analysis for censored survival data. However, due to the very high-dimensional space of the predictors, i.e. the genes with expression levels measured by microarray experiments, the standard maximum Cox partial likelihood method cannot be applied directly to obtain the parameter estimates. Besides the high-dimensionality, the expression levels of some genes are often highly correlated, which creates the problem of high collinearity. To deal with the problem of collinearity, the most popular approach is the penalized partial likelihood, including both the L_2 penalized estimation, which is often called the ridge regression, and the L_1 penalized estimation, which was proposed by Tibshirani (1996) and is called the least absolute shrinkage and selection operator (Lasso) estimation. Such a Lasso procedure minimizes the negative log partial likelihood subject to the sum of the absolute value of the coefficients being less than a constant s . Compared to the L_2 penalized procedure with constraints on the sum of the square of the coefficients, the Lasso procedure provides a method for variable selection. These penalized procedures have been investigated mainly in the setting where the sample size is greater than the number of predictors. Li and Luan (2003) were the first to investigate the L_2 penalized estimation of the Cox model in the high-dimensional low-sample size settings and applied their method to relate the gene expression profile to survival data. To avoid the inversion of large matrices, they used kernel tricks to reduce the computation to involving only inversion of the matrix of the size of the sample size. They demonstrated that such a procedure can be applied to build a model for predicting patients' future survival times.

*To whom correspondence should be addressed.

One limitation of the L_2 penalized estimation of the Cox model as presented in Li and Luan (2003) is that it uses all the genes in the prediction and does not provide a way of selecting relevant genes for prediction. However, from the biological point of view, one should expect that only a small subset of the genes are relevant to predicting the phenotypes. Including all the genes in the predictive model introduces noise and is expected to lead to poor predictive performance. Due to the high-dimensionality, the standard variable selection methods such as stepwise and backward selection cannot be applied. Tibshirani (1997) further extended the Lasso procedure for variable selection for the Cox proportional hazard models and proposed using the quadratic programming procedure for maximizing the L_1 penalized partial likelihood in order to obtain the parameter estimates. However, such a quadratic programming procedure cannot be applied directly to the settings when the sample size is much smaller than the number of potential predictors, such as in the setting of microarray data analysis.

Recently, Efron *et al.* (2004) proposed the least angle regression (LARS) procedure for variable selection in the linear regression setting. The LARS selects predictors by its current correlation or angle with the response, where the current correlation is defined as the correlation between the predictor and the current residuals. If the active set is defined as the set of indices corresponding to covariates with the greatest absolute current correlations, as the constraint constant s increases, the predictors are chosen one by one without deletion into the active set. The special feature of LARS is that before a new predictor is chosen to the active set as s increases, the corresponding increment of the coefficients only depends on all predictors in the active set. Efron *et al.* (2004) further pointed out the link between LARS and Lasso, showing that LARS can be modified to provide solutions for Lasso. Instead of solving Lasso discretely by quadratic programming, modified LARS can give the whole solution path of all predictors. With this powerful algorithm, Lasso can be extended to perform subset selection in the high-dimension and low-sample settings. We propose in this paper to use the LARS algorithm to obtain the solutions for the Cox model with L_1 penalty in the setting of very high-dimensional covariates such as the gene expression data obtained by microarrays. We call such an estimation procedure the LARS–Cox procedure.

The rest of the paper is organized as follows. We first present the model and briefly review the Lasso estimation of the regression coefficients and present a modified LARS procedure for the Lasso estimation. We then evaluate the LARS–Cox procedure by simulation studies and applications to a real dataset of diffuse large B-cell lymphoma (DLBCL) survival times and gene expression data (Rosenwald *et al.*, 2002). Comparisons of results with methods proposed previously by using simulations and analysis of real datasets of patients with DLBCL are also presented. Finally, we give a brief discussion of the methods and conclusions.

STATISTICAL MODELS AND METHODS

Cox proportional hazards model and Lasso estimation

Suppose we have a sample size of n from which to estimate the relationship between the survival time and the gene expression levels X_1, \dots, X_p of p genes. Due to censoring, for $i = 1, \dots, n$, the i th datum in the sample is denoted by $(t_i, \delta_i, x_{i1}, x_{i2}, \dots, x_{ip})$, where δ_i is the censoring indicator and t_i is the survival time if $\delta_i = 1$ or censoring time if $\delta_i = 0$, and $x_i = \{x_{i1}, x_{i2}, \dots, x_{ip}\}'$ is the vector of the gene expression level of p genes for

the i th sample. Our aim is to build the following Cox regression model for the hazard of cancer recurrence or death at time t :

$$\begin{aligned}\lambda(t) &= \lambda_0(t) \exp(\beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p) \\ &= \lambda_0(t) \exp(\beta' X),\end{aligned}\quad (1)$$

where $\lambda_0(t)$ is an unspecified baseline hazard function, $\beta = \{\beta_1, \dots, \beta_p\}$ is the vector of the regression coefficients and $X = \{X_1, \dots, X_p\}$ is the vector of gene expression levels with the corresponding sample values of $x_i = \{x_{i1}, \dots, x_{ip}\}$ for the i th sample. We define $f(X) = \beta' X$ to be the linear risk score function.

