

Memorial Sloan-Kettering Cancer Center
Memorial Sloan-Kettering Cancer Center, Dept. of Epidemiology
& Biostatistics Working Paper Series

Year 2016

Paper 32

Variable Selection for Case-Cohort Studies
with Failure Time Outcome

Andy Ni* Jianwen Cai[†]

Donglin Zeng[‡]

*Department of Epidemiology and Biostatistics, memorial sloan kettering cancer center,
nia@mskcc.org

[†]Department of Biostatistics, University of North Carolina at Chapel Hill, cai@bios.unc.edu

[‡]Department of Biostatistics, University of North Carolina at Chapel Hill,
dzeng@email.unc.edu

This working paper is hosted by The Berkeley Electronic Press (bepress) and may not be commercially reproduced without the permission of the copyright holder.

<http://biostats.bepress.com/mskccbiostat/paper32>

Copyright ©2016 by the authors.

Variable Selection for Case-Cohort Studies with Failure Time Outcome

Andy Ni, Jianwen Cai, and Donglin Zeng

Abstract

Case-cohort designs are widely used in large cohort studies to reduce the cost associated with covariate measurement. In many such studies the number of covariates is very large, so an efficient variable selection method is necessary. In this paper, we study the properties of variable selection using the smoothly clipped absolute deviation penalty in a case-cohort design with a diverging number of parameters. We establish the consistency and asymptotic normality of the maximum penalized pseudo-partial likelihood estimator, and show that the proposed variable selection procedure is consistent and has an asymptotic oracle property. Simulation studies compare the finite sample performance of the procedure with Akaike information criterion- and Bayesian information criterion-based tuning parameter selection methods. We make recommendations for use of the procedures in case-cohort studies, and apply them to the Busselton Health Study.

Variable Selection for Case-Cohort Studies with Failure Time Outcome

BY AI NI *, JIANWEN CAI, AND DONGLIN ZENG

3101 McGavran-Greenberg Hall, CB 7420, Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina 27599-7420, U.S.A.

aini01@gmail.com cai@bios.unc.edu dzeng@email.unc.edu

SUMMARY

Case-cohort designs are widely used in large cohort studies to reduce the cost associated with covariate measurement. In many such studies the number of covariates is very large, so an efficient variable selection method is necessary. In this paper, we study the properties of variable selection using the smoothly clipped absolute deviation penalty in a case-cohort design with a diverging number of parameters. We establish the consistency and asymptotic normality of the maximum penalized pseudo-partial likelihood estimator, and show that the proposed variable selection procedure is consistent and has an asymptotic oracle property. Simulation studies compare the finite sample performance of the procedure with Akaike information criterion- and Bayesian information criterion-based tuning parameter selection methods. We make recommendations for use of the procedures in case-cohort studies, and apply them to the Busselton Health Study.

Some key words: Case-cohort design; Diverging number of parameters; Oracle property; Smoothly clipped absolute deviation; Survival analysis; Variable selection.

1. INTRODUCTION

Large-scale epidemiological studies and disease prevention trials often follow thousands of subjects for a long period. The assembly of covariates for the entire study cohort can be prohibitively expensive, especially when it requires biological samples or expensive bioassays. Moreover, the rate of the event of interest is usually low in these studies, especially for such events as cardiovascular disease, stroke, or death. We refer to subjects who develop the event of interest during the study as cases and the others as noncases. If the covariates were to be measured for everyone in the study, most of the cost would be spent on noncases, who do not contribute as much information as cases. To reduce the cost and effort in collecting expensive covariates without losing much efficiency, Prentice (1986) proposed the case-cohort design, where the complete covariate information is only obtained from a random subcohort of the sample, plus all cases.

Various estimation methods have been developed for case-cohort studies under the proportional hazard model (Cox, 1972). Prentice (1986) and Self & Prentice (1988) proposed a pseudo-partial likelihood method that modifies the risk set to account for subcohort sampling. Barlow (1994) introduced a time-dependent weight to estimate the risk set from the subcohort sample and developed a robust variance estimator for the regression parameters. Kalbfleisch & Lawless

* Ai Ni is currently at Memorial Sloan Kettering Cancer Center

(1988) proposed a more efficient weight that uses the complete covariate history of all cases. Bor-
 gan et al. (2000) further studied several types of weights under the stratified case-cohort design.
 40 Kulich & Lin (2004) established the asymptotic properties of the efficiently weighted estima-
 tor (Kalbfleisch & Lawless, 1988). Kang & Cai (2009) extended this estimator to studies with
 multivariate failure time outcomes, and Kim et al. (2013) further improved its efficiency in the
 presence of multivariate failure time outcomes. In this paper, we focus on the efficient weighting
 45 proposed by Kalbfleisch & Lawless (1988) in a univariate unstratified case-cohort design.

In large epidemiological studies that use the case-cohort design, many covariates are usually
 collected, and one research goal is often to identify a subset related to the event of interest.
 With the inclusion of interactions and polynomial terms, the number of candidate covariates
 can be very large. As Huber (1973) argued, in the context of variable selection the number of
 50 parameters should be considered as increasing to infinity with sample size n . In this paper, we
 consider the scenario where the model size d_n diverges to infinity but at a slower rate than the
 sample size. Traditional variable selection methods such as stepwise and best subset selection
 are computationally intensive and unstable. Since the introduction of lasso by Tibshirani (1996),
 penalty-based variable selection procedures have achieved great success. Under certain regularity
 55 conditions, they can simultaneously select variables and estimate their coefficients. Many penalty
 functions have been proposed, among which the smoothly clipped absolute deviation (Fan &
 Li, 2001), adaptive lasso (Zou, 2006), adaptive elastic net (Zou & Zhang, 2009), and minimax
 concave (Zhang, 2010) penalties have been shown to possess the oracle property, namely, as
 $n \rightarrow \infty$, the procedure correctly identifies the true model with probability tending to one and
 60 estimates the standard errors of nonzero parameters as efficiently as if the true model is known.
 Fan & Li (2002) applied the smoothly clipped absolute deviation penalty to the proportional
 hazard model and proved its oracle property. Cai et al. (2005) further extended the penalized
 partial likelihood procedure to multivariate models with a diverging number of parameters, but
 to our knowledge, the properties of penalized variable selection have not been studied under the
 65 case-cohort design where not all covariates are fully observed.

2. PSEUDO-PARTIAL LIKELIHOOD FOR CASE-COHORT DESIGN

Suppose there are n independent subjects in a cohort. Let $Z_i(t)$ be the $d_n \times 1$, possibly time-
 dependent, covariate vector for subject i at time t . Since d_n goes to infinity with n , all quantities
 that are functions of the covariates depend on n . For notational simplicity, however, we suppress
 70 the subscript n for them. Without loss of generality, we partition the real-valued true paramter
 vector β_{n0} as $(\beta_{n0,I}^T, \beta_{n0,II}^T)^T$, where $\beta_{n0,I}$ and $\beta_{n0,II}$ are the nonzero and zero components of
 β_{n0} , respectively. Denote by k_n the dimension of $\beta_{n0,I}$, which is also allowed to diverge with n
 and k_n/d_n converges to a constant $c \in [0, 1]$.

Let T and C be respectively the time to the outcome of interest and the censoring time. Let
 75 $X = \min(T, C)$ be the observed time and $\Delta = I(T \leq C)$ be the censoring indicator, where
 $I(\cdot)$ is an indicator function. We assume that T and C are independent, conditional on Z .
 Define for subject i the counting process $N_i(t) = I(X_i \leq t, \Delta_i = 1)$ and the at-risk process
 $Y_i(t) = I(X_i \geq t)$. Let $\lambda_i(t)$ denote the hazard function for subject i . Cox (1972) proposed the
 proportional hazard model where $\lambda_i\{t \mid Z_i(t)\} = \lambda_0(t) \exp\{\beta^T Z_i(t)\}$, in which $\lambda_0(t)$ is an un-
 80 specified baseline hazard function.

