

THE DANTZIG SELECTOR FOR CENSORED LINEAR REGRESSION MODELS

Yi Li, Lee Dicker and Sihai Dave Zhao

University of Michigan, Rutgers University, University of Pennsylvania

Abstract: The Dantzig variable selector has recently emerged as a powerful tool for fitting regularized regression models. To our knowledge, most work involving the Dantzig selector has been performed with fully-observed response variables. This paper proposes a new class of adaptive Dantzig variable selectors for linear regression models when the response variable is subject to right censoring. This is motivated by a clinical study to identify genes predictive of event-free survival in newly diagnosed multiple myeloma patients. Under some mild conditions, we establish the theoretical properties of our procedures, including consistency in model selection and the optimal efficiency of estimation. The practical utility of the proposed adaptive Dantzig selectors is verified via extensive simulations. We apply our new methods to the aforementioned myeloma clinical trial and identify important predictive genes.

Key words and phrases: Buckley-James imputation, censored linear regression, dantzig selector, oracle property.

1. Introduction

Technical advances in biomedicine have produced an abundance of high-throughput data. This has resulted in major statistical challenges and brought attention to the variable selection and estimation problem, where the goal is to discover relevant variables among many potential candidates and obtain high prediction accuracy. For example, variable selection is essential when performing gene expression profiling for cancer patients in order to better understand cancer genomics and design effective gene therapy (Anderson et al. (2006); Pawitan et al. (2005)).

Penalized likelihood methods, represented by the LASSO, have been extensively studied as a means of simultaneous estimation and variable selection (Tibshirani (1996)). It is known that the LASSO estimator can discover the right sparse representation of the model (Zhao and Yu (2006)), but the LASSO estimator is in general biased (Zou (2006)), especially when the true coefficients are relatively large. Several remedies, including the smoothly clipped absolute deviation (SCAD) (Fan and Li (2001)), and the adaptive LASSO (ALASSO) (Zou (2006)), have been proposed to discover the sparsity of the true models, while

producing consistent estimates for nonzero regression coefficients. Though these methods differ to a great extent, they are all cast in the framework of penalized likelihoods or penalized objective functions.

More recently the Dantzig selector (Candés and Tao (2007)), has emerged to enrich the class of regularization techniques. The Dantzig selector can be implemented as a linear programming problem, making the computational burden manageable. Though under some general conditions the LASSO and Dantzig may produce the same solution path (James, Radchenko, and Lv (2008)), they differ conceptually in that the Dantzig stems directly from an estimating equation, whereas the LASSO stems from a likelihood or an objective function.

The Dantzig selector has been most thoroughly studied with fully observed outcome variables. But in many clinical studies, the outcome variable, e.g. the CD4 counts in an AIDS trial or patients' survival times, may not be fully observed. In a myeloma clinical trial that motivates this research, the goal was to identify genes predictive of a patient's event-free survival.

While the vast majority of work in variable selection for censored outcome data has focused on the Cox proportional hazards model (e.g., Tibshirani (1997); Li and Luan (2003); Li and Gui (2004); Gui and Li (2005a,b); Antoniadis, Fryzlewicz, and Letue (2010)), a linear regression model offers a viable alternative as it directly links the outcome to the covariates. Hence, its regression coefficients have an easier interpretation than those of the Cox model, especially when the response does not pertain to a survival time. Some recent work on regularized linear regression models for censored data can be found in Ma, Kosorok, and Fine (2006), Johnson, Lin, and Zeng (2008), Wang et al. (2008), Cai, Huang, and Tian (2009), Engler and Li (2009), and Johnson (2009).

Most of these methods operate under the penalization framework. Given that a censored linear regression does not pertain to a likelihood function, the Dantzig selector may be a natural choice. Johnson, Long, and Chung (2011) approached the problem using a penalized estimation equation approach, but Johnson (2009) noted that their procedure gives only an *approximate* root- n consistent estimator. To our knowledge, it remains unclear whether the Dantzig selector can also be used to estimate linear regression models with censored outcome data. Johnson, Long, and Chung (2011) studied such a procedure but did not provide theoretical support. It is therefore of interest to (i) explore the utility of the Dantzig selector in censored linear regression models, (ii) rigorously evaluate its theoretical properties, and (iii) compare its numerical properties to similar methods developed under the lasso/penalization-based framework.

We propose a new class of Dantzig variable selectors for linear regression models when the response variable is subject to right censoring. Dicker (2011) proposed the adaptive Dantzig selector for the linear model, and here we develop a similar procedure for use with censored outcomes. First, our method carries

out simultaneous variable selection and estimation, and is motivated from the estimating equation perspective, which may be important for some semiparametric models whose likelihood functions are difficult to specify. Second, the proposed selectors possess the oracle property when the tuning parameters follow some appropriate rates, providing the theoretical justification for the proposed procedures. Third, the complex regularization problem has been reduced to a linear programming problem, resulting in computationally efficient algorithms.

The rest of the paper is structured as follows. Section 2 reviews the Dantzig selector for noncensored linear regression models, as well as its connection with penalized likelihood methods. Section 3 considers its extension to linear regression models when the response variable is subject to censoring. In Section 4, we discuss the large sample properties and prove the consistency of variable selection and the optimal efficiency of the estimators. We discuss the choice of tuning parameters for the finite sample situations in Section 5. We report on numerical simulations in Section 6, and apply the proposal to a myeloma study in Section 7. We conclude the paper with a discussion in Section 8. Proofs are relegated to a web supplement.

2. Penalized Likelihood Methods and the Dantzig Selector

We begin by considering a linear regression model with p predictors

$$Y_i = \sum_{j=1}^p X_{ij}\beta_j + \epsilon_i, \quad (2.1)$$

where ϵ_i are iid mean zero residuals for $i = 1, \dots, n$. Denote the truth by $\beta_0 = (\beta_{01}, \dots, \beta_{0p})$ and set $A = \{j : \beta_{0j} \neq 0\}$. The goal of the model selection in this context is to identify A , often referred to as the “true model.”

A variable selector $\hat{\beta}$ for β_0 is considered to have reasonable large sample behavior if (i) $P(\{j : \hat{\beta}_j \neq 0\} = A) \rightarrow 1$ as the sample size $n \rightarrow \infty$, and (ii) $\sqrt{n}(\hat{\beta}_A - \beta_A) \rightarrow N(0, \Sigma^*)$ where β_A is the subvector of β extracted by the subset A of $\{1, \dots, p\}$ and Σ^* is some $|A| \times |A|$ covariance matrix (here, $|A|$ denotes the cardinality of the set A). Property (i) is often considered to be the consistency property, while property (ii) involves the efficiency of the estimator. If properties (i) and (ii) hold and Σ^* is optimal (by some criterion), the variable selection procedure is said to have the oracle property.