Based on the available sample data, the Cox's partial likelihood (Cox, 1972) can be written as

$$L(\beta) = \prod_{r \in D} \frac{\exp(\beta' x_r)}{\sum_{j \in R_r} \exp(\beta' x_j)},$$

where D is the set of indices of the events (e.g. deaths) and R_r denotes the set of indices of the individuals at risk at time $t_r - 0$. Let $l(\beta) = \log L(\beta)$; then the Lasso estimate of β (Tibshirani, 1996, 1997) can be expressed as

$$\hat{\beta}(s) = \operatorname{argmax} l(\beta), \quad \text{subject to } \sum_{j=1}^p |\beta_j| \leq s,$$

where s is a tuning parameter determining how many covariates with coefficients are zero.

Tibshirani (1997) proposed the following iterative procedure to reformulate this optimization problem with constraints as a Lasso problem for linear regression models. Specifically, let $\eta = \beta' x$, $\mu = \partial l / \partial \eta$, $A = -\partial^2 l / \partial \eta \partial \eta^T$ and $z = \eta + A^- \mu$, where $x = (x_1, \dots, x_n)$ is the gene expression matrix. Here since the sum of all elements in each row (or column) of the matrix A is 0, A is clearly a singular matrix. We can however use the generalized inverse. Alternatively, Tibshirani proposed replacing the information matrix A with a diagonal matrix D , which has the same diagonal elements as A . However, in most of our applications, n is usually small and calculation of the generalized inverse is computationally feasible. In addition, due to the high-dimensionality of the predictors, it is important to make the algorithm as accurate as possible. With this reparameterization, a one-term Taylor series expansion for $l(\beta)$ has the form of

$$(z - \eta)^T A(z - \eta).$$

Although there are multiple choices of A^- , it is easy to show that if $\operatorname{rank}(A) = n - 1$, for any A^- that satisfies $AA^-A = A$ and $z = \eta + A^- \mu$, $(z - \eta)^T A(z - \eta)$ is invariant to the choice of the generalized inverse of A .

The iterative procedure of Tibshirani (1997) involves the following four steps:

- (1) Fix s and initialize $\hat{\beta} = 0$.
- (2) Compute η , μ , A and z based on the current value of $\hat{\beta}$.
- (3) Minimize $(z - \beta' x)^T A(z - \beta' x)$ subject to $\sum |\beta_j| \leq s$.
- (4) Repeat step 2 and 3 until $\hat{\beta}$ does not change.

Tibshirani (1997) proposed using quadratic programming for solving Step 3. However, in the high-dimension and low-sample size setting, i.e. in the case when $p \gg n$, the quadratic programming algorithm cannot be directly applied. We propose in the next section a simple modification of the LARS algorithm of Efron *et al.* (2004) for Step 3.

A LARS–Cox procedure for obtaining the Lasso for the Cox model

We propose an efficient algorithm called LARS–Cox to solve Step 3 of the algorithm, which is based on the recently proposed LARS algorithm (Efron *et al.*, 2004). In their paper, Efron *et al.* (2004) proved that for the linear regression models, starting from zero, the Lasso solution paths grow piecewise linearly in a predictable way and they also proposed the LARS algorithm to efficiently solve the entire Lasso solution path using the same order of computations as a single ordinary least square fit. We propose applying the LARS

algorithm for solving Step 3 of the Lasso algorithm. To do so, we first apply the Cholesky decomposition to obtain $T = A^{1/2}$ such that $T'T = A$, and define $y = Tz$ and $\hat{x} = Tx$; then Step 3 of the iterative procedure presented in the previous section can be rewritten as

$$\begin{aligned} \text{Step 3: } & \text{minimize } (y - \beta'\hat{x})^T (y - \beta'\hat{x}) \\ & \text{subject to } \sum |\beta_j| \leq s, \end{aligned}$$

which is precisely the Lasso of y on \hat{x} and can be efficiently solved by the LARS algorithm for a given s .

To determine the value of the tuning parameter s or the number of genes to be used in the final model, one can choose s which minimizes the cross-validated partial likelihood (CVPL) (Verwij and Van Houwelingen, 1993; Huang and Harrington, 2002), which is defined as

$$CVPL(s) = -\frac{1}{n} \sum_{i=1}^n \left[l(\hat{f}^{(-i)}(s)) - l^{(-i)}(\hat{f}^{(-i)}(s)) \right],$$

where $\hat{f}^{(-i)}(s)$ is the estimate of the score function based on the LARS-Cox procedure with tuning parameter s from the data without the i th subject. The terms $l(f)$ and $l^{(-i)}(f)$ are the log partial likelihoods with all the subjects and without the i th subject, respectively. The optimal value of s is chosen to maximize the sum of the contributions of each subject to the log partial likelihood. This CVPL is a special case of a more general cross-validated likelihood approach for model selections (Smyth, 2001; Van der Laan *et al.*, 2003) and has been demonstrated to perform well in prediction in the context of the penalized Cox regression (Huang and Harrington, 2002).

Evaluation of the predictive performance: time-dependent ROC curves and area under the curves

In order to assess how well the model predicts the outcome, we propose employing the idea of time-dependent receiver-operator characteristics (ROC) curves for censored data and area under the curve (AUC) as our criteria. These methods were recently developed by Heagerty *et al.* (2000) in the context of the medical diagnosis. For a given score function $f(X)$, we can define time-dependent sensitivity and specificity functions as

$$\begin{aligned} \text{sensitivity}(c, t|f(X)) &= Pr\{f(X) > c | \delta(t) = 1\}, \\ \text{specificity}(c, t|f(X)) &= Pr\{f(X) \leq c | \delta(t) = 0\}, \end{aligned}$$

and define the corresponding ROC($t|f(X)$) curve for any time t as the plot of sensitivity($c, t|f(X)$) versus $1 - \text{specificity}(c, t|f(X))$ with the cutoff point c varying, and the AUC as the area under the ROC($t|f(X)$) curve, denoted by AUC($t|f(X)$). Here $\delta(t)$ is the event indicator at time t . A nearest neighbor estimator for the bivariate distribution function is used for estimating these conditional probabilities accounting for possible censoring (Akritas, 1994). Note that a larger AUC at time t based on a score function $f(X)$ indicates better predictability of time to event at time t as measured by sensitivity and specificity evaluated at time t . In our application presented in the next section, we study several different methods of constructing the score function $f(X)$ in the Cox model (1) and compare their predictive performance based on the AUCs.