Under the case-cohort design, suppose we randomly select a subcohort of fixed size \tilde{n} from the
 full cohort. Let ξ_i denote the indicator for the i th subject being selected into the subcohort, and
 let $\alpha = \tilde{n}/n = \text{pr}(\xi_i = 1) \in (0, 1]$ denote the selection probability for the i th subject. Here we
 consider simple random sampling without replacement. Under this sampling scheme (ξ_1, \dots, ξ_n)

are correlated. The covariate histories are not observed for censored subjects outside the sub-cohort. If complete covariate histories are available for all the cases, one can use the following pseudo-partial likelihood to estimate the regression coefficients β (Kalbfleisch & Lawless, 1988):

$$\tilde{\ell}_n(\beta) = \sum_{i=1}^n \int_0^\tau \left[\beta^T Z_i(t) - \log \sum_{j=1}^n \rho_j(t) Y_j(t) \exp\{\beta^T Z_j(t)\} \right] dN_i(t), \quad (1)$$

where τ is the time at the end of study, and $\rho_i(t) = \Delta_i + (1 - \Delta_i)\xi_i\hat{\alpha}^{-1}(t)$, $\hat{\alpha}(t) = \sum_{i=1}^n (1 - \Delta_i)\xi_i Y_i(t) / \{\sum_{i=1}^n (1 - \Delta_i)Y_i(t)\}$ is a time-dependent estimator of the true sampling probability α . The corresponding pseudo-partial score equation is

$$\tilde{\ell}'_n(\beta) = \sum_{i=1}^n \int_0^\tau \left\{ Z_i(t) - \frac{\tilde{S}^{(1)}(\beta, t)}{\tilde{S}^{(0)}(\beta, t)} \right\} dN_i(t) = 0,$$

where $\tilde{S}^{(k)}(\beta, t) = n^{-1} \sum_{i=1}^n \rho_i(t) Y_i(t) Z_i(t)^{\otimes k} e^{\beta^T Z_i(t)}$ for $k = 0, 1, 2$. For a vector a , $a^{\otimes 0} = 1$, $a^{\otimes 1} = a$, and $a^{\otimes 2} = aa^T$.

3. VARIABLE SELECTION WITH A PENALIZED PSEUDO-PARTIAL LIKELIHOOD

3.1. Penalized Pseudo-Partial Likelihood

We define a penalized pseudo-partial likelihood as

$$\tilde{Q}_n(\beta) = \tilde{\ell}_n(\beta) - n \sum_{j=1}^{d_n} P_{\lambda_{nj}}(|\beta_j|), \quad (2)$$

where $P_{\lambda_{nj}}(|\beta_j|)$ is a nonnegative penalty function with $P_{\lambda_{nj}}(0) = 0$. The nonnegative tuning parameter λ_{nj} controls the model complexity. We use the smoothly clipped absolute deviation penalty (Fan & Li, 2001) with covariate-specific tuning parameters λ_{nj} , which allows different regression coefficients to have different penalty functions. The smoothly clipped absolute deviation penalty is

$$P_{\lambda_{nj}}(\theta) = \begin{cases} \lambda_{nj}\theta, & \theta \leq \lambda_{nj}, \\ -\frac{\theta^2 - 2a\lambda_{nj}\theta + \lambda_{nj}^2}{2(a-1)}, & \lambda_{nj} < \theta \leq a\lambda_{nj}, \\ \frac{(a+1)\lambda_{nj}^2}{2}, & \theta > a\lambda_{nj}, \end{cases}$$

for some $a > 2$ and $\theta > 0$. The first derivative of the penalty is

$$P'_{\lambda_{nj}}(\theta) = \lambda_{nj}I(\theta \leq \lambda_{nj}) + \frac{(a\lambda_{nj} - \theta)_+}{a-1}I(\theta > \lambda_{nj}).$$

3.2. Regularity Conditions

For each n , we define

$$S_n^{(k)}(\beta_n, t) = \frac{1}{n} \sum_{i=1}^n Y_i(t) Z_i(t)^{\otimes k} e^{\beta_n^T Z_i(t)}, \quad k = 0, 1, 2,$$

$$s_n^{(k)}(\beta_n, t) = E\{S_n^{(k)}(\beta_n, t)\}, \quad k = 0, 1, 2,$$

$$e_n(\beta_n, t) = s_n^{(1)}(\beta_n, t) / s_n^{(0)}(\beta_n, t),$$

$$V_n(\beta_n, t) = \frac{S_n^{(2)}(\beta_n, t)S_n^{(0)}(\beta_n, t) - S_n^{(1)}(\beta_n, t)^{\otimes 2}}{S_n^{(0)}(\beta_n, t)^2},$$

$$\tilde{V}_n(\beta_n, t) = \frac{\tilde{S}_n^{(2)}(\beta_n, t)\tilde{S}_n^{(0)}(\beta_n, t) - \tilde{S}_n^{(1)}(\beta_n, t)^{\otimes 2}}{\tilde{S}_n^{(0)}(\beta_n, t)^2},$$

$$I_n(\beta_n) = E \left\{ \int_0^\tau V_n(\beta_n, t) S_n^{(0)}(\beta_n, t) d\Lambda_0(t) \right\},$$

$$\Gamma_n(\beta_n) = \text{var} \{ n^{-1/2} \tilde{\ell}'_n(\beta_n) \}.$$

105 We require the following regularity conditions:

Condition 1. $\int_0^\tau \lambda_0(t) dt < \infty$ and $E\{Y(\tau)\} > 0$;

Condition 2. $|Z_{ij}(0)| + \int_0^\tau |dZ_{ij}(t)| < C_1 < \infty$ almost surely for some constant C_1 , $i = 1, \dots, n$, and $j = 1, \dots, d_n$;

110 *Condition 3.* there exists a neighborhood \mathcal{B}_n of β_{n0} such that for all $\beta_n \in \mathcal{B}_n$ and $t \in [0, \tau]$, $\partial s_n^{(0)}(\beta_n, t) / \partial \beta_n = s_n^{(1)}(\beta_n, t)$, and $\partial^2 s_n^{(0)}(\beta_n, t) / \partial \beta_n \partial \beta_n^T = s_n^{(2)}(\beta_n, t)$. The functions $s_n^{(k)}(\beta_n, t)$ ($k = 0, 1, 2$) are continuous and bounded, and $s_n^{(0)}(\beta_n, t)$ is bounded away from zero on $\mathcal{B}_n \times [0, \tau]$;

Condition 4. there exist positive constants C_2, C_3, C_4 , and C_5 such that

$$0 < C_2 < \lambda_{\min}\{I_n(\beta_{n0})\} \leq \lambda_{\max}\{I_n(\beta_{n0})\} < C_3 < \infty,$$

$$0 < C_4 < \lambda_{\min}\{\Gamma_n(\beta_{n0})\} \leq \lambda_{\max}\{\Gamma_n(\beta_{n0})\} < C_5 < \infty,$$

where $\lambda_{\min}(\cdot)$ and $\lambda_{\max}(\cdot)$ are the minimum and maximum eigenvalues of a matrix;

115 *Condition 5.* $\min_{1 \leq j \leq k_n} |\beta_{nj0}| / \lambda_{nj} \rightarrow \infty$ as $n \rightarrow \infty$; and

Condition 6. $\liminf_{n \rightarrow \infty} \liminf_{\theta \rightarrow 0^+} P'_{\lambda_{nj}}(\theta) / \lambda_{nj} > 0$ for $j = 1, \dots, d_n$.