Concise notation is to be used for referring to sub-vectors and sub-matrices. For a subset $T \subset \{1, \dots, p\}$ and $\beta \in \mathbf{R}^p$, let $\beta_T = (\beta_j)_{j \in T}$ be the $|T| \times 1$ vector whose entries are those of β indexed by T . For an $n \times p$ matrix, \mathbf{X} , \mathbf{X}_T is the $n \times |T|$ matrix whose columns are those of \mathbf{X} that are indexed by T . Additionally, let \mathbf{X}_i and $\mathbf{X}_{\cdot j}$ denote the i^{th} row and j^{th} column of \mathbf{X} , respectively, for $i = 1, \dots, n$ and $j = 1, \dots, p$. Denote the complement of T in $\{1, \dots, p\}$ by \bar{T} .

Let $\|\boldsymbol{\beta}\|_r = (\sum_{i=1}^p |\beta_i|^r)^{1/r}$ for $0 < r < \infty$, $\|\boldsymbol{\beta}\|_0 = \#\{j : \beta_j \neq 0\}$ and $\|\boldsymbol{\beta}\|_\infty = \max_{1 \leq j \leq p} |\beta_j|$, and also let $\text{sgn}(\boldsymbol{\beta})$ have $\text{sgn}(\boldsymbol{\beta})_j = \text{sgn}(\beta_j)$ with $\text{sgn}(0) = 0$. For a diagonal matrix $\mathbf{W} = \text{diag}(w_1, \dots, w_p)$, we take $\mathbf{W}_{T,T} = \text{diag}(w_j; j \in T)$.

2.1. Penalized likelihood methods

The LASSO is a benchmark penalized likelihood procedure. It works by minimizing an L_2 loss function $\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$ subject to an L_1 constraint: $\|\boldsymbol{\beta}\|_1 = \sum_j |\beta_j| \leq s$, where $\mathbf{Y} = (Y_1, \dots, Y_n)$ is the response vector, \mathbf{X} is the $n \times p$ design matrix, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ is the $p \times 1$ vector of coefficients and s is a nonnegative tuning parameter. Equivalently, the LASSO estimate can be obtained by minimizing

$$\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{j=1}^p |\beta_j|, \quad (2.2)$$

where λ is a nonnegative tuning parameter. The LASSO performs variable selection, but in general does not possess the oracle property. A remedy is to utilize an adaptive LASSO that minimizes

$$\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{j=1}^p w_j |\beta_j|, \quad (2.3)$$

where w_j is a data-driven weight.

2.2. Adaptive Dantzig selector

The Dantzig selector also belongs to the class of regularization methods in regression. The estimator is the solution to

$$\begin{aligned} &\text{minimize} && \|\boldsymbol{\beta}\|_1 \\ &\text{subject to} && \|\mathbf{X}'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\|_\infty \leq \lambda. \end{aligned}$$

Thus it strikes a balance between nearly solving the score equation, $\mathbf{X}'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) = 0$ and minimizing the L_1 norm of $\boldsymbol{\beta}$. Connections between the Dantzig selector and the LASSO have been discussed in James, Radchenko, and Lv (2008), where it is shown that under some general conditions the Dantzig selector and the LASSO produce the same solution path. In general the Dantzig selector does not have the oracle property.

As a remedy, a modified Dantzig selector, analogous to the adaptive LASSO, was proposed by Dicker (2011):

$$\begin{aligned} &\text{minimize} && \sum_j w_j |\beta_j| \\ &\text{subject to} && |\mathbf{X}'_{\cdot j}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})| \leq \lambda w_j, j = 1, \dots, p. \end{aligned}$$

The adaptive Dantzig selector and adaptive LASSO are also related: the adaptive Dantzig selector and adaptive LASSO are equivalent to instances of the Dantzig

selector and LASSO, respectively, where \mathbf{X} is replaced with $\mathbf{X}\mathbf{W}^{-1}$, $\boldsymbol{\beta}$ is replaced with $\mathbf{W}\boldsymbol{\beta}$, and $\mathbf{W} = \text{diag}(w_1, \dots, w_p)$. The key to the adaptive Dantzig selector is to strike a balance between minimizing the weighted L_1 norm, which promotes sparsity, and approximately solving the weighted normal equations. Weights w_j in the adaptive Dantzig selector need to be chosen according to the principles that determine weights in the adaptive LASSO. When the response vector \mathbf{Y} is fully observed, Dicker (2011) established the oracle property of the adaptive Dantzig selector for an appropriately chosen tuning parameter λ . It is unclear, however, whether this property hold when the response \mathbf{Y} is subject to censoring.

3. Adaptive Dantzig Selector for Censored Linear Regression

Consider a slightly modified version of (2.1),

$$Y_i = \mathbf{X}_i' \boldsymbol{\beta} + \epsilon_i,$$

where $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})'$ is the covariate vector for the i^{th} subject and ϵ_i are iid with an unspecified distribution denoted by $F(\cdot)$, with survival function $S(\cdot) = 1 - F(\cdot)$. The mean of ϵ_i , denoted by α , is not necessarily 0. Here $\boldsymbol{\beta}_0$ denotes the true $\boldsymbol{\beta}$ and $A = \{j; \beta_{0j} \neq 0\}$ is the true model. Suppose that Y_i may be right censored by a competing observation C_i and that only $Y_i^* = Y_i \wedge C_i$ and $\delta_i = I(Y_i^* = Y_i)$ are observed for each subject. We assume that Y_i is independent of C_i conditional on \mathbf{X}_i . When the response variable pertains to survival time, both Y_i and C_i are commonly measured on the log scale and the model is called the accelerated failure time model (Kalbfleisch and Prentice (2002)).

Let $e_i(\boldsymbol{\beta}) = Y_i^* - \boldsymbol{\beta}' \mathbf{X}_i$, and consider

$$\tilde{Y}_i(\boldsymbol{\beta}) = E(Y_i | Y_i^*, \delta_i, \mathbf{X}_i, \boldsymbol{\beta}) = Y_i^* + (1 - \delta_i) \frac{\int_{e_i(\boldsymbol{\beta})}^{\infty} S(s, \boldsymbol{\beta}) ds}{S(e_i(\boldsymbol{\beta}), \boldsymbol{\beta})}.$$

Clearly,

$$E\left\{\tilde{Y}_i(\boldsymbol{\beta}) | \mathbf{X}_i, \boldsymbol{\beta}\right\} = \alpha + \mathbf{X}_i' \boldsymbol{\beta}.$$