EVALUATION OF THE METHODS BY SIMULATION STUDIES

We performed simulation studies to evaluate how well the LARS-Cox procedure performs in the high-dimensional and low-sample size settings. We focus on whether the important covariates that are related to survival endpoints can be selected by the LARS-Cox procedure and how well the model can be used for predicting the survival time for future patients.

In our simulation studies, we assume that 20 out of a total of 500 genes are related to time to cancer recurrence through a Cox

regression model with 10 coefficients generated from a uniform $U(-1, -0.1)$ distribution and 10 coefficients generated from a uniform $U(0.1, 1)$ distribution (see first column of Table 1 for the coefficients generated). A Weibull distribution with the shape parameter of 5 and the scale parameter of 2 is used for the baseline hazard function, and a uniform $U(2, 10)$ is used for simulating the censoring times. Based on this setting, we would expect about 40% censoring.

In order to generate gene expression data for 500 predictors (genes), we first generate a 100×500 dataset X from a uniform $U(-1.5, 1.5)$ distribution. We assume that the first 20 genes with expression levels X_1, X_2, \dots, X_{20} are related to patients' risk cancer recurrence through a Cox model. In order to generate gene expression data for the rest of the 480 genes which are not related to the survival but of which some may be correlated with the 20 relevant genes, we first use Gram-Schmidt orthonormalization to construct its normalized orthogonal basis $\{\alpha_1, \dots, \alpha_{20}, \gamma_1, \dots, \gamma_{80}\}$, where $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_{20}\}$ is an orthogonal basis of the linear space A expanded by X_1, X_2, \dots, X_{20} and $\gamma = \{\gamma_1, \gamma_2, \dots, \gamma_{80}\}$ is a set of orthogonal basis of B , which is the orthogonal complement space of A . By Cauchy's inequality, it is easy to show that if $\{\alpha_1, \dots, \alpha_{20}, \gamma_1, \dots, \gamma_{80}\}$ is a set of normalized orthogonal bases, then for any 20×80 matrix T , we have $\text{corr}(\alpha y, (\gamma + \alpha T)x) \leq \lambda / \sqrt{1 + \lambda^2}$, for $\forall x \in R^{80}, y \in R^{20}$, where λ^2 is the largest eigenvalue of $T'T$. Based on this result, we can generate the expression levels of genes which are unrelated to survival from the linear space $C = \{\gamma + \alpha T\}$ with an appropriate choice of the maximum eigenvalue of $T'T$ in order to control the maximum possible correlation between vectors in spaces A and C . We considered the maximum possible correlation of 0, 0.71, 0.82 and 0.87 in our simulations. These numbers are chosen to assess how gradual changes in correlations between irrelevant and non-irrelevant genes affect the LARS-Cox procedure in identifying relevant genes. Note that the actual observed possible correlations between the relevant (to the risk of event) 20 genes and the irrelevant genes are much smaller than these values, and for most simulated datasets, the observed maximum of the pair-wise sample correlations between genes is smaller than half of the theoretical maximum correlation.

Effects of between-gene correlations on identifying relevant genes

For each chosen maximum possible correlation between the relevant genes and non-relevant genes, we generated 100 datasets with a sample size of 100 individuals. For each replication, we applied the LARS-Cox procedure to build a model which included 20 genes by selecting an appropriate s value in the LARS-Cox estimation. Table 1 summarizes the frequencies that the 20 relevant genes are among the first 20 genes that are selected by the LARS-Cox procedure. We observe the following interesting results. First, as expected, the predictors with larger coefficients are more likely to be selected by the LARS-Cox procedure. Second, it is interesting to observe that when the maximum possible correlation between the relevant and non-relevant genes increases, i.e. when the linear space spanned by the non-relevant genes gets close to the linear space expanded by those relevant genes, the chance of the relevant genes with smaller coefficients being selected gets smaller. This is because that at each step, the LARS-Cox procedure only selects the gene with the largest absolute correlation in the model. Of course, the chance of these relevant genes being selected also depends on the sample size. For example, for the maximum possible correlation of 0.85, more

Table 1. Simulation results based on 100 replications^a

Coefficient	Maximum correlation ^b				
	0(SS = 100)	0.71(SS = 100)	0.82(SS = 100)	0.87 (SS = 100)	(SS = 200)
$\beta_1 = 0.19$	50	15	3	0	3
$\beta_2 = 0.95$	100	100	92	75	91
$\beta_3 = 0.96$	100	100	95	80	94
$\beta_4 = 0.91$	100	99	87	71	92
$\beta_5 = 0.19$	53	15	2	0	7
$\beta_6 = 0.25$	60	23	2	0	5
$\beta_7 = 0.69$	100	95	67	45	56
$\beta_8 = 0.33$	88	42	6	4	13
$\beta_9 = 0.34$	88	50	16	6	2
$\beta_{10} = 0.33$	91	53	13	1	4
$\beta_{11} = -0.92$	100	100	92	61	84
$\beta_{12} = -0.16$	40	7	5	1	0
$\beta_{13} = -0.83$	100	98	86	59	84
$\beta_{14} = -0.62$	100	91	58	26	44
$\beta_{15} = -0.65$	100	96	60	32	46
$\beta_{16} = -0.47$	98	76	38	11	22
$\beta_{17} = -0.72$	100	95	70	39	62
$\beta_{18} = -0.24$	66	19	6	5	8
$\beta_{19} = -0.41$	100	68	24	5	14
$\beta_{20} = -0.23$	64	23	3	4	4