Condition 1 ensures a finite baseline cumulative hazard and a non-empty risk set at the end of the study. Condition 2 requires the stochastic process of each time-dependent covariate to have bounded variation almost surely. Condition 3 essentially requires $\exp\{\beta_n^T Z_i(t)\}$ to be integrable under a diverging dimension so that integration and differentiation with respect to $S_n^{(k)}(\beta_n, t)$ ($k = 0, 1$) can be interchanged. Condition 4 ensures that the covariance matrices of the score function under both regular and case-cohort designs are positive definite and have uniformly bounded eigenvalues for all n . It assumes a non-singular Hessian matrix of the objective function used for variable selection. The same condition has been assumed in the variable selection literature (Peng & Fan, 2004; Cai et al., 2005; Cho & Qu, 2013). Condition 5 specifies the rate at which the proposed procedure can distinguish nonzero parameters from zero ones. As $n \rightarrow \infty$, the size of nonzero parameters detectable by the procedure can approach zero, but at a slower rate than the tuning parameter. This condition is required for the development of the asymptotic properties of the proposed procedure, and has been assumed by many authors (Peng & Fan, 2004; Wang et al., 2009; Cho & Qu, 2013; Fan & Tang, 2013). In real-world biomedical research, there usually exists a fixed minimum clinically important effect size. Any effect smaller than this size can be effectively treated as zero. Thus, Condition 5 is a reasonable requirement. Condition 6 implies that those zero parameters, whose finite sample estimates are about the scale of λ_{nj} 's, will be automatically shrunk to zero. This helps to achieve the oracle property of variable selection.

3.3. Asymptotic Properties

135

Throughout this paper we use $O_p(\cdot)$ and $o_p(\cdot)$ to denote probability order relations and $O(\cdot)$ and $o(\cdot)$ to denote almost sure order relations. Let $a_n = \max_{1 \leq j \leq k_n} \{|P'_{\lambda_{nj}}(|\beta_{nj0}|)|\}$ and $b_n = \max_{1 \leq j \leq k_n} \{|P''_{\lambda_{nj}}(|\beta_{nj0}|)|\}$. We first prove the existence of a penalized pseudo-partial likelihood estimator that converges at rate $O_p\{d_n^{1/2}(n^{-1/2} + a_n)\}$, and then establish its oracle property. The proofs of Theorem 1 and 2 are provided in the Appendix.

140

THEOREM 1. *Under Conditions 1 to 5, if $b_n \rightarrow 0$ and $d_n^4/n \rightarrow 0$ as $n \rightarrow \infty$, then with probability tending to one there exists a local maximizer $\hat{\beta}_n$ of $\tilde{Q}_n(\beta_n) = \tilde{\ell}_n(\beta_n) - n \sum_{j=1}^{d_n} P_{\lambda_{nj}}(|\beta_{nj}|)$, such that $\|\hat{\beta}_n - \beta_{n0}\| = O_p\{d_n^{1/2}(n^{-1/2} + a_n)\}$.*

From Theorem 1 one can obtain a $(n/d_n)^{1/2}$ -consistent penalized pseudo-partial likelihood estimator, provided that $a_n = O(n^{-1/2})$, which is the case for the smoothly clipped absolute deviation penalty under Condition 5. This consistency rate is the same as that of the maximum likelihood estimator for the exponential family (Portnoy, 1988). For Theorem 2, we define

145

$$\Sigma_n = \text{diag}\{P''_{\lambda_{1n}}(|\beta_{n01}|), \dots, P''_{\lambda_{k_n n}}(|\beta_{n0k_n}|)\}, \tag{3}$$

$$B_n = \{P'_{\lambda_{1n}}(|\beta_{n01}|)\text{sgn}(\beta_{n01}), \dots, P'_{\lambda_{k_n n}}(|\beta_{n0k_n}|)\text{sgn}(\beta_{n0k_n})\}^T. \tag{4}$$

THEOREM 2. *Under Conditions 1 to 6, if $b_n \rightarrow 0$, $d_n^5/n \rightarrow 0$, $\lambda_{nj} \rightarrow 0$, $\lambda_{nj}(n/d_n)^{1/2} \rightarrow \infty$, and $a_n = O(n^{-1/2})$ as $n \rightarrow \infty$, the $(n/d_n)^{1/2}$ -consistent local maximizer $\hat{\beta}_n = (\hat{\beta}_{n,I}^T, \hat{\beta}_{n,II}^T)^T$ must satisfy that $\hat{\beta}_{n,II} = 0$ with probability tending to one and for any nonzero $k_n \times 1$ constant vector u_n with $\|u_n\| = 1$,*

$$n^{1/2}u_n^T \Gamma_{n11}^{-1/2} (I_{n11} + \Sigma_n) \{\hat{\beta}_{n,I} - \beta_{n0,I} + (I_{n11} + \Sigma_n)^{-1} B_n\} \rightarrow N(0, 1)$$

in distribution, where Σ_n and B_n are defined in (3) and (4) respectively, I_{n11} consists of the first $k_n \times k_n$ components of $I_n(\beta_{n0})$, and Γ_{n11} consists of the first $k_n \times k_n$ components of $\Gamma_n(\beta_{n0})$.

150

Due to the diverging dimension of $\beta_{n0,I}$, Theorem 2 establishes the asymptotic normality of some linear combination of standardized estimators. However, by choosing a particular u_n , it can give the asymptotic distribution for each individual estimator. Thus, it provides a theoretical basis for inference on individual coefficients. The matrix $I_n(\beta_{n0})$ can be consistently estimated by $\hat{I}_n(\hat{\beta}_n) = n^{-1} \sum_{i=1}^n \int_0^\tau \tilde{V}_n(\hat{\beta}_n, t) dN_i(t)$. The estimator of matrix $\Gamma_n(\beta_{n0})$ is given in the Supplementary material. For the smoothly clipped absolute deviation penalty, $a_n = 0$, $\Sigma_n = 0$, and $B_n = 0$ for large n under Condition 5. Therefore, the result of Theorem 2 reduces to

$$n^{1/2}u_n^T \Gamma_{n11}^{-1/2} I_{n11} (\hat{\beta}_{n,I} - \beta_{n0,I}) \rightarrow N(0, 1)$$

in distribution as $n \rightarrow \infty$. The conditions $d_n^4/n \rightarrow 0$ and $d_n^5/n \rightarrow 0$ in the above theorems describe the divergence rate of d_n relative to the sample size. They do not impose any one-to-one relationship between finite d_n and n .

4. CONSIDERATIONS IN PRACTICAL IMPLEMENTATION

155

4.1. Local Quadratic Approximation and Variance Estimation

Since the smoothly clipped absolute deviation penalty function is not differentiable at the origin, in practical implementation the Newton–Raphson algorithm cannot be directly applied to maximize (2). Instead, we use a modified Newton–Raphson algorithm with a local quadratic

160 approximation to the penalty function. The unpenalized pseudo-partial likelihood (1) can be seen as a special case of the penalized pseudo-partial likelihood (2) with $P_{\lambda_{nj}}(|\beta_{nj}|) = 0$ for all $j = 1, \dots, d_n$. Applying Theorem 1 with $\lambda_{nj} = 0$ for all $j = 1, \dots, d_n$, we know there exists a $(n/d_n)^{1/2}$ -consistent maximizer of (1). The concavity of (1) ensures that the maximizer is unique. We use this maximizer as the initial value $\beta_n^{(0)}$ for the modified Newton–Raphson algorithm. If $|\beta_{nj}^{(0)}|$ is less than a pre-specified small positive constant c_j , then we set $\hat{\beta}_{nj} = 0$. 165 Otherwise, the penalty function is locally approximated by a quadratic function, $P_{\lambda_{nj}}(|\beta_{nj}|) \approx P_{\lambda_{nj}}\{|\beta_{nj}^{(0)}|\} + P'_{\lambda_{nj}}\{|\beta_{nj}^{(0)}|\}\{2|\beta_{nj}^{(0)}|\}^{-1}[\beta_{nj}^2 - \{\beta_{nj}^{(0)}\}^2]$, which has the same value and first derivative as the original penalty at $\beta_{nj}^{(0)}$. It follows that $P'_{\lambda_{nj}}(|\beta_{nj}|) \approx [P'_{\lambda_{nj}}\{|\beta_{nj}^{(0)}|\}/|\beta_{nj}^{(0)}|]\beta_{nj}$. This approximation is local in the sense that it is only good in the neighborhood of $\beta_{nj}^{(0)}$. With the 170 approximated penalty function, one Newton–Raphson step is performed and the updated nonzero estimate is used as the new initial value. The process is iterated until convergence or until all parameters are estimated as zero. Hunter & Li (2005) showed that the local quadratic approximation is an extension of the expectation-maximization algorithm and has the same properties.