The Buckley-James estimating equation is

$$\sum_{i=1}^n (X_{ij} - \bar{X}_j) \left\{ \hat{Y}_i(\boldsymbol{\beta}) - \mathbf{X}_i' \boldsymbol{\beta} \right\} = 0, \quad j = 1, \dots, p, \quad (3.1)$$

where $\bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}$ for $j = 1, \dots, p$ and

$$\hat{Y}_i(\boldsymbol{\beta}) = Y_i^* + (1 - \delta_i) \frac{\int_{e_i(\boldsymbol{\beta})}^{\infty} \hat{S}(s, \boldsymbol{\beta}) ds}{\hat{S}\{e_i(\boldsymbol{\beta}), \boldsymbol{\beta}\}} \quad (3.2)$$

is the empirical version of $\tilde{Y}_i(\boldsymbol{\beta})$. Here, $\hat{S}(\cdot, \boldsymbol{\beta})$ is the one-sample Nelson-Aalen estimator based on $(e_i(\boldsymbol{\beta}), \delta_i)$,

$$\hat{S}(t, \boldsymbol{\beta}) = \exp \left\{ - \sum_{i=1}^n \int_{-\infty}^t \frac{dN_i(u, \boldsymbol{\beta})}{\bar{Y}(u, \boldsymbol{\beta})} \right\}, \quad (3.3)$$

where $N_i(u, \boldsymbol{\beta}) = I\{e_i(\boldsymbol{\beta}) \leq u, \delta_i = 1\}$ and $\bar{Y}(u, \boldsymbol{\beta}) = \sum_i I\{e_i(\boldsymbol{\beta}) \geq u\}$. Under mild conditions, Lai and Ying (1991) have shown that the Buckley-James estimator that solves (3.1) is \sqrt{n} -consistent. To facilitate the ensuing development, note that (3.1) can be written as

$$\mathbf{X}'\mathbf{P}_n\{\hat{\mathbf{Y}}(\boldsymbol{\beta}) - \mathbf{X}\boldsymbol{\beta}\} = 0, \quad (3.4)$$

where $\mathbf{P}_n = \mathbf{I}_n - \mathbf{1}\mathbf{1}'/n$, \mathbf{I}_n is an $n \times n$ identity matrix, $\mathbf{1}$ is an $n \times 1$ vector with all elements 1, and $\hat{\mathbf{Y}}(\boldsymbol{\beta}) = (\hat{Y}_1(\boldsymbol{\beta}), \dots, \hat{Y}_n(\boldsymbol{\beta}))'$.

Solving (3.4) does not directly render an automatic variable selection procedure, but the adaptive Dantzig selector presents an appealing solution to the problem. Applying it to (3.4) gives

$$\begin{aligned} & \text{minimize} && \sum_j w_j |\beta_j| \\ & \text{subject to} && |\mathbf{X}'_j \mathbf{P}_n \{\hat{\mathbf{Y}}(\boldsymbol{\beta}) - \mathbf{X}\boldsymbol{\beta}\}| \leq \lambda w_j, \quad j = 1, \dots, p. \end{aligned}$$

This is no longer a simple linear programming problem and one strategy is to use an iterative algorithm, as in Wang et al. (2008), though such methods have numerical and theoretical difficulties (Johnson (2009)).

3.1. Low-dimensional setting

We propose the use of a \sqrt{n} -consistent initial estimator $\hat{\boldsymbol{\beta}}^{(0)}$ to construct an imputed version of the true response \mathbf{Y} , $\hat{\mathbf{Y}}(\hat{\boldsymbol{\beta}}^{(0)})$, then to employ a version of the adaptive Dantzig selector replacing (3.4) with $\mathbf{X}'\mathbf{P}_n\{\hat{\mathbf{Y}}(\hat{\boldsymbol{\beta}}^{(0)}) - \mathbf{X}\boldsymbol{\beta}\} = 0$. If $p < n$, we can obtain $\boldsymbol{\beta}_0$ using unpenalized Buckley-James estimation, or rank-based procedures with Gehan or logrank weights (Jin, Lin, Wei, and Ying (2003)). A similar one-step imputation strategy was used by Johnson (2009), though the imputed $\hat{\mathbf{Y}}(\hat{\boldsymbol{\beta}}^{(0)})$ was used to construct a loss function with a LASSO penalty.

Our version of the adaptive Dantzig selector is to

$$\begin{aligned} & \text{minimize} && \sum_j w_j |\beta_j| \\ & \text{subject to} && |\mathbf{X}'_j \mathbf{P}_n \{\hat{\mathbf{Y}}(\hat{\boldsymbol{\beta}}^{(0)}) - \mathbf{X}\boldsymbol{\beta}\}| \leq \lambda w_j, \quad j = 1, \dots, p. \end{aligned} \quad (3.5)$$

The w_j are data driven weights and should be chosen to vary inversely with the magnitude of β_{0j} . For instance, we can take $w_j = |\hat{\beta}_j^{(0)}|^{-\gamma}$ for some $\gamma > 0$. Then when $|\hat{\beta}_j^{(0)}|$ is large, (3.5) requires us to nearly solve the j^{th} score equation, where $\hat{\mathbf{Y}}(\hat{\boldsymbol{\beta}}^{(0)})$ is treated as a fully observed outcome vector, that heavily penalizes non-zero estimates of β_{0j} when $|\hat{\beta}_j^{(0)}|$ is small. When $\delta_i \equiv 1$ for all i , (3.5) reduces

to the adaptive Dantzig selector for linear regression models and the result is an effective variable selection procedure. Still, censoring presents difficulties that need investigation.

3.2. High-dimensional setting

In the high-throughput datasets that now characterize modern medicine, typically $p > n$, so it is difficult to obtain an initial \sqrt{n} -consistent estimate $\hat{\beta}^{(0)}$. Our strategy is to first reduce the number of the covariates to be smaller than n using a sure screening procedure. We then calculate the $\hat{\beta}^{(0)}$ using the retained covariates and proceed as in Section 3.1.

We employ the screening procedure of Zhu et al. (2011) that can provide sure screening for any single-index model, including the AFT model. To choose the number of covariates to retain after screening, we follow the combined soft- and hard-thresholding rule of Zhu et al. (2011) that chooses up to $n/\log(n)$ covariates using a procedure involving randomly generated auxiliary variables.

Recently, Johnson, Long, and Chung (2011) studied Buckley-James estimation using the Dantzig selector with an initial \sqrt{n} -consistent estimator for the high-dimensional problem. It is similar to our proposal, but our work has several advantages. We propose an adaptive Dantzig selector that has practical and theoretical advantages over the nonadaptive version. By using the procedure of Zhu et al. (2011) to select the subset of important covariates, we can take advantage of their sure screening property that states that, under certain conditions on the design matrix, the probability that the selected covariates contains the truly important covariates approaches 1.

4. Theoretical Results

A difficulty in extending the Dantzig selector to the censored regression setting is that $\hat{Y}_i(\hat{\beta}^{(0)})$ is a surrogate for the unobserved outcome Y_i . In the ensuing development, we first quantify the “distance” between the surrogate and the true outcome, then show that their average difference is bounded by a random variable of order $n^{-1/2}$. Given this random bound, we show that the existing Dantzig selector results for the non-censored case can be extended to the censored case, leading to the oracle property.

4.1. Quantify the “Distance” between the imputed and “True” responses

Before stating the main result of the section, we need a lemma useful for bounding the difference between the surrogate and the true outcomes.

Lemma 1. *Assume conditions 1–4 of Ying (1993, p.80), and suppose the derivative of the hazard function $\lambda(s)$ with respect to s is continuous for $-\infty < s < \infty$. Then,*

$$\begin{aligned} \hat{S}(s_1, \beta_1) - S(s_0) &= S(s_0) \{ (\beta_1 - \beta_0)^T \mathcal{A}(s_0, \beta_0) - \lambda(s_0, \beta_0)(s_1 - s_0) \\ &\quad + n^{-1/2} Z(s_0) \} + o\{ \max(n^{-1/2}, |s_1 - s_0| + \|\beta_1 - \beta_0\|) \}, \end{aligned}$$

with probability 1 uniformly for any $(s_1, \beta_1) \in \mathcal{B} = \{(s, \beta) : |s - s_0| + \|\beta - \beta_0\| < Cn^{-1/2}\}$, where $C > 0$ is any arbitrary constant, \mathcal{A} is a $p \times 1$ nonrandom function, $\lambda(s)$ is the hazard function for $S(s)$ and the stochastic process $Z(s)$ is a version of $W(v(s))$, and where $W(\cdot)$ is the Wiener process and $v(\cdot)$ is defined at (S1.2) in the web supplement.