^aThe first column shows the true coefficients of the 20 genes which are related to the risk of cancer recurrence. Columns 2–5 show the frequency of each of these 20 relevant genes being selected by the LARS–Cox procedure under four different correlation structures. The sample sizes are 100 patients for all the simulations. For the maximum possible correlation of 0.87, a sample size of 200 patients was also considered and the results are presented in the last column.

^bSS = sample size.

relevant genes are selected if the sample size is increased to 200 (see the last column of Table 1).

Predictive performance and comparison with other methods

We then examined the predictive performance of the proposed method. We simulated a sample size of 100 patients as the training dataset to build the predictive model and evaluated the predictive performance based on another new dataset of 100 patients (test dataset). For each simulation, we generated 500 gene expression levels for each patients with the maximum possible between-gene correlation of 0.82. For each replication, we built a predictive model based on the training set. We applied the CVPL to choose the tuning parameter s used in the model. We also considered three other methods, including the L_2 penalized procedure proposed by Li and Luan (2003), the principal-components based partial Cox regression (PC-PCR) procedure proposed by Li and Gui (2004) and the supervised principal components (SPCA) procedure proposed in Bair and Tibshirani (2004). For each method, we build the model based on the training dataset, and predict the risk scores for the 100 patients in the test set. We repeated this procedure 100 times. We used the time-dependent AUC as a criterion to assess the predictive performance.

Figure 1(a)–(d) shows the averages of the estimated AUCs over 100 replications using the predictive score for the test sets for each method together with the estimated 95% point-wise confidence intervals. The plot indicates a very good predictive performance of the LARS–Cox procedure. The AUC is over 75% at the beginning of the follow-ups and remains high even at later times. As a comparison,

the other three procedures did not perform as well as indicated by the estimated AUCs [Fig. 1(b)–(d)]. It should however be noted that due to the small sample sizes we simulated, the comparisons are not statistically significant as indicated by the slight overlaps of the 95% confidence intervals. Note that both the L_2 penalized procedure and the PC-PCR procedure use all the genes in building the predictive models. Clearly, neither of these procedures performed as well as the LARS–Cox procedure in predicting the survival times for future patients as measured by the AUCs. We also performed the L_2 procedure and the PC-PCR procedure using genes selected based on univariate Cox regression analysis and did not observe any improvement in their predictive performances. The SPCA procedure, although it performs gene selection by univariate analysis, did not perform as well as the LARS–Cox procedure. These results indicate that selecting genes by performing univariate analysis may not be the best choice in building predictive models. In contrast, the LARS–Cox procedure selects genes by considering all the genes together.

As another way of comparing these three different methods, for each replication, we divided the patients in the test set into high- and low-risk groups based on having positive or negative predictive risk scores and tested the statistical significance in the risk of cancer recurrence between the two groups. We observed that for a p -value of $<10^{-5}$, all 100 replications showed significant differences in risk between the high- and low-risk groups defined by the LARS–Cox predicted scores, as compared to only 38, 22 and 36 replications showing significant differences in risk between the high- and low-risk groups defined by the risk scores predicted by the L_2 penalized procedure, the PC-PCR procedure and the SPCA procedure.

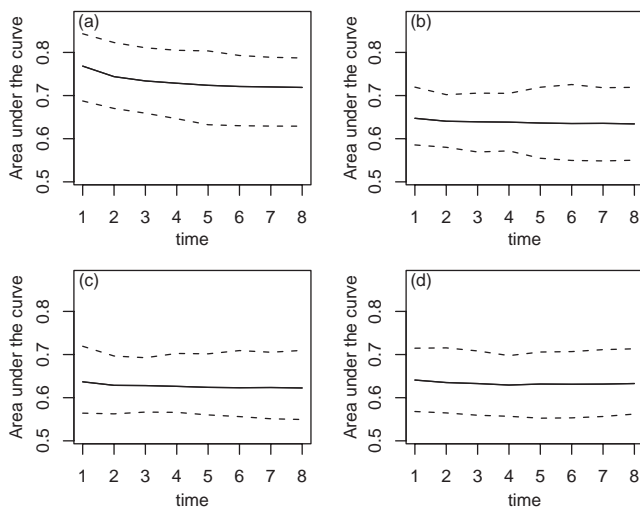


Fig. 1. Results of simulations: AUCs for the test samples based on the LARS–Cox procedure (a), the L_2 penalized estimation (b), the PC-PCR procedure (c) and the SPCA procedure (d). For each plot, the three lines are the average AUCs over 100 replications together with 95% confidence intervals.

In summary, the results from our simulation studies indicate that the LARS–Cox procedure can indeed select genes that are related to censored phenotypes, especially those genes with relatively strong effects, although genes with smaller effects on survival are difficult to identify, especially when the correlations between the gene expression levels are high. When the correlations between the gene expression levels of the relevant genes and non-relevant genes are high, the CVPL procedure tends to select more genes in building the predictive models. However, we observed a better predictive performance of the LARS–Cox procedure than other procedures.