The sandwich estimate of the covariance matrix for $\hat{\beta}_n$ can be directly obtained from the 175 last iteration of the above algorithm as $\text{cov}(\hat{\beta}_n) = \{\tilde{\ell}''_n(\hat{\beta}_n) - n\Sigma_\lambda(\hat{\beta}_n)\}^{-1}n\hat{\Gamma}_n(\hat{\beta}_n)\{\tilde{\ell}''_n(\hat{\beta}_n) - n\Sigma_\lambda(\hat{\beta}_n)\}^{-1}$, where $\Sigma_\lambda(\beta_n) = \text{diag}\{P'_{\lambda_{1n}}\{|\beta_{n1}^{(0)}|\}/|\beta_{n1}^{(0)}|, \dots, P'_{\lambda_{d_nn}}\{|\beta_{nd_n}^{(0)}|\}/|\beta_{nd_n}^{(0)}|\}$. The sandwich estimate of the covariance matrix is only applicable to the nonzero parameter estimates.

4.2. Selection of Tuning Parameters

The tuning parameter λ in the smoothly clipped absolute deviation penalty function $P_\lambda(\cdot)$ controls the magnitude of the penalty on each regression coefficient and thereby controls the complexity of the selected model. Typical methods of selecting tuning parameters are data-driven procedures such as K-fold cross-validation and generalized cross-validation (Craven & Wahba, 1979). We follow Fan & Li (2002) and Cai et al. (2005) and use generalized cross-validation. The effective number of parameters measures the degrees of freedom in a regularized regression model (Hastie et al., 2009). For the proportional hazards model, the effective number of parameters is defined as $e(\lambda_{1n}, \dots, \lambda_{d_nn}) = \text{tr}\{[\tilde{\ell}''_n(\hat{\beta}_n) - n\Sigma_\lambda(\hat{\beta}_n)]^{-1}\tilde{\ell}''_n(\hat{\beta}_n)\}$ (Fan & Li, 2002). The generalized cross-validation statistic is defined as

$$\text{GCV}(\lambda_{1n}, \dots, \lambda_{d_nn}) = \frac{-\tilde{\ell}_n(\hat{\beta}_n)}{n\{1 - e(\lambda_{1n}, \dots, \lambda_{d_nn})/n\}^2},$$

180 which is guaranteed to be positive since the log-pseudo-partial likelihood in the numerator is negative. The optimal tuning parameters are chosen as $\text{argmin}_{(\lambda_{1n}, \dots, \lambda_{d_nn})} \text{GCV}(\lambda_{1n}, \dots, \lambda_{d_nn})$. This d_n -dimensional optimization problem is difficult to solve in practice. We follow Cai et al. (2005) and take $\lambda_{nj} = \lambda_n \hat{\text{se}}\{\beta_{nj}^{(0)}\}$, where $\hat{\text{se}}\{\beta_{nj}^{(0)}\}$ is the estimated standard error of the unpenalized pseudo-partial likelihood estimator used in Section 4.1. Then the optimization problem reduces to a one-dimensional search for the optimal λ_n .

185 When $e(\lambda_n)/n$ is small, as is the case under the conditions for Theorems 1 and 2, we can write $\log \text{GCV}(\lambda_n) = \log\{-\tilde{\ell}_n(\hat{\beta}_n)/n\} - 2 \log\{1 - e(\lambda_n)/n\} \approx \log\{-\tilde{\ell}_n(\hat{\beta}_n)/n\} + 2e(\lambda_n)/n$. This expression is analogous to the Akaike information criterion (Akaike, 1973), so we denote $\log \text{GCV}(\lambda_n)$ as $\text{AIC}(\lambda_n)$, and define $\lambda_n^{\text{AIC}} = \text{argmin}_{\lambda_n} \text{AIC}(\lambda_n)$. Following the idea of the Bayesian information criterion (Schwarz, 1978), we define another tuning parameter selection criteria, where the optimal tuning parameter, denoted by λ_n^{BIC} , minimizes 190 $\text{BIC}(\lambda_n) = \log\{-\tilde{\ell}_n(\hat{\beta}_n)/n\} + \log(n)e(\lambda_n)/n$. Wang et al. (2007) and Zhang et al. (2010)

showed in linear and generalized linear models with a finite number of parameters that λ_n^{AIC} overfits the model with a positive probability whereas λ_n^{BIC} consistently identifies the true model. Such a result has not been established in the Cox proportional hazards model to our best knowledge. In the simulation section that follows, we investigate the performance of λ_n^{AIC} and λ_n^{BIC} . Following Fan & Li (2001), we set the second tuning parameter a in the penalty function to 3.7 in our simulation.

In practice, researchers can perform a grid search to identify λ_n^{AIC} and λ_n^{BIC} . The lower limit of the search range is zero and the upper limit is the smallest λ_n that gives an empty model. From our simulation experience, the upper limit rarely exceeds 2. Moreover, the model selection results are fairly robust to the fineness of the search grid.

5. NUMERICAL STUDY AND APPLICATION

5.1. Simulation Study

Independent failure times are generated from the proportional hazards model. We set the baseline hazard $\lambda_0(t) = 2$ and the model dimension $d_n = \lceil 5n_c^{1/5-1/500} \rceil$ to reflect its dependence on sample size, where n_c is the expected number of cases for a given censoring rate and $\lceil x \rceil$ rounds x to the nearest integer. We relate the model dimension to the number of cases rather than the sample size directly because the former better represents the amount of information in the dataset. We follow Tibshirani (1997) and consider two scenarios for the true parameter: a few large effects and many small effects. In the first scenario, $\beta_{n0} = (0.35, 0, 0, 0.6, 0, 0, -0.8, 0, 0, 0.6, 0, 0, -0.8, 0, 0, \dots)$. Thus a third of the components of β_{n0} are nonzero and the smallest nonzero effect in absolute value is 0.35, which corresponds to a hazard ratio of 1.4. In the second scenario, all components of β_{n0} equal 0.1, which corresponds to a hazard ratio of 1.1. In both scenarios, we generate the design matrix Z as a mixture of correlated binary and continuous variables. First, a d_n -dimensional multivariate standard normal variable Z^* is generated with $\text{corr}(Z_i^*, Z_j^*) = 0.5^{|i-j|}$. Then the first three components of Z^* are kept continuous while the next three components are dichotomized at zero, and this pattern is repeated for the rest of Z^* . Thus, half of the covariates become binary with parameter 0.5. Censoring times C_i are generated from a uniform distribution $U(0, c)$, with c adjusted to achieve the desired censoring percentage.

Various sample sizes, censoring rates, and noncase-to-case ratios are considered for both scenarios. Performance of the penalized variable selection with tuning parameter λ_n^{AIC} and λ_n^{BIC} is assessed. As a benchmark, we include the hard threshold variable selection procedure, where the unpenalized full model is fit and the components of the unpenalized estimates with a significant Wald test at 0.05 level are included in the final model. We also include the oracle procedure where the correct subset of covariates is used to fit the model. As the censoring rate is typically high in case-cohort studies, we set it to 80% and 90%, with 1000 replications for each setting.

We define model error for a given model as $\text{ME}(\hat{\mu}) = E\{E(T | z) - \hat{\mu}(z)\}^2$. Under the proportional hazard model with constant baseline hazard λ_0 , $\text{ME}(\hat{\mu}) = \lambda_0^{-2} E\{\exp(-\hat{\beta}_n^T z) - \exp(-\beta_{n0}^T z)\}^2$. The relative model error of a given model is defined as the ratio of its model error to that of the unpenalized full model. We use the median and the median absolute deviation of the relative model error to evaluate the prediction performance of different procedures. We also calculate the average number of parameters correctly estimated as zero, the average number of parameters erroneously estimated as zero, and the overall rate of identifying the true model as measures of variable selection performance. Point estimates, empirical and model-based stan-

standard errors, and the empirical 95% confidence interval coverages are calculated for $\beta_{n01} = 0.35$ in the first scenario.