Proposition 1. *Under the regularity conditions of Lemma 1, $(1/n) \sum_{i=1}^n \mathbf{X}_i \{\hat{Y}_i(\hat{\beta}^{(0)}) - Y_i\} = O_p(n^{-1/2})$ if $\hat{\beta}^{(0)} = \beta_0 + O_p(n^{-1/2})$.*

The result implies $(\mathbf{P}_n \mathbf{X})'(\hat{\mathbf{Y}} - \mathbf{Y}) = O_p(n^{1/2})$, where \mathbf{X} is replaced by its centralized version; this facilitates the proof of consistency of model selection. As the validity of Proposition 1 requires $\hat{\beta}^{(0)}$ to be \sqrt{n} -consistent, taking $\hat{\beta}^{(0)}$ to be the Buckley-James estimate, which is \sqrt{n} -consistent suffices.

4.2. Selection-consistent adaptive Dantzig selector

We use $\hat{\mathbf{Y}}$ to denote $\hat{\mathbf{Y}}(\hat{\beta}^{(0)})$. Observe that the adaptive Dantzig selector for data with a censored response, (3.5), can be rewritten as

$$\begin{aligned} &\text{minimize} && \|\mathbf{W}\beta\|_1 \\ &\text{subject to} && \|\mathbf{Z}'(\hat{\mathbf{Y}} - \mathbf{Z}\mathbf{W}\beta)\|_\infty \leq \lambda, \end{aligned} \tag{4.1}$$

where $\mathbf{W} = \text{diag}(w_1, \dots, w_p)$ and $\mathbf{Z} = \mathbf{P}_n \mathbf{X} \mathbf{W}^{-1}$. We refer to the solution $\hat{\beta}$ as the selection-consistent adaptive Dantzig selector (SADS). The optimization problem (SADS) is a linear programming problem that, with its dual problem, allows $\hat{\beta}$ to be characterized in terms of primal and dual feasibility and complementary slackness conditions.

Lemma 2. *If there is $\hat{\mu} \in \mathbf{R}^p$ such that,*

$$\|\mathbf{Z}'(\hat{\mathbf{Y}} - \mathbf{Z}\mathbf{W}\hat{\beta})\|_\infty \leq \lambda, \tag{4.2}$$

$$\|\mathbf{Z}'\mathbf{Z}\hat{\mu}\|_\infty \leq 1, \tag{4.3}$$

$$\hat{\mu}'\mathbf{Z}'\mathbf{Z}\mathbf{W}\hat{\beta} = \|\mathbf{W}\hat{\beta}\|_1, \tag{4.4}$$

$$\hat{\mu}'\mathbf{Z}'(\hat{\mathbf{Y}} - \mathbf{Z}\mathbf{W}\hat{\beta}) = \lambda\|\hat{\mu}\|_1, \tag{4.5}$$

then $\hat{\beta} \in \mathbf{R}^p$ solves (SADS).

Proposition 2. *Suppose β_0 is the true parameter value, $A = \{j; \beta_{0j} \neq 0\}$, and $(1/n)\mathbf{X}'\mathbf{P}_n\mathbf{X}$ converges in probability to some positive definite matrix. If*

$$\frac{\lambda}{\sqrt{n}}w_j \xrightarrow{P} \infty \text{ if } j \notin A \text{ and } \lambda w_j = O_P(\sqrt{n}) \text{ if } j \in A,$$

then, with probability tending to 1, a solution to (SADS), $\hat{\beta}$, and the corresponding $\hat{\mu}$ are given by

$$\hat{\mu}_A = (\mathbf{Z}'_A \mathbf{Z}_A)^{-1} \text{sgn}(\beta_0)_A, \quad (4.6)$$

$$\hat{\mu}_{\bar{A}} = 0, \quad (4.7)$$

$$\begin{aligned} \hat{\beta}_A &= \mathbf{W}_{A,A}^{-1} \left\{ (\mathbf{Z}'_A \mathbf{Z}_A)^{-1} \mathbf{Z}'_A \hat{\mathbf{Y}} - \lambda (\mathbf{Z}'_A \mathbf{Z}_A)^{-1} \text{sgn}(\hat{\mu})_A \right\} \\ &= (\mathbf{X}'_A \mathbf{P}_n \mathbf{X}_A)^{-1} \mathbf{X}'_A \mathbf{P}_n \hat{\mathbf{Y}} - \lambda (\mathbf{X}'_A \mathbf{P}_n \mathbf{X}_A)^{-1} \mathbf{W}_{A,A} \text{sgn}(\hat{\mu})_A, \end{aligned} \quad (4.8)$$

$$\hat{\beta}_{\bar{A}} = 0. \quad (4.9)$$

Corollary 1 (Consistency of model selection). *If the conditions of Proposition 2 hold and $\hat{\beta}$ is any sequence of solutions to (SADS), then $P(\{j; \hat{\beta}_j \neq 0\} = A) \rightarrow 1$.*

To ensure that the conditions in Proposition 2 hold, one selects data-driven weights w_j and an appropriate λ . Examples of weights and λ such that these conditions hold include $w_j = |\hat{\beta}^{(0)}|^{-\gamma}$, where $\hat{\beta}^{(0)}$ is \sqrt{n} -consistent for β_0 and $\gamma > 0$, and λ such that $n^{-1/2}\lambda = O(1)$ and $n^{(\gamma-1)/2}\lambda \rightarrow \infty$. Proposition 2 makes no uniqueness claims about solutions to (SADS), but it can be shown that in “most” cases (SADS) has a unique solution (Dicker (2011)).

4.3. Oracle adaptive Dantzig selector

The estimator at (4.8) and (4.9) solves (SADS) in probability. This expression may be leveraged to obtain the large-sample distribution of \sqrt{n} -standardized (SADS) estimates. However, though the solution to (SADS) is selection consistent, it may not achieve optimal efficiency. To remedy this, we propose the oracle adaptive Dantzig selector (OADS).

Let $T = \{j; \hat{\beta}_j \neq 0\}$ be the index set of non-zero estimated coefficients from the SADS estimator $\hat{\beta}$. Take OADS estimator $\hat{\beta}^{(0,T)}$ so that $\hat{\beta}_T^{(0,T)} = 0$ and $\hat{\beta}_T^{(0,T)}$ is the Buckley-James estimate obtained by solving (3.4) with \mathbf{X} replaced by \mathbf{X}_T . This is similar to the Gauss-Dantzig selector of Candés and Tao (2007), in which ordinary linear regression is performed on the covariates selected by the Dantzig selector.