APPLICATION TO PREDICTION OF SURVIVAL TIME OF PATIENTS WITH DLBCL

To further demonstrate the utility of the LARS–Cox procedure in relating microarray gene expression data to censored survival phenotypes, we re-analyzed a recently published dataset of DLBCL by Rosenwald *et al.* (2002). This dataset includes a total of 240 patients with DLBCL, including 138 patient deaths during the follow-ups with a median death time of 2.8 years. Rosenwald *et al.* (2002) divided the 240 patients into a training set of 160 patients and a validation set or test set of 80 patients and built a multivariate Cox model. The variables in the Cox model included the average gene expression levels of smaller sets of genes in four different gene expression signatures together with the gene expression level of BMP6. It should be noted that in order to select the gene expression signatures, they performed a hierarchical clustering analysis for genes across all the samples (including both test and training samples). In order to compare our results with those in Rosenwald *et al.* (2002), we used the same training and test datasets in our analysis.

The gene expression measurements of 7399 genes are available for analysis. However, there are a large number of missing gene expression values in the dataset. Among the 7399 genes, only 434 genes have no missing values. We first applied a nearest neighbor technique (Troyanskaya *et al.*, 2001) to estimate those missing values. Specifically, for each gene, we first identified eight genes which

Table 2. GenBank ID and descriptions of the top 10 genes selected by the LARS–Cox procedure based on the 160 patients in the training dataset^a

GenBank ID	Signature	Description
AA760674		Cytochrome oxidase assembly protein (yeast)
X00452	MHC	Major histocompatibility complex, class II, DQ alpha 1
AA729055	MHC	Major histocompatibility complex, class II, DR alpha
AA714513	MHC	Major histocompatibility complex, class II, DR beta 5
AA729003		T-cell leukemia/lymphoma 1A
AA805575	Germ	Thyroxine-binding globulin precursor
AA598653	Lymph	Osteoblast specific factor 2 (fasciclin I-like)
LC_29222	Lymph	
X59812	Lymph	Cytochrome P450, subfamily XXVIIA polypeptide 1
L19872		Hydrocarbon receptor

^aGerm = Germinal-center B-cell signature, MHC = MHC class II signature, Lymph = Lymph-node signature. Genes AA760674, AA729003 and L19872 do not belong to these signature groups. No description was provided for gene LC_29222 by Rosenwald *et al.* (2002).

are the nearest neighbors according to Euclidean distance. We then filled the remaining with the average of the nearest neighbors. Our method is slightly different from that of Troyanskaya *et al.* (2001) in that the nearest neighbors are not restricted to those 434 genes with no missing data. We also tried the method of Troyanskaya *et al.* (2001) for filling the missing value, and the results of survival time prediction with the two methods were very close.

Selection of genes related to risk of death

We applied the LARS–Cox procedure to first build a predictive model using the training data of 160 patients and all the 7399 features. As the tuning parameter increases, more genes are selected and these genes are chosen in order of their relevances in predicting survival. The genes entered first in the model would provide a good list of candidate genes for further investigation. Table 2 shows the GenBank ID and a brief description of the first 10 genes selected. It is interesting to note that seven of these genes belong to the three gene expression signature groups defined in Rosenwald *et al.* (2002). These three signature groups include Germinal-center B-cell signature, MHC class II signature and Lymph-node signature. No genes in the proliferation signature group defined by Rosenwald *et al.* (2002) were among the top 10 genes selected by LARS–Cox. However, ribosomal protein S12 from the proliferation group was among the top 20 gene selected by our method.

The other three genes which do not belong to the signature groups of Rosenwald *et al.* (2002) may also be related to lymphoma or risk of death from lymphoma; however, evidence for their direct involvement in any mechanism is currently lacking. It should also be noted that there is always a possibility that genes are selected because of coexpression with other genes, or for reasons that cannot be explained mechanistically. Among these three genes, AA729003 is a protein coding TCL1A gene which has been demonstrated to be a powerful oncogene and when it is over-expressed in both B and T cells, it predominantly yields mature B cell lymphomas (Pekarsky *et al.*,

1999). The gene L19872 is a Aryl hydrocarbon receptor (AHR), which is a ligand-activated transcription factor involved in the regulation of biological responses to planar aromatic hydrocarbons. The AHR has been shown to regulate xenobiotic-metabolizing enzymes such as cytochrome P450, which belongs to the lymph-node signature group. Finally, the gene AA760674 is a COX15 homolog, which is the terminal component of the mitochondrial respiratory chain that catalyzes the electron transfer from reduced cytochrome c to oxygen (Petruzzella *et al.*, 1998). It has been reported that mutation in the COX15 gene can cause Leigh syndrome (Oquendo *et al.*, 2004); however, its involvement in cancers is not clear.

Evaluation of the predictive performance

We also examined how well the model built by the LARS–Cox procedure predicts the survival of a future patient. Using the training set of 160 patients, we built a predictive Cox model with the LARS–Cox procedure using the CVPL to select the optimal tuning value s . The minimum CVPL was obtained when $s = 0.28$, which corresponds to selecting four genes in the model. We also observed that the CVPL value increases by only 0.001 when the tuning parameter s increases from 0.28 to 0.33, which corresponds to nine genes in the model. As a matter of fact, for s ranging from 0.28 to 0.33, the predictive performances of the resulting models are very comparable. We chose the most parsimonious model with four genes. These four genes are AA805575, LC_29222, X00452 and X59812 (see Table 2 for a description of these four genes), belonging to three of the four signature groups defined in Rosenwald *et al.* (2002).