Table 1 summarizes the simulation results under the scenario of a few large effects. The penalized method with tuning parameter λ_n^{BIC} has by far the best performance in all settings in terms of the relative model error and the rate of identifying the true model. The inferior performance of λ_n^{AIC} is apparently due to overfitting as shown by the low average number of correctly identified zero parameters; this is consistent with the theoretical findings of Wang et al. (2007) and Zhang et al. (2010). For both λ_n^{AIC} and λ_n^{BIC} , more noncases in the case-cohort and lower censoring rate are associated with better prediction and variable selection performance. Table 2 summarizes the parameter estimation of $\beta_{n01} = 0.35$ under the same settings as Table 1, but only using simulation replications where β_{n01} is correctly identified as nonzero. Conditional on $\hat{\beta}_{n1} \neq 0$, all procedures produce approximately unbiased point and standard error estimates and the coverage is close to the nominal level. The normality of the sampling distributions of $\hat{\beta}_{n1}$ was assessed by Q-Q plots; see the Supplementary Material. The sampling distribution of $\hat{\beta}_{n1}$ is a mixture of a point mass at zero and a left-truncated distribution that is well approximated by a truncated normal distribution. As the rate of identifying the true model increases, the point mass at zero vanishes and the sampling distribution of $\hat{\beta}_{n1}$ becomes normal.

Table 3 summarizes the simulation results under the scenario of many small effects where all $\beta_{n0} = 0.1$. In this scenario the oracle model is just the unpenalized full model with the relative model error being unity by definition, which is not very informative and hence not included in the table. With many small but nonzero effects, none of the three methods can identify all the effects with a high probability, reflected by the near-zero rate of identifying the true model for all settings, which is not shown in the table. The inference results are not satisfactory either; they are not shown due to space limitations. Nevertheless, λ_n^{AIC} produces the smallest relative model error, suggesting that it has the best prediction performance among the three methods. Moreover, λ_n^{AIC} correctly identifies the largest number of small effects as nonzero. The Bayesian information criterion tends to select sparse models, so it may not perform as well as the Akaike information criterion when there are many small nonzero parameters. The relative model error is not comparable across different settings because it depends on the model error of the full model, which has large variation under this scenario.

5.2. Analysis of Busselton Health Study

We use the proposed variable selection procedures to analyze the Busselton Health Study data (Cullen, 1972; Knuiman et al., 2003). The study is a series of cross-sectional health surveys conducted in the town of Busselton in Western Australia. Every 3 years from 1966 to 1981, general health information for adult participants was collected by questionnaire and clinical visits. In this analysis we are interested in identifying risk factors for stroke. In particular, the main risk factor of interest is the serum ferritin level. We also consider several other risk factors in the variable selection process: age, body mass index, blood pressure treatment, systolic blood pressure, cholesterol, triglycerides, hemoglobin, and smoking status. All variables were measured at baseline. The full cohort of this analysis consists of 1401 subjects aged 40 to 89 years who participated in the Busselton Health Survey in 1981 and had no history of diagnosed coronary heart disease or stroke at that time. Subjects were followed until December 31, 1998, and their time to stroke, if one took place, was recorded. They were treated as censored if they left Western Australia during the follow-up period. There were 118 incidences of stroke in the full cohort during the follow-up period. To reduce costs and preserve stored serum, a case-cohort design was used where the serum ferritin level was only measured for a randomly selected subcohort plus all stroke cases. The random subcohort size was 450, and the case-cohort size was 513.

Table 1. Model selection performance with a few large effects

Method	Noncase : Case = 1:1					Noncase : Case = 2:1				
	α	RME median (MAD)	Zero Parm. C I		RITM (%)	α	RME median (MAD)	Zero Parm. C I		RITM (%)
$n = 3000, 80\% \text{ censored}, d_n = 18$										
HT	0.25	0.67 (0.21)	11.2	0.0	45.4	0.50	0.65 (0.21)	11.3	0.0	52.1
SCAD(AIC)		0.63 (0.20)	10.7	0.0	30.3		0.49 (0.22)	11.5	0.0	61.6
SCAD(BIC)		0.39 (0.20)	12.0	0.2	83.7		0.37 (0.18)	12.0	0.0	95.2
Oracle		0.34 (0.16)	12.0	0.0	100.0		0.36 (0.17)	12.0	0.0	100.0
$n = 3000, 90\% \text{ censored}, d_n = 15$										
HT	0.11	0.88 (0.30)	9.2	0.5	25.1	0.22	0.75 (0.29)	9.3	0.2	42.7
SCAD(AIC)		0.92 (0.14)	6.4	0.1	1.2		0.82 (0.20)	7.6	0.0	8.3
SCAD(BIC)		0.74 (0.38)	9.3	0.5	33.3		0.49 (0.30)	9.8	0.3	63.9
Oracle		0.32 (0.18)	10.0	0.0	100.0		0.33 (0.17)	10.0	0.0	100.0
$n = 6000, 90\% \text{ censored}, d_n = 18$										
HT	0.11	0.71 (0.24)	11.1	0.1	39.6	0.22	0.64 (0.21)	11.3	0.0	48.4
SCAD(AIC)		0.89 (0.12)	7.9	0.0	1.2		0.80 (0.16)	9.5	0.0	9.4
SCAD(BIC)		0.49 (0.24)	11.5	0.1	58.6		0.38 (0.18)	11.9	0.0	87.8
Oracle		0.36 (0.17)	12.0	0.0	100.0		0.33 (0.15)	12.0	0.0	100.0
$n = 10000, 90\% \text{ censored}, d_n = 20$										
HT	0.11	0.69 (0.20)	12.1	0.0	36.4	0.22	0.65 (0.20)	12.2	0.0	48.0
SCAD(AIC)		0.88 (0.14)	8.9	0.0	1.2		0.80 (0.18)	10.2	0.0	8.0
SCAD(BIC)		0.47 (0.21)	12.5	0.0	60.8		0.39 (0.18)	12.9	0.0	92.8
Oracle		0.34 (0.15)	13.0	0.0	100.0		0.35 (0.17)	13.0	0.0	100.0

α : subcohort sampling probability; RME: relative model error; MAD: median absolute deviation; C: average number of 0 parameters correctly identified as 0; I: average number of nonzero parameters incorrectly identified as 0; RITM: rate of identifying true model; HT: hard threshold; SCAD(AIC): smoothly clipped absolute deviation with λ_n^{AIC} ; SCAD(BIC): smoothly clipped absolute deviation with λ_n^{BIC} .

Table 5 summarizes the baseline characteristics of the full cohort and the subcohort. The average ferritin level is not available for the full cohort due to the case-cohort design. The summary statistics of the baseline characteristics are similar between the full cohort and sub-cohort, suggesting that the subcohort is representative of the full cohort. 285

We apply the hard threshold method and penalized variable selection with tuning parameter, λ_n^{AIC} and λ_n^{BIC} to the Busselton Health Study. In order to avoid missing any potentially important effects, we also include the quadratic terms of all continuous covariates as well as interactions between ferritin and all covariates in the initial model. The total number of parameters is 28. All continuous covariates are standardized using the means and standard deviations from the subcohort in Table 4. To decrease their skewness we log-transform ferritin and triglycerides values before standardization. The tuning parameter selector identifies $\lambda_n^{\text{AIC}} = 0.244$ and $\lambda_n^{\text{BIC}} = 0.305$. Table 5 shows the models identified by the three methods. Due to space limitations, only terms that are selected by at least one method are shown. The use of λ_n^{AIC} selects seven terms and λ_n^{BIC} selects four. Both methods select age, sex, blood pressure treatment, and squared systolic blood pressure as important risk factors for stroke. The use of λ_n^{AIC} additionally selects the linear term of systolic blood pressure, linear and squared terms of triglycerides. The hard threshold method only selects age and blood pressure treatment. 290

Table 2. Estimation performance for $\beta_{n01} = 0.35$ with a few large effects