Proposition 3 (Oracle property). *Assume the conditions of Proposition 2. Let $T = \{j; \hat{\beta}_j \neq 0\}$, where $\hat{\beta}$ is the SADS estimator for β_0 , and let $\beta_{0,A}$ be the*

non-zero subvector of β_0 . If $\hat{\beta}^{(0,A)}$ so that $\hat{\beta}_A^{(0,A)} = 0$ and $\hat{\beta}_A^{(0,A)}$ is the Buckley-James estimate obtained by solving (3.4) with \mathbf{X} replaced by \mathbf{X}_A , then the OADS estimator $\hat{\beta}^{(0,T)}$ satisfies

$$P\left(\hat{\beta}^{(0,T)} = \hat{\beta}^{(0,A)}\right) \rightarrow 1, \quad \sqrt{n}\left(\hat{\beta}_A^{(0,T)} - \beta_{0,A}\right) \rightarrow N(0, \Sigma),$$

where $\Sigma = \Omega^{-1}\Lambda\Omega^{-1}$ with

$$\begin{aligned} \Omega &= \int_{-\infty}^{-\infty} \left[\Gamma^{(2)}(t, \beta_0) - \frac{\{\Gamma^{(1)}(t, \beta_0)\}^{\otimes 2}}{\Gamma^{(0)}(t, \beta_0)} \right] \frac{\int_t^\infty (1 - F(s)) ds}{1 - F(t)} \left\{ \frac{d \log f(t)}{dt} \right. \\ &\quad \left. + \frac{f(t)}{1 - F(t)} \right\} dF(t), \\ \Lambda &= \int_{-\infty}^{-\infty} \left[\Gamma^{(2)}(t, \beta_0) - \frac{\{\Gamma^{(1)}(t, \beta_0)\}^{\otimes 2}}{\Gamma^{(0)}(t, \beta_0)} \right] \left\{ \frac{\int_t^\infty (1 - F(s)) ds}{1 - F(t)} \right\}^2 dF(t), \end{aligned}$$

and $\Gamma^{(r)}(t, \beta_0)$ for $r = 0, 1$ defined in (S1.1) in the web supplement.

The Σ in Proposition 3 is the asymptotic variance of the Buckley-James estimator given the true subset of covariates; see Lai and Ying (1991)).

In practice we propose using the covariance matrix estimated from the second-stage Buckley-James fit to estimate the covariance of the nonzero components of $\hat{\beta}^{(0,T)}$, while we set the zero components to have zero variance. This ad-hoc estimator ignores the variability coming from the imputation of $\hat{Y}(\hat{\beta}^{(0)})$ as well as the variability from the first-stage SADS model selection, and so in general underestimates the true variance of the OADS estimator. However, as the probability of selecting the true model increases, this variance estimator approaches Σ , the variance of the oracle estimator.

5. Tuning Parameter Selection

One needs to select an appropriate tuning parameter λ in order to obtain good performance. For the uncensored linear regression (2.1), Tibshirani (1996) and Fan and Li (2001) proposed the generalized cross-validation (GCV) statistic

$$GCV^*(\lambda) = \frac{AR(\lambda)}{\{1 - d(\lambda)/n\}^2}$$

where $AR(\lambda)$ is the average residual sum of squares $(1/n)\|\mathbf{Y} - \mathbf{X}\hat{\beta}(\lambda)\|_2^2$, $\hat{\beta}(\lambda)$ is the estimate of β under λ and $d(\lambda)$ is the effective number of parameters (Zou, Hastie, and Tibshirani (2007)).

When the data are censored, we adopt an inverse reweighting scheme to account for it. Assume the censoring C_i are iid and have a common survival

function G_i , a reasonable assumption for clinical trials where most censoring is administrative. As suggested by Johnson, Lin, and Zeng (2008), we approximate the unobserved $AR(\lambda)$ by

$$\widehat{AR}(\lambda) = \frac{\sum_{i=1}^n \delta_i \{Y_i^* - \hat{\alpha}^{(0)} - \mathbf{X}_i' \hat{\beta}(\lambda)\}^2 / \hat{G}(Y_i^*)}{\sum_{i=1}^n \delta_i / \hat{G}(Y_i^*)},$$

where $\hat{G}(\cdot)$ is the Kaplan-Meier estimator for $G(\cdot)$, and $\hat{\alpha}^{(0)} = (1/n) \sum_{i=1}^n \{Y_i(\hat{\beta}^{(0)}) - \mathbf{X}_i' \hat{\beta}^{(0)}\}$. Conditional on (Y_i, C_i, \mathbf{X}_i) , the expected value of $\delta_i / G(Y_i^*)$ is one, and hence, the expected values of the numerator and the denominator of $\widehat{AR}(\lambda)$ are the expected values of $\sum_{i=1}^n \{Y_i - \hat{\alpha}^{(0)} - \mathbf{X}_i' \hat{\beta}(\lambda)\}^2$ and n , respectively. This implies that $\widehat{AR}(\lambda)$ and $AR(\lambda)$ have the same limit, justifying the utility of the inverse reweighting scheme. To obtain an estimate of the effective number of parameters for the SADS estimator, we follow Zou, Hastie, and Tibshirani (2007). The expressions (4.8) and (4.9) suggest that $\hat{d}(\lambda) = \text{trace}\{\mathbf{X}_T(\mathbf{X}_T' \mathbf{P}_n \mathbf{X}_T)^{-1} \mathbf{X}_T' \mathbf{P}_n\} = \|\mathbf{T}\|_0$, where $T = \{j; \hat{\beta}_j \neq 0\}$, is a consistent estimator for $d(\lambda)$. In our data analysis and simulation studies, we select λ to yield the smallest

$$GCV(\lambda) = \frac{\widehat{AR}(\lambda)}{\{1 - \hat{d}(\lambda)/n\}^2}. \quad (5.1)$$

Similar GCV schemes have been proposed by Wang et al. (2008) and Johnson, Lin, and Zeng (2008).

6. Simulations and Comparisons

6.1. Simulation set-up

We examined the finite sample performance of the proposed methods in low- and high-dimensional settings. Mimicking the simulation setup of Tibshirani (1997) and Cai, Huang, and Tian (2009), for $i = 1, \dots, n$ we generated the true response Y_i (after the exponential transform) from the exponential distribution with rate $\lambda_i = \exp(-\beta_0' \mathbf{X}_i)$, so $Y_i = \beta_0' \mathbf{X}_i + e_i$. In the low-dimensional setting, we let $p = 9$, and to model weak and moderate associations between the predictors and the response we took $\beta_0 = (0.35, 0.35, 0, 0, 0, 0.35, 0, 0, 0)'$ and $\beta_0 = (0.7, 0.7, 0, 0, 0, 0.7, 0, 0, 0)'$. In the high-dimensional setting, we let $p = 10,000$ and considered $\beta_0 = (\beta_0^{low}, \beta_0^{low}, \mathbf{0})$, where β_0^{low} are the 9×1 -dimensional true parameter vectors from the low-dimensional setting.