We obtained the estimates of the coefficients of these four genes using the LARS–Cox procedure, denoted by vector $\hat{\beta}$. The estimated coefficients for all four genes were negative, indicating that high expression levels of these genes reduce the risk of death among the patients with DLBCL. This agrees with what Rosenwald *et al.* (2002) has found (see Table 2 of their paper). Based on the estimated model with four genes, we estimated the risk scores ($f(X) = \hat{\beta}'X$) for the 80 patients in the test dataset based on their gene expression levels of the four genes in the predictive model. The time-dependent AUCs for 1–20 years after diagnosis based on the estimated scores for the patients in the test set are around 0.67 in the first 10 years of follow-up, indicating a reasonable predictive performance.

To further examine whether clinically relevant groups can be identified by the model, we used zero as a cutoff point of the risk scores and divided the test patients into two groups based on whether they have positive or negative risk scores. Figure 2(b) shows the Kaplan–Meier curves for the two groups of patients, showing very significant differences (p -value = 0.0004) in overall survival between the high-risk group (36 patients) and the low-risk group (44 patients).

A comparison with other methods

As a comparison, we also analyzed the lymphoma dataset using three other methods, the PCR method of Li and Gui (2004), the L_2 penalized method of Li and Luan (2003) and the SPCA method of Bair and Tibshirani (2004). Figure 2(b)–(d) shows the survival curves of the two groups of patients in the test dataset defined by the scores estimated by each of the three methods. We observe that the two risk groups defined by the LARS–Cox estimated model showed more significant difference in risk of death than the groups defined by the other three models (p -value of 0.0004 versus 0.003, 0.003 and 0.034). Finally, the AUCs based on the risk scores estimated by the LARS–Cox procedure are also higher than those from the other

three procedures; however, the results are not statistically significant (the time-dependent AUCs and their 95% confidence intervals are provided in the Supplemental material).

As another evaluation of the methods, we performed another set of analyses using the training and testing datasets as defined in Bair and Tibshirani (2004); (data available from <http://www-stat.stanford.edu/~tibs/superpc/staudt.html>) for the lymphoma dataset. Again, if we used zero as the cutoff point of the predicted scores to divide the 80 patients in the test set into high- and low-risk groups, we observed that the LARS–Cox procedure gives a slightly more significant difference in risk between the two groups. The log-rank test p -values are 0.007, 0.06, 0.007 and 0.015 for the LARS–Cox, L_2 penalized procedure, PC-PCR procedure and the SPCA procedure, respectively, again indicating that the LARS–Cox procedure performs well on this new training/testing partition of the lymphoma dataset. The corresponding survival curves are provided in the Supplemental material.

DISCUSSION AND CONCLUSIONS

It is clinically relevant and very important to predict patients' time to cancer relapse or time to death due to cancer after treatment using gene expression profiles of the cancerous cells prior to the treatment. Powerful statistical methods for such prediction allow microarray gene expression data to be used most efficiently. In this paper, we have proposed and studied the LARS–Cox procedure for censored survival data in order to identify important predictive genes for survival using microarray gene expression data. To solve the computational difficulty, we modified the recently developed LARS procedure (Efron *et al.*, 2004) to obtain the solutions for the Lasso estimation of the Cox model. Since the risk of cancer recurrence or death due to cancer may result from the interplay between many genes, methods which can utilize the data of many genes, as in the case of our proposed procedure, are expected to show better performance in predicting risk. Our simulation studies demonstrated that the procedure can indeed be used to identify genes which are related to censored survival outcomes and to build a parsimonious model for predicting future patients' survival. We have also demonstrated the applicability of our methods by analyzing time to death of the diffuse large B-cell lymphoma patients and obtained satisfactory results, as evaluated by both applying the model to the test dataset and time-dependent ROC curves.

While we did not compare the new procedure with all the other procedures available, we did compare the LARS–Cox procedure with several other previously proposed methods in predictive performance and found that the new procedure performed better than the L_2 penalized procedure, PC-PCR procedure and the SPCA procedure (Li and Luan, 2003; Li and Gui, 2004; Bair and Tibshirani, 2004) in predicting the future patients' survival. We would however expect that the LARS–Cox procedure performs better than other dimension-reduction based procedures such as the partial least squares (Park *et al.*, 2002) or the principal components Cox regression because the LARS–Cox procedure automatically selects and utilizes only the relevant genes in building the predictive model. A comprehensive comparison of different methods warrants further research. It is worth mentioning that the L_1 penalized regression was also demonstrated to perform better than other procedures in the settings of microarray gene expression data and linear models (Segal *et al.*, 2003).