Method	n_c	Noncase : Case = 1:1				Noncase : Case = 2:1				
		$\hat{\beta}_{n1}$	se_e ($\times 10^{-2}$)	se_m ($\times 10^{-2}$)	95% CI_e	n_c	$\hat{\beta}_{n1}$	se_e ($\times 10^{-2}$)	se_m ($\times 10^{-2}$)	95% CI_e
$n = 3000, 80\% \text{ censored}, d_n = 18$										
HT	998	0.36	7.00	6.66	92.6	1000	0.35	5.85	5.55	92.7
SCAD(AIC)	1000	0.35	6.68	5.95	92.0	1000	0.35	5.28	4.87	92.7
SCAD(BIC)	991	0.35	5.96	5.88	94.8	1000	0.35	5.12	4.84	93.3
Oracle	1000	0.35	6.06	5.89	94.5	1000	0.35	5.08	4.84	93.5
$n = 3000, 90\% \text{ censored}, d_n = 15$										
HT	888	0.40	10.9	11.0	92.8	971	0.37	9.26	9.20	94.4
SCAD(AIC)	981	0.38	11.9	10.2	89.8	997	0.36	9.24	8.29	92.2
SCAD(BIC)	916	0.38	10.3	9.83	92.5	964	0.36	8.19	8.04	94.7
Oracle	1000	0.36	10.8	9.87	92.1	1000	0.35	8.37	8.05	93.8
$n = 6000, 90\% \text{ censored}, d_n = 18$										
HT	992	0.37	8.27	7.95	92.5	1000	0.36	7.01	6.53	92.2
SCAD(AIC)	1000	0.36	8.40	7.32	91.2	1000	0.36	6.73	5.92	91.0
SCAD(BIC)	992	0.36	7.68	7.09	92.5	996	0.35	6.06	5.74	93.8
Oracle	1000	0.35	7.64	7.10	93.0	1000	0.35	6.03	5.74	94.0
$n = 10000, 90\% \text{ censored}, d_n = 20$										
HT	1000	0.36	6.51	6.29	93.2	1000	0.35	5.27	5.10	94.4
SCAD(AIC)	1000	0.36	6.31	5.83	91.6	1000	0.35	5.11	4.63	94.0
SCAD(BIC)	1000	0.36	5.93	5.67	94.0	1000	0.35	4.55	4.50	94.8
Oracle	1000	0.36	5.74	5.67	95.0	1000	0.35	4.53	4.50	94.8

n_c : number of simulation replications where $\hat{\beta}_{n1} \neq 0$; se_e : empirical standard error; se_m : model-based standard error; 95% CI_e : empirical 95% confidence interval coverage; HT: hard threshold; SCAD (AIC): smoothly clipped absolute deviation with λ_n^{AIC} ; SCAD (BIC): smoothly clipped absolute deviation with λ_n^{BIC} . Results in this table are based on replications where $\hat{\beta}_{n1} \neq 0$.

300 To shed some light on which model provides the best fit to the data, we performed five-fold cross-validation. The average log-pseudo-partial likelihood from the test datasets is used as the validation statistic. The hard threshold method and penalized variable selection with λ_n^{AIC} or λ_n^{BIC} give validation statistics of -621.5 , -627.7 , and -614.0 , respectively. Therefore, we consider the model with λ_n^{BIC} as the best fit to the Busselton data. According to this model, increased
 305 age, maleness, blood pressure treatment, and increased systolic blood pressure are associated with higher risk of stroke. There is no evidence that serum ferritin level is associated with stroke.

6. DISCUSSION

One potential limitation of the theorems presented in this study is that they only establish the consistency and oracle property for a local maximizer of the penalized objective function.
 310 Due to its non-concavity, there may be multiple maximizers for the penalized objective function. However, based on Section 3.5 of Fan & Li (2001) and the small bias in the estimates in Table 2, it is reasonable to assume that the maximizer identified using the unpenalized estimator as the initial value is the $(n/d_n)^{1/2}$ -consistent local maximizer described in Theorems 1 and 2.

In this paper the quantity $\hat{\alpha}(t)$ used in the weight function $\rho(t)$ is calculated at each failure time
 315 point, and so is time-dependent. When cases are rare, $\hat{\alpha}(t)$ is almost constant across t . However,

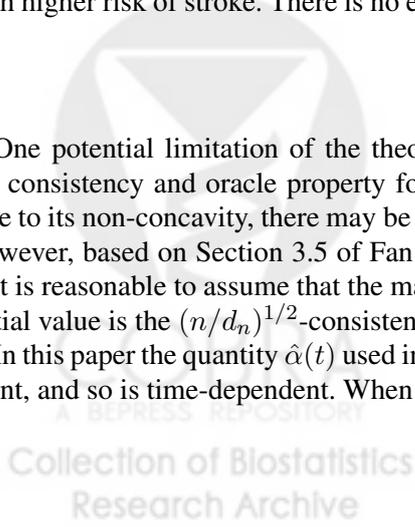


Table 3. Model selection performance with many small effects (all $\beta_{n0} = 0.1$)

Method	Noncase : Case = 1:1			Noncase : Case = 2:1		
	α	RME median (MAD)	Nonzero estimates	α	RME median (MAD)	Nonzero estimates
<i>n</i> = 3000, 80% censored, <i>d_n</i> = 18						
HT	0.25	2.90 (1.50)	4.0	0.50	3.59 (1.82)	5.2
SCAD(AIC)		1.79 (0.88)	6.0		3.15 (1.59)	5.5
SCAD(BIC)		5.62 (2.39)	1.3		8.94 (3.46)	1.1
<i>n</i> = 3000, 90% censored, <i>d_n</i> = 15						
HT	0.11	1.89 (1.00)	2.6	0.22	2.91 (1.63)	3.5
SCAD(AIC)		0.99 (0.29)	6.0		1.67 (0.78)	5.4
SCAD(BIC)		2.48 (1.23)	1.8		4.92 (2.08)	1.5
<i>n</i> = 6000, 90% censored, <i>d_n</i> = 18						
HT	0.11	2.82 (1.45)	3.4	0.22	3.48 (1.69)	4.5
SCAD(AIC)		1.08 (0.28)	8.6		1.41 (0.54)	8.3
SCAD(BIC)		3.17 (1.52)	3.0		5.36 (2.47)	2.6
<i>n</i> = 10000, 90% censored, <i>d_n</i> = 20						
HT	0.11	3.85 (2.02)	6.0	0.22	4.49 (2.37)	7.7
SCAD(AIC)		1.26 (0.39)	11.6		1.84 (0.81)	11.4
SCAD(BIC)		4.91 (2.49)	4.7		8.38 (3.75)	4.2

α : subcohort sampling probability; RME: relative model error; MAD: median absolute deviation; Nonzero estimates: average number of parameters not estimated as 0; HT: hard threshold; SCAD (AIC): smoothly clipped absolute deviation with λ_n^{AIC} ; SCAD (BIC): smoothly clipped absolute deviation with λ_n^{BIC} .

Table 4. Baseline characteristics of the Busselton Health Study

Variables	Full cohort (<i>n</i> =1401)	Subcohort (\tilde{n} =450)
	Mean (SD) or %	Mean (SD) or %
Age (yrs)	58.0 (10.8)	58.9 (10.9)
Body mass index	25.9 (3.9)	25.9 (4.0)
Blood pressure treatment (%)	17.2	18.4
Systolic blood pressure (mmHg)	132.2 (20.0)	132.9 (20.2)
Cholesterol (mmol/L)	6.14 (1.14)	6.24 (1.17)
Triglycerides (mmol/L)	1.52 (0.97)	1.55 (0.97)
Hemoglobin (g/100ml)	141.9 (12.0)	142.0 (11.5)
Smoking (%)		
Never	49.5	51.6
Former	32.4	32.0
Current	18.1	16.4
Ferritin (μ g/L)	–	148.1 (140.8)
log(ferritin)	–	4.57 (1.01)

using time-dependent $\hat{\alpha}(t)$ is more general and allows the sampling probability to vary with time *t*. Therefore, we use $\hat{\alpha}(t)$ in the paper. A potential practical issue is that $\hat{\alpha}(t)$ may not be reliable if the number of noncases in the random subcohort becomes very small, though this is highly unlikely due to the use of case-cohort design for studies of rare disease. In the unlikely situation where there is no noncase left in the subcohort, $\hat{\alpha}(t)$ is not well-defined. To avoid computational difficulties, one can define $(1 - \Delta)\xi/\hat{\alpha}(t) = 0$ if $\hat{\alpha}(t) = 0$. In fact, when $\hat{\alpha}(t) = 0$, $1 - \Delta$ is necessarily 0 for all subjects remaining in the subcohort.