We generated covariates $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})'$ from a multivariate normal with mean zero and a compound symmetry covariance matrix $\Sigma = (\sigma_{jj'})_{p \times p} = (\rho)$, and the e_i with the standard extreme value distribution. In low dimensions, we took ρ to be 0, 0.5, and 0.9, corresponding to zero, moderate, and

strong collinearity among the predictors. In high dimensions, good performance is difficult to achieve with a p as large as ours, so we took ρ equal to 0, 0.3, or 0.5.

The censoring variable C_i (after exponential transform) was generated from a uniform $[0, \xi]$, where ξ was chosen to achieve about 40% censoring. The initial estimate $\hat{\beta}^{(0)}$ was obtained via the Buckley-James procedure. We simulated sample sizes of $n = 50$ or $n = 200$ and generated 200 independent datasets under each setting.

6.2. Comparisons of competing methods

For each scenario, we evaluated the selection-consistent adaptive Dantzig selector $\hat{\beta}$ (SADS) and the oracle adaptive Dantzig selector $\hat{\beta}^{(0,T)}$. For the OADS estimator, we tuned the SADS estimator and then fit a Buckley-James estimate to the selected covariates. We used the unpenalized Buckley-James procedure to obtain the initial \sqrt{n} -consistent estimates $\hat{\beta}^{(0)}$.

We also evaluated the adaptive penalized Buckley-James estimator (APBJ) of Johnson (2009) that uses $\hat{\beta}^{(0)}$ to impute outcomes $\hat{Y}(\hat{\beta}^{(0)})$ and then applies the adaptive LASSO penalty to the least-square loss function constructed using the $\hat{Y}(\hat{\beta}^{(0)})$. Here we followed Johnson (2009) and used the Gehan estimator (Gehan (1965)) to obtain the initial $\hat{\beta}^{(0)}$. In the high-dimensional setting, we used the screening procedure of Zhu et al. (2011) on the APBJ estimator and our adaptive Dantzig selectors.

We evaluated the accuracy and precision of the parameter estimates based on the mean squared errors $\text{MSE} = E(\|\hat{\beta} - \beta_0\|^2)$. We recorded the zero regression coefficients incorrectly set to non-zero and non-zero regression coefficients incorrectly set to zero, giving the average number of false positives (FP) and false negatives (FN). We also recorded the probability, across the 200 simulations, of selecting the correct model. Finally, we compared the predictive abilities of the fitted models by estimating their C-statistics (Uno et al. (2011)) on independent test datasets.

The results for the low-dimensional setting are summarized in Table 1. The Dantzig selector-based estimators appeared to have better model selection performances. When $n = 50$, all methods performed alike in terms of model selection, and had a difficult time selecting the correct model. When $n = 200$ and $\beta_{0j} = 0.7$, the SADS and OADS estimators were able to select the correct model up to 74% of the time. The APBJ method of Johnson (2009) performed worse than the others except when $\rho = 0.9$.

The APBJ method of Johnson (2009) appeared to have the best estimation accuracy in general, as its average mean squared errors were usually lower than those of OADS and SADS. The OADS estimator could outperform the SADS

Table 1. Comparisons of methods with different signal strengths in low dimensions.

Method	$\beta_{0j} = 0.7$					$\beta_{0j} = 0.35$				
	MSE	FP	FN	% Correct	C-stat	MSE	FP	FN	% Correct	C-stat
$n = 50, \rho = 0$										
SADS	0.46	1.36	0.24	0.24	0.72	0.42	1.58	1.36	0.04	0.59
OADS	0.54	1.36	0.24	0.24	0.72	0.55	1.58	1.36	0.04	0.59
APBJ	0.44	1.42		0.23	0.24	0.72	0.39	1.60	1.23	0.04
$n = 200, \rho = 0$										
SADS	0.11	0.56	0.01	0.74	0.75	0.10	0.80	0.34	0.32	0.64
OADS	0.08	0.56	0.01	0.74	0.75	0.11	0.80	0.34	0.32	0.64
APBJ	0.11	0.87		0	0.64	0.75	0.10	0.82	0.34	0.30
$n = 50, \rho = 0.5$										
SADS	1.03	1.60	0.55	0.10	0.78	0.61	1.27	1.62	0	0.66
OADS	1.14	1.60	0.55	0.10	0.78	0.85	1.27	1.62	0	0.66
APBJ	0.94	1.77		0.48	0.10	0.78	0.54	1.29	1.51	0.02
$n = 200, \rho = 0.5$										
SADS	0.18	0.88	0.06	0.57	0.81	0.18	0.93	0.74	0.12	0.69
OADS	0.17	0.88	0.06	0.57	0.81	0.21	0.93	0.74	0.12	0.69
APBJ	0.18	1.10		0.02	0.48	0.81	0.16	0.96	0.55	0.20
$n = 50, \rho = 0.9$										
SADS	4.31	1.90	1.41	0	0.82	2.25	1.52	1.96	0	0.71
OADS	5.04	1.90	1.41	0	0.81	3.32	1.52	1.96	0	0.70
APBJ	3.68	2.12	1.26	0	0.82	1.85	1.52	1.88	0.01	0.71
$n = 200, \rho = 0.9$										
SADS	1.20	1.55	0.78	0.04	0.83	0.59	0.65	1.98	0	0.72
OADS	1.35	1.55	0.78	0.04	0.83	0.86	0.65	1.98	0	0.72
APBJ	1.05	1.47	0.71	0.08	0.83	0.53	0.74	1.77	0.01	0.72

estimator, but apparently only when the probability of selecting the true model was sufficiently high. In such situations the OADS outperformed the APBJ estimator, but when this was not the case, such as in the simulation settings with $n = 50$, it was detrimental to fit a Buckley-James estimator to the covariates selected by SADS.

The methods did not exhibit appreciable differences in predictive abilities. The selection performance and mean squared errors of all methods improved with increasing sample size and degraded with increasing correlation between the covariates. The predictive abilities were not affected much by the sample size, and improved with increasing correlation.

The results for the high-dimensional setting are summarized in Table 2. All methods performed poorly in variable selection. When $n = 200$ and $\beta_{0j} = 0.7$, they were able to achieve fairly good estimation accuracy, as the MSE's of the three methods were lower than $\|\beta_0\|_2^2 = 2.94$ for $\rho = 0$ and $\rho = 0.3$. The APBJ estimator gave the most accurate parameter estimates. We see that the different methods gave fitted models with very similar predictive abilities. When $n = 50$

Table 2. Comparisons of methods with different signal strengths in high dimensions.