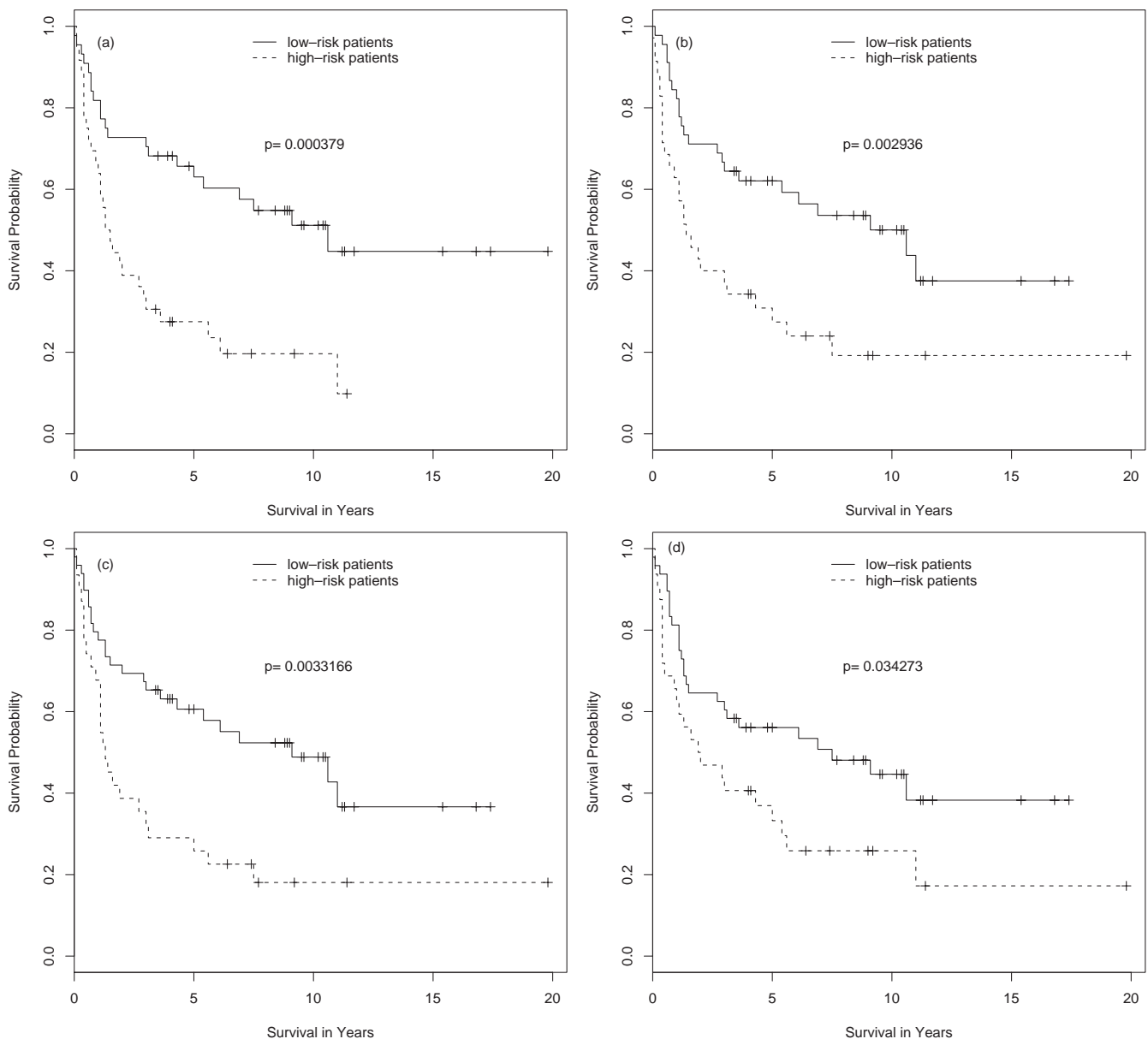


Fig. 2. Results of analyses of the lymphoma dataset: the Kaplan–Meier curves for the high- and low-risk groups defined by the estimated scores with zero as the cutoff for the 80 patients in the test dataset. The scores are estimated based on the models estimated by the LARS–Cox procedure (a), L_2 penalized procedure (b), the PC–PCR procedure (c) and the SPCA procedure (d). The number of patients in the high-risk groups are 36, 35, 31 and 32 respectively.

The proposed LARS–Cox procedure has no computational or methodological limitation in terms of the number of genes that can potentially be used in building models for the prediction of patients' time to clinical event. The method can in principle select $n - 1$ genes, where n is the sample size, although the cross-validation procedure often results in a much smaller number of predictors in the model. However, when the number of predictors is close to the sample size, there is a risk of over-fitting. In addition, as pointed out by Osborne *et al.* (1998), as s increases, when the number of nonzero coefficients gets close to the number of observations, Lasso may not have a unique solution. This implies that the number of genes selected

by the procedure cannot be more than the sample size. In addition, the LARS–Cox tends to select only one variable from a group of highly correlated genes. These points have also been pointed out by Zou and Hastie (2003) for the Lasso. If the LARS–Cox procedure is used mainly for selecting important and relevant genes, one may want to include all these highly correlated genes, if one of them is selected. However, if the goal is to build a model with a good predictive accuracy, this should not be a problem since simple models are preferred for the scientific insight into the relationship between survival and gene expressions. One way to extend the LARS–Cox procedure in order to identify correlated genes is that at each LARS

variable selection step, we selected not only one single gene with the largest absolute current inner product, but a group of such genes with similar current inner products. An alternative is to use the elastic net penalty as recently proposed by Zou and Hastie (2004) for the penalized estimation. How well such extensions perform in prediction of future survival time deserves further investigation.

The LARS–Cox procedure assumes the Cox proportional hazards model, which is the most popular model for censored survival data. However, the proportional hazards assumption may not hold for gene expression data or for all diseases. It is possible to develop robust procedures under mis-specified proportional hazards models along the lines of Lin and Wei (1989). In addition, model checking techniques analogous to those of Lin *et al.* (1993) can be derived. As an alternative, we can consider similar L_1 penalized estimation for the accelerated failure time models (Wei, 1992) or more general semi-parametric transformation models (Cheng *et al.*, 1995). We are currently pursuing these alternative models.