320

Table 5. *Estimated coefficients and standard errors from Busselton Health Study data*

Variable	Hard threshold $\hat{\beta}$ (sê)	SCAD (AIC) $\hat{\beta}$ (sê)	SCAD (BIC) $\hat{\beta}$ (sê)
Age (yrs)	0.92 (0.27)	0.87 (0.15)	0.85 (0.14)
Sex (1=female)	0 (-)	-0.61 (0.26)	-0.65 (0.25)
Blood pressure treatment	0.83 (0.34)	0.83 (0.29)	0.89 (0.25)
Systolic blood pressure	0 (-)	0.21 (0.15)	0 (-)
Systolic blood pressure ²	0 (-)	0.092 (0.067)	0.16 (0.044)
log(triglycerides)	0 (-)	-0.24 (0.18)	0 (-)
log ² (triglycerides)	0 (-)	0.18 (0.093)	0 (-)

All continuous covariates were standardized using the means and standard deviations based on the random sub-cohort before the variable selection procedure.

SCAD (AIC): smoothly clipped absolute deviation with λ_n^{AIC} ; SCAD (BIC): smoothly clipped absolute deviation with λ_n^{BIC} .

There is a strong line of research on the convergence of and post-selection inference of penalized estimators (Leeb & Pötscher, 2005; Leeb & Pötscher, 2006; Pötscher & Leeb, 2009). In particular, Pötscher & Leeb (2009) showed that the penalized estimators are not uniformly consistent, and that their asymptotic distributions are non-normal if the true parameter lies within a shrinking neighborhood of zero with rate $(d_n/n)^{1/2}$. The lack of local regularity is a theoretical limitation of the penalized variable selection methods. However, in this paper Condition 5 together with the requirement that $\lambda_{nj}(n/d_n)^{1/2} \rightarrow \infty$ for all j ensures that the nonzero parameters are uniformly larger than $O\{(d_n/n)^{1/2}\}$, and therefore avoids the aforementioned irregularity. Our simulation study suggests that the performance of the proposed variable selection method depends on the true effect size. In practice, since these sizes are unknown, we suggest conducting penalized variable selection with both Akaike and Bayesian information criteria-based tuning parameter selection, and then using cross-validation to choose the best model, as done in Section 5.2. Theoretical justification of these model selection approaches will be further investigated. Moreover, the regularity conditions required for our asymptotic results may not be testable under finite samples. Therefore, it will be important to replicate findings from one particular finite data analysis. One possibility to examine the consistency of findings is to use bootstrap data or apply resampling-based variable selection approach such as stability selection (Meinshausen & Bühlmann, 2010).

In the Busselton data analysis we standardized all continuous covariates, for several reasons. First, this makes the regression coefficients comparable. Second, it reduces the correlation between the linear and quadratic terms and between the main effect and interaction terms, which generally results in more robust and precise parameter estimates. More importantly, penalized regression procedures are not invariant to covariate scaling, and standardization makes the penalization fair for all covariates (Tibshirani, 1997). For these reasons, we recommend standardizing continuous covariates before carrying out penalized regression.

ACKNOWLEDGEMENT

We thank Professor Matthew Knuiman and the Busselton Population Medical Research Foundation for permission to use the data in the illustration. This work was partially supported by grants from the U.S. National Institutes of Health.

SUPPLEMENTARY MATERIAL

Supplementary material available at *Biometrika* online includes the proofs of Lemmas, the estimation of the covariance matrix $\Gamma_n(\beta_{n0})$, and the Q-Q plots of the estimate $\hat{\beta}_{n1}$ under the simulation scenario of a few large effects.

355

APPENDIX

Proof of Theorems

Throughout the proofs, we write $\tilde{\ell}'_n(\beta_{n0})_j = \partial \tilde{\ell}_n(\beta_{n0}) / \partial \beta_{nj}$, $\tilde{\ell}''_n(\beta_{n0})_{jk} = \partial^2 \tilde{\ell}_n(\beta_{n0}) / \partial \beta_{nj} \partial \beta_{nk}$, and $\tilde{\ell}'''_n(\beta_{n0})_{jkl} = \partial^3 \tilde{\ell}_n(\beta_{n0}) / \partial \beta_{nj} \partial \beta_{nk} \partial \beta_{nl}$. We also let $\tilde{V}_{nj}(s, t)$, $V_{nj}(s, t)$, $\tilde{S}_{nj}^{(2)}(\beta_{n0}, t)$, and $S_{nj}^{(2)}(\beta_{n0}, t)$ be the (j, k) th component of corresponding matrices. For a matrix $A = \{a_{ij}\}$, $(i, j = 1, \dots, n)$, the norm is defined as $\|A\| = (\sum_{i=1}^n \sum_{j=1}^n a_{ij}^2)^{1/2}$. The following lemma will be used repeatedly.

360

LEMMA 1. Let $W_n(t)$ and $G_n(t)$ be two sequences of processes with bounded variation almost surely, and $G_n(t)$ progressively measurable and cadlag. For some constant τ , assume that $\sup_{0 \leq t \leq \tau} \|W_n(t) - W(t)\| \rightarrow 0$ in probability for some bounded process $W(t)$, $W_n(t)$ is monotone on $[0, \tau]$, and $G_n(t)$ converges to a zero-mean process with continuous sample paths in the metric space $BV[0, \tau]$, the bounded variation function space in $[0, \tau]$. Then both $\sup_{0 \leq t \leq \tau} \left\| \int_0^t \{W_n(s) - W(s)\} dG_n(s) \right\|$ and $\sup_{0 \leq t \leq \tau} \left\| \int_0^t G_n(s) d\{W_n(s) - W(s)\} \right\|$ converge to zero in probability as $n \rightarrow \infty$.

365

The proof of this lemma follows straightforwardly from that of Lemma 1 in Lin (2000) by noticing that a process with bounded variation can be decomposed into two monotone processes.

370

We also need the following lemmas, whose proofs are provided in the Supplementary material.

LEMMA 2. Let $\xi = (\xi_1, \dots, \xi_n)$ be a random vector containing \tilde{n} ones and $n - \tilde{n}$ zeros, with each permutation equally likely. Let $X_{ni}(t)$ ($i = 1, \dots, n$) be a triangular array of real-valued random processes on $[0, \tau]$ with $E\{X_{ni}(t)\} = \mu_n(t)$, $\text{var}\{X_{ni}(0)\} < \infty$ and $\text{var}\{X_{ni}(\tau)\} < \infty$ for all i and n . Let $X_n(t) = \{X_{n1}(t), \dots, X_{nn}(t)\}$ and ξ be independent. Suppose that almost all paths of $X_{ni}(t)$ have finite variation. Then $n^{-1/2} \sum_{i=1}^n \xi_i \{X_{ni}(t) - \mu_n(t)\}$ converges weakly to a tight zero-mean Gaussian process and therefore $n^{-1} \sum_{i=1}^n \xi_i \{X_{ni}(t) - \mu_n(t)\}$ converges in probability to zero uniformly in t .

375

LEMMA 3. Given that ξ is independent of Δ and $Y(t)$, $n^{1/2} \{\hat{\alpha}^{-1}(t) - \alpha^{-1}\}$ converges weakly to a zero-mean Gaussian process.

LEMMA 4. Under Conditions 1 to 3, for any nonzero $d_n \times 1$ constant vector u_n with $\|u_n\| = C < \infty$ and $\|u_n\|_0 = c_n > 0$ where $\|\cdot\|_0$ denotes the number of nonzero components of a vector, $n^{1/2} \{\tilde{S}_n^{(0)}(\beta_{n0}, t) - S_n^{(0)}(\beta_{n0}, t)\}$, $(n/c_n)^{1/2} u_n^T \{\tilde{S}_n^{(1)}(\beta_{n0}, t) - S_n^{(1)}(\beta_{n0}, t)\}$, and $n^{1/2} c_n^{-1} u_n^T \{\tilde{S}_n^{(2)}(\beta_{n0}, t) - S_n^{(2)}(\beta_{n0}, t)\} u_n$ all converge weakly to tight zero-mean Gaussian processes.