Method	$\beta_{0j} = 0.7$					$\beta_{0j} = 0.35$				
	MSE	FP	FN	% Correct	C-stat	MSE	FP	FN	% Correct	C-stat
$n = 50, \rho = 0$										
SADS	4.81	7.66	5.50	0	0.55	1.87	7.78	5.91	0	0.54
OADS	4.97	7.66	5.50	0	0.55	2.00	7.78	5.91	0	0.54
APBJ	4.66	8.13	5.51	0	0.55	1.74	8.01	5.91	0	0.54
$n = 200, \rho = 0$										
SADS	1.37	20.01	0.79	0	0.72	1.26	21.34	3.58	0	0.55
OADS	1.64	20.01	0.79	0	0.70	1.43	21.34	3.58	0	0.55
APBJ	1.34	19.31	0.79	0	0.72	1.18	22.02	3.51	0	0.55
$n = 50, \rho = 0.3$										
SADS	5.92	7.18	5.79	0	0.73	2.08	5.96	5.92	0	0.67
OADS	6.01	7.18	5.79	0	0.73	2.21	5.96	5.92	0	0.67
APBJ	5.66	7.30	5.80	0	0.73	1.88	6.29	5.92	0	0.67
$n = 200, \rho = 0.3$										
SADS	2.40	15.62	2.86	0	0.81	1.15	12.64	4.66	0	0.71
OADS	2.57	15.62	2.86	0	0.81	1.20	12.64	4.66	0	0.71
APBJ	2.30	13.45	2.87	0	0.81	1.03	10.57	4.69	0	0.71
$n = 50, \rho = 0.5$										
SADS	6.90	7.23	5.87	0	0.80	2.52	5.98	5.94	0	0.72
OADS	6.97	7.23	5.87	0	0.80	2.71	5.98	5.94	0	0.72
APBJ	6.50	7.50	5.88	0	0.80	2.29	6.05	5.94	0	0.72
$n = 200, \rho = 0.5$										
SADS	3.52	17.32	4.06	0	0.85	1.41	12.33	5.28	0	0.76
OADS	3.70	17.32	4.06	0	0.84	1.48	12.33	5.28	0	0.76
APBJ	3.34	15.70	4.07	0	0.84	1.25	9.54	5.35	0	0.76

and $\rho = 0$, the C-statistics were very low because of the noise involved in the screening step, but with larger n and higher ρ the C-statistics were around 80% even though few of the truly important covariates were selected.

We also studied the performance of our covariance estimator for the OADS estimators. As the SADS stage does a better job of selecting the true model, our variance estimator approaches the true variance of the oracle estimator. In Table 3 we compare a few average estimated covariance matrices to their corresponding empirical oracle covariance matrices in low dimensions. When $n = 200$, $\beta_{0j} = 0.7$, and $\rho = 0$, the OADS estimator selected the correct model 74% of the time, so the two covariance matrices were similar. With $n = 50$, $\beta_{0j} = 0.7$, and $\rho = 0.5$, even when the model was selected only 10% of the time, our covariance estimator still performed fairly well. In the worst case setting of $n = 50$, $\beta_{0j} = 0.35$, and $\rho = 0$, however, when the OADS estimator never selected the correct model, our estimator was very different from the truth. Thus while our ad-hoc proposal is reasonable for easy or moderately difficult settings, a more appropriate variance estimator is an interesting subject for further research.

Table 3. Covariance estimators.

	Oracle cov.			Ave. OADS cov.		
	β_1	β_2	β_6	β_1	β_2	β_6
	$n = 200, \beta_{0j} = 0.7, \rho = 0, 74\%$ correct					
β_1	0.014	0.004	0.004	0.016	0.002	0.002
β_2	0.004	0.013	0.004	0.002	0.016	0.002
β_6	0.004	0.004	0.016	0.002	0.002	0.016
	$n = 50, \beta_{0j} = 0.7, \rho = 0.5, 10\%$ correct					
β_1	0.127	-0.036	-0.031	0.100	-0.010	-0.011
β_2	-0.036	0.098	-0.020	-0.010	0.108	-0.014
β_6	-0.031		-0.020 0.094	-0.011	-0.014	0.098
	$n = 50, \beta_{0j} = 0.35, \rho = 0.9, 0\%$ correct					
β_1	0.457	-0.206	-0.224	0.144	-0.018	-0.017
β_2	-0.206	0.406	-0.171	-0.018	0.132	-0.011
β_6	-0.224	-0.171	0.416	-0.017	-0.011	0.158

Table 4. Validation C-statistics on TT3.

Method	Model size	C-statistic
SADS	4	0.6067
OADS	4	0.6160
APBJ	13	0.6255

7. Example of Myeloma Patients' Survival Prediction

Multiple myeloma is a progressive hematologic (blood) disease, characterized by excessive numbers of abnormal plasma cells in the bone marrow and overproduction of intact monoclonal immunoglobulin. Myeloma patients are typically characterized by wide clinical and pathophysiologic heterogeneities, with survival ranging from a few months to more than 10 years. Gene expression profiling of multiple myeloma patients has offered an effective way of understanding the cancer genomics and designing gene therapy. Identifying risk groups with a high predictive power could contribute to selecting patients for personalized medicine.

We studied event-free survival from newly diagnosed multiple myeloma patients enrolled in trials UARK 98-026 and UARK 2003-33 (Zhan et al. (2006), Shaughnessy et al. (2007)). The trials compared the results of two treatment regimes, total therapy II (TT2) and total therapy III (TT3). There were 340 patients in TT2, with 191 events and an average follow-up of 47.1 months, and 214 patients in TT3, with 55 events and an average follow-up of 35.6 months. Gene expression values for 54,675 probesets were measured for each subject using Affymetrix U133Plus2.0 microarrays. We retrieved the data from the MicroArray Quality Control Consortium II (Shi et al. (2010)) GEO entry (GSE24080).

We used our adaptive Dantzig selector methods to develop risk scores by fitting AFT models to the TT2 patients. We then estimated the C-statistics

Table 5. Parameter estimates by various selectors.

Probeset	Gene name	SADS	OADS	APBJ
205072_s_at	XRCC4	0.057	0.285	0.140
208966_x_at	IFI16	-0.836	-0.678	-0.485
225450_at	AMOTL1	-0.042	-0.180	-0.082
233750_s_at	C1orf25	-0.207	-0.386	-0.191
204700_x_at	C1orf107			-0.055
1565951_s_at	CHML			-0.135
1568907_at	Unknown			0.080
201897_s_at	CKS1B			0.009
222437_s_at	VPS24			0.297
222443_s_at	RBM8A			0.255
209052_s_at	WHSC1			-0.040
228817_at	ALG9			0.236
225834_at	FAM72A /// FAM72B /// GCUD2			-0.046

(Uno et al. (2011)) of those models on the TT3 patients, and compared the results to the APBJ estimator. Table 4 contains the results, and we see that the SADS, OADS, and APBJ estimators have similar predictive performances, with validation C-statistics of around 61%. Our adaptive Dantzig selectors achieved this using 4 probesets, while the APBJ estimator used 13. Table 5 reports the final models and the parameter estimates.

8. Discussion

Several issues merit further investigations. Our asymptotic setup in this paper is that the number of predictors is fixed while the sample size approaches infinity. We have appealed to the sure screening procedure of Zhu et al. (2011), but an asymptotic theory with a diverging p seems to be more applicable to problems involving a huge number of predictors, such as microarray analysis and document/image classification.

More research is needed on the evaluation of the variation of the estimator for small or moderate sample size. We proposed an ad-hoc variance estimator that gives reasonable performance when the signal-to-noise ratio is not too weak. Another possibility is to use a perturbation resampling technique, as in Minnier, Tian, and Cai (2011), though this lacks theoretical justification when applied to Dantzig selector-type regularization.