In summary, an important application of microarray technology is to relate gene expression profiles to various clinical phenotypes of patients such as time to cancer recurrence or overall survival time. The statistical model built to relate gene expression profiles to the censored survival time should have the property of high predictive accuracy and parsimony. The proposed LARS–Cox procedure in this paper can be very useful in building a parsimonious predictive model that can be used for classifying future patients into clinically relevant high- and low-risk groups based on the gene expression profile and survival times of previous patients. The procedure can also be applied to select important genes which are related to patients' survival outcome.

ACKNOWLEDGEMENTS

This research was supported by NIH grant ES009911. We thank the two reviewers for their very helpful comments.

REFERENCES

- Akritis, M.G. (1994) Nearest neighbor estimation of a bivariate distribution under random censoring. *Ann. Statist.*, **22**, 1299–1327.
- Alizadeh, A.A. *et al.* (2000) Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling. *Nature*, **403**, 503–511.
- Alon, U. *et al.* (1999) Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proc. Natl Acad. Sci., USA*, **96**, 6745–6750.
- Bair, E. and Tibshirani, R. (2004) Semi-supervised methods for predicting patient survival from gene expression papers. *PLoS Biol.*, **2**, 5011–5022.
- Cheng, S.C. *et al.* (1995) Analysis of transformation models with censored data. *Biometrika*, **82**, 835–845.
- Cox, D.R. (1972) Regression models and life-tables. *J. R. Statist. Soc. B*, **34**, 187–220.
- Efron, B. *et al.* (2004) Least angle regression. *Ann. Statist.*, **32**, 407–499.
- Garber, M.E. *et al.* (2001) Diversity of gene expression in adenocarcinoma of the lung. *Proc. Natl Acad. Sci. USA*, **98**, 13784–13789.
- Golub, T.R. *et al.* (1999) Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science*, **286**, 531–537.
- Heagerty, P.J. *et al.* (2000) Time dependent ROC curves for censored survival data and a diagnostic marker. *Biometrics*, **56**, 337–344.
- Huang, J. and Harrington, D. (2002) Penalized partial likelihood regression for right-censored data with bootstrap selection of the penalty parameter. *Biometrics*, **58**, 781–791.
- Li, H. and Gui, J. (2004). Partial Cox regression analysis for high-dimensional microarray gene expression data. *Bioinformatics*, **20** (suppl. 1), i208–i215.
- Li, H. and Luan, Y. (2003) Kernel Cox regression models for linking gene expression profiles to censored survival data. *Pacific Symp. Biocomput.*, **8**, 65–76.
- Lin, D.Y. and Wei, L.J. (1989) The robust inference for the Cox proportional hazards model. *J. Am. Statist. Assoc.*, **84**, 1074–1078.
- Lin, D.Y. *et al.* (1993) Checking the Cox model with cumulative sums of martingale-based residuals. *Biometrika*, **80**, 557–572.
- Oquendo, C.E. *et al.* (2004) Functional and genetic studies demonstrate that mutation in the COX15 gene can cause Leigh syndrome. *J. Med. Genet.*, **41**, 540–544.
- Osborne, M.R., Presnell, B. and Turlach, B.A. (1998). On the Lasso and its dual. *Research Report*, Department of Statistics, University of Adelaide.
- Park, P.J. *et al.* (2002) Linking expression data with patient survival times using partial least squares. *Bioinformatics*, **18**, S120–S127.
- Pekarsky, Y. *et al.* (1999) Abnormalities at 14q32.1 in T cell malignancies involve two oncogenes. *Proc. Natl Acad. Sci. USA*, **96**, 2949–2951.
- Petrzzella, V. *et al.* (1998). Identification and characterization of human cDNAs specific to BCS1, PET112, SCO1, COX15, and COX11, five genes involved in the formation and function of the mitochondrial respiratory chain. *Genomics*, **54**, 494–504.
- Rosenwald, A. *et al.* (2002) The use of molecular profiling to predict survival after chemotherapy for diffuse large-B-Cell lymphoma. *N. Eng. J. Med.*, **346**, 1937–1947.
- Segal, M.R. *et al.* (2003) Regression approaches for microarray data analysis. *J. Comput. Biol.*, **10**, 961–980.
- Smyth, P. (2001) Model selection of probabilistic clustering using cross-validated likelihood. *Statist. Comput.*, **10**, 63–72.
- Sorlie, T. *et al.* (2001) Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. *Proc. Natl Acad. Sci.*, **98**, 10869–10874.
- Tibshirani, R. (1996) Regression shrinkage and selection via the Lasso. *J. R. Statist. Soc. B*, **58**, 267–288.
- Tibshirani, R. (1997) The Lasso method for variable selection in the Cox model. *Statist. Med.*, **16**, 385–395.
- Troyanskaya, O. *et al.* (2001) Missing value estimation methods for DNA microarrays. *Bioinformatics*, **17**, 520–525.
- Van der Laan, M.J., Dudoit, S. and Keles, S. (2003) Asymptotic optimality of likelihood-based cross validation. *Technical Report*, Division of Biostatistics, University of California, Berkeley.
- Verwij, P.J.M. and Van Houwelingen, J.C. (1993) Cross validation in survival analysis. *Statist. Med.*, **12**, 2305–2314.
- Wei, L.J. (1992) The accelerated failure time model. a useful alternative to the Cox regression model in survival analysis. *Statist. Med.*, **11**, 1871–1879.
- Zou, H. and Hastie, T. (2004) Regression shrinkage and selection via the elastic net, with applications to microarrays. *J. R. Statist. Soc. B*, in press.