380

LEMMA 5. Under Conditions 1 to 4, for any nonzero $d_n \times 1$ constant vector u_n with $\|u_n\| = 1$, $n^{-1/2} u_n^T \Gamma_n^{-1/2}(\beta_{n0}) \tilde{\ell}'_n(\beta_{n0})$ converges to a standard normal distribution, where $\Gamma_n(\beta_{n0})$ is the covariance matrix of $n^{-1/2} \tilde{\ell}'_n(\beta_{n0})$.

385

LEMMA 6. Under Conditions 1 to 4, $n^{-1/2} \{\tilde{\ell}''_n(\beta_{n0})_{jk} + n I_n(\beta_{n0})_{jk}\}$ is $O_p(1)$ for $j, k = 1, \dots, d_n$, where $I_n(\beta_{n0})_{jk}$ is the (j, k) th component of $I_n(\beta_{n0})$ as defined in Section 3.2.

LEMMA 7. Under Conditions 1 to 6, if $d_n^4/n \rightarrow 0$, $\lambda_{nj} \rightarrow 0$, and $\lambda_{nj} n^{1/2} d_n^{-1/2} \rightarrow \infty$, with probability tending to one, for any given $\beta_{n,I}$ satisfying $\|\beta_{n,I} - \beta_{n0,I}\| = O(d_n^{1/2} n^{-1/2})$ and any constant C , we have $\tilde{Q}_n\{(\beta_{n,I}^T, 0^T)^T\} = \max_{\|\beta_{n,I}\| \leq C d_n^{1/2} n^{-1/2}} \tilde{Q}_n\{(\beta_{n,I}^T, \beta_{n,I}^T)^T\}$.

390

Proof of Theorem 1. Let β_{n0} be the true parameters, and $\alpha_n = d_n^{1/2}(n^{-1/2} + a_n)$. It suffices to show that, for any $\varepsilon > 0$ and any constant vector u_n with $\|u_n\| = C$, there exists a large enough C such that $\text{pr}\{\sup_{\|u_n\|=C} \tilde{Q}_n(\beta_{n0} + \alpha_n u_n) < \tilde{Q}_n(\beta_{n0})\} \geq 1 - \varepsilon$. This implies that there exists a local maximizer $\hat{\beta}_n$ such that $\|\hat{\beta}_n - \beta_{n0}\| = O_p(\alpha_n)$. Since $P_{\lambda_{nj}}(0) = 0$ and $P_{\lambda_{nj}}(\cdot) \geq 0$, we have

$$\begin{aligned} & \tilde{Q}_n(\beta_{n0} + \alpha_n u_n) - \tilde{Q}_n(\beta_{n0}) \\ & \leq \{\tilde{\ell}_n(\beta_{n0} + \alpha_n u_n) - \tilde{\ell}_n(\beta_{n0})\} - n \sum_{j=1}^{k_n} \{P_{\lambda_{nj}}(|\beta_{n0j} + \alpha_n u_{nj}|) - P_{\lambda_{nj}}(|\beta_{n0j}|)\} = I_1 - I_2. \end{aligned}$$

We first consider I_1 . By Taylor expansion we have

$$I_1 = \alpha_n u_n^T \tilde{\ell}'_n(\beta_{n0}) + \frac{1}{2} \alpha_n^2 u_n^T \tilde{\ell}''_n(\beta_{n0}) u_n + \frac{1}{6} \alpha_n^3 \sum_{i=1}^n \sum_{j,k,l=1}^{d_n} \tilde{\ell}'''_i(\beta_{n0}^*)_{jkl} u_{nj} u_{nk} u_{nl} = I_{11} + I_{12} + I_{13},$$

where β_{n0}^* lies between β_{n0} and $\beta_{n0} + \alpha_n u_n$. From Lemma A5 we have $\tilde{\ell}'_n(\beta_{n0})_j = O_p(n^{1/2})$ for $j = 1, \dots, d_n$. Therefore,

$$|I_{11}| = |\alpha_n u_n^T \tilde{\ell}'_n(\beta_{n0})| \leq \alpha_n \|u_n\| \|\tilde{\ell}'_n(\beta_{n0})\| = \alpha_n \|u_n\| O_p\{(d_n n)^{1/2}\} = \|u_n\| O_p(\alpha_n^2 n).$$

The term I_{12} can be written as $\alpha_n^2 u_n^T \{\tilde{\ell}''_n(\beta_{n0}) + n I_n(\beta_{n0})\} u_n / 2 - \alpha_n^2 u_n^T n I_n(\beta_{n0}) u_n / 2 = J_1 - J_2$. By the Cauchy–Schwarz inequality and $\tilde{\ell}''_n(\beta_{n0})_{jk} + n I_n(\beta_{n0})_{jk} = O_p(n^{1/2})$ for $j, k = 1, \dots, d_n$, and Lemma A6, we have $|J_1| \leq \alpha_n^2 \|u_n\|^2 \|\tilde{\ell}''_n(\beta_{n0}) + n I_n(\beta_{n0})\| / 2 = \|u_n\|^2 O_p(\alpha_n^2 n^{1/2} d_n) = \|u_n\|^2 o_p(\alpha_n^2 n)$. By spectral decomposition of $I_n(\beta_{n0})$ and Condition 4, $|J_2| \geq \alpha_n^2 \|u_n\|^2 n \lambda_{\min}\{I_n(\beta_{n0})\} / 2 \geq \|u_n\|^2 (\alpha_n^2 n) C_2 / 2$. Under Conditions 1 to 3, $\partial \tilde{V}_{nj k}(\beta_n^*, t) / \partial \beta_{nl}$ has bounded variation in t for $i = 1, \dots, n, j, k, l = 1, \dots, d_n$. Therefore $\tilde{\ell}'''_i(\beta_n^*)_{jkl} = -\int_0^T \partial \tilde{V}_{nj k}(\beta_n^*, t) / \partial \beta_{nl} dN_i(t)$ is $O_p(1)$. Along with $\alpha_n = d_n^{1/2}(n^{-1/2} + a_n)$, $d_n^4/n \rightarrow 0$ and $d_n^2 a_n \rightarrow 0$, we have $|I_{13}| = O_p(d_n^{3/2} n \alpha_n^3 \|u_n\|^3) = O_p\{d_n^2 (n^{-1/2} + a_n)\} n \alpha_n^2 \|u_n\|^3 = \|u_n\|^3 o_p(\alpha_n^2 n)$. Therefore, for large enough $\|u_n\|$, $|J_2|$ dominates $|I_{11}|$, $|J_1|$, and $|I_{13}|$.

We now consider I_2 . By Taylor expansion and the Cauchy–Schwarz inequality

$$\begin{aligned} |I_2| &= \left| n \sum_{j=1}^{k_n} P'_{\lambda_{nj}}(|\beta_{n0j}|) \text{sgn}(\beta_{n0j}) \alpha_n u_{nj} + \frac{1}{2} n \sum_{j=1}^{k_n} P''_{\lambda_{nj}}(|\beta_{n0j}|) \alpha_n^2 u_{nj}^2 \{1 + o(1)\} \right| \\ &\leq n \left| \sum_{j=1}^{k_n} P'_{\lambda_{nj}}(|\beta_{n0j}|) \alpha_n u_{nj} \right| + \frac{1}{2} n \left| \sum_{j=1}^{k_n} P''_{\lambda_{nj}}(|\beta_{n0j}|) \alpha_n^2 u_{nj}^2 \{1 + o(1)\} \right| \\ &\leq n \alpha_n a_n k_n^{1/2} \|u_n\| + \frac{1}{2} n \alpha_n^2 b_n \|u_n\|^2 \{1 + o(1)\} \\ &= \|u_n\| O_p(\alpha_n^2 n). \end{aligned}$$

The last equality holds because $a_n = O_p(\alpha_n d_n^{-1/2})$ and $b_n \rightarrow 0$ under Condition 5. Therefore, $|J_2|$ dominates $|I_2|$ for large enough C . Since J