A potential advantage of the Dantzig selector over penalized likelihood methods such as LASSO is that it can be extended to settings in which no explicit likelihoods or loss functions are available, and may be more computationally and theoretically appealing than the penalized estimating equation method of Johnson, Lin, and Zeng (2008). We believe that our work can be extended to handle Dantzig selectors in the framework of more general estimating equations.

Acknowledgement

This work was partially supported by U.S. National Cancer Institute grant R01 CA95747.

References

- Anderson, E., Miller, P., Ilsley, D., Marshall, W., Khvorova, A., Stein, C. and Benimetskaya, L. (2006). Gene profiling study of G3139- and Bcl-2-targeting siRNAs identifies a unique G3139 molecular signature. *Cancer Gene Therapy* **13**, 406-414.
- Antoniadis, A., Fryzlewicz, P. and Letue, F. (2010). The Dantzig selector in Cox's proportional hazards model. *Scand. J. Statist.* **37**, 531-552.
- Cai, T., Huang, J. and Tian, L. (2009). Regularized estimation for the accelerated failure time model. *Biometrics* **65**, 394-404.
- Candés, E. and Tao, T. (2007). The Dantzig selector: statistical estimation when p is much larger than n . *Ann. Statist.* **35**, 2313-2351.
- Dicker, L. (2011). Regularized regression methods for variable selection and estimation. Ph.D. thesis, Harvard University, USA.
- Engler, D. and Li, Y. (2009). Survival analysis with high dimensional covariates: an application in microarray studies. *Statist. Appl. Genet. Mol. Biol.* **8**.
- Fan, J. and Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.* **96**, 1348-1360.
- Gehan, E. A. (1965). A generalized Wilcoxon test for comparing arbitrarily single-censored samples. *Biometrika* **90**, 341-353.
- Gui, J. and Li, H. (2005a). Threshold gradient descent method for censored data regression, with applications in pharmacogenomics. *Pacific Symposium on Biocomputing* **10**, 272-283.
- Gui, J. and Li, H. (2005b). Penalized Cox regression analysis in the highdimensional and low-sample size settings, with applications to microarray gene expression data. *Bioinformatics* **21**, 3001-3008.
- James, G., Radchenko, P. and Lv, J. (2008). DASSO: connections between the Dantzig selector and lasso. *J. Roy. Statist. Soc. Ser. B* **71**, 127-142.
- Jin, Z., Lin, D., Wei, L. J. and Ying, Z. (2003). Rank-based inference for the accelerated failure time model. *Biometrika* **90**, 341-353.
- Johnson, B., Lin, D. and Zeng, D. (2008). Penalized estimating functions and variable selection in semiparametric regression models. *J. Amer. Statist. Assoc.* **103**, 672-680.
- Johnson, B. A. (2009). On lasso for censored data. *Electronic J. Statist.* **3**, 485-506.
- Johnson, B. A., Long, Q. and Chung, M. (2011). On path restoration for censored outcomes. *Biometrics* **67**, 1379-1388.
- Kalbfleisch, J. and Prentice, R. (2002). *The Statistical Analysis of Failure Time Data*. Wiley, New York.
- Lai, T. and Ying, Z. (1991). Large sample theory of a modified Buckley-James estimator for regression analysis with censored data. *Ann. Statist.* **19**, 1370-1402.
- Li, H. and Gui, J. (2004). Partial Cox regression analysis for high-dimensional microarray gene expression data. *Bioinformatics* **20**, 208-215.
- Li, H. and Luan, Y. (2003). Kernel Cox regression models for linking gene expression profiles to censored survival data. *Pacific Symposium of Biocomputing* **8**, 65-76.
- Ma, S., Kosorok, M. and Fine, J. (2006). Additive risk models for survival data with high-dimensional covariates. *Biometrics* **62**, 202-210.

- Minnier, J., Tian, L. and Cai, T. (2011). A perturbation method for inference on regularized regression estimates. *J. Amer. Statist. Assoc.* **106**, 1371-1382.
- Pawitan, Y. et al. (2005). Gene expression profiling spares early breast cancer patients from adjuvant therapy: derived and validated in two population-based cohorts. *Breast Cancer Research* **7**, 953-964.
- Shaughnessy, J. D., Zhan, F., Buring, B., Huang, Y., Colla, S., Hanamura, I., Stewart, J. P., Kordsmeier, B., Randolph, C., Williams, D. R., Xiao, Y., Xu, H., Epstein, J., Anaissie, E., Krishna, S. G., Cottler-Fox, M., Hollmig, K., Mohiuddin, A., Pineda-Roman, M., Tricot, G., van Rhee, F., Sawyer, J., Alsayed, Y., Walker, R., Zangari, M., Crowley, J. and Barlogie, B. (2007). A validated gene expression model of high-risk multiple myeloma is defined by deregulated expression of genes mapping to chromosome 1. *Blood* **109**, 2276-2284.
- Shi, L., Campbell, G., Jones, W. D., et al. (2010). The MAQC-II Project: A comprehensive study of common practices for the development and validation of microarray-based predictive models. *Nature Biotechnology* **28**, 827-838.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58**, 267-288.
- Tibshirani, R. (1997). The lasso method for variable selection in the Cox model. *Statist. Medicine* **16**, 385-395.
- Uno, H., Cai, T., Pencina, M. J., D'Agostino, R. B. and Wei, L. J. (2011). On the C-statistics for evaluating overall adequacy of risk prediction procedures with censored survival data. *Statist. Medicine* **30**, 1105-1117.
- Wang, S., Nan, B., Zhu, J. and Beer, D. G. (2008). Doubly penalized Buckley-James method for survival data with high-dimensional covariates. *Biometrics* **64**, 132-140.
- Ying, Z. (1993). A large sample study of rank estimation for censored regression data. *Ann. Statist.* **21**, 76-99.
- Zhan, F., Huang, Y., Colla, S., Stewart, J., Hanamura, I., Gupta, S., Epstein, J., Yaccoby, S., Sawyer, J., Burington, B., Anaissie, E., Hollmig, K., Pineda-Roman, M., Tricot, G., van Rhee, F., Walker, R., Zangari, M., Crowley, J., Barlogie, B. and Shaughnessy, J. D. (2006). The molecular classification of multiple myeloma. *Blood* **108**, 2020-2028.
- Zhao, P. and Yu, B. (2006). On model selection consistency of lasso. *J. Mach. Learn. Res.* **7**, 2541-2563.
- Zhu, L. P., Li, L., Li, R. and Zhu, L. X. (2011). Model-free feature screening for ultrahigh-dimensional data. *J. Amer. Statist. Assoc.* **106**, 1464-1475.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.* **101**, 1418-1429.
- Zou, H., Hastie, T. and Tibshirani, R. (2007). On the "degrees of freedom" of the lasso. *Ann. Statist.* **35**, 2173-2192.

Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA.

E-mail: yili@umich.edu

Department of Statistics and Biostatistics, Rutgers University, Piscataway, NJ 08854, USA.

E-mail: ldicker@stat.rutgers.edu

Department of Biostatistics and Epidemiology, University of Pennsylvania, Philadelphia, PA 19104, USA.

E-mail: sihai@mail.med.upenn.edu

(Received September 2011; accepted December 2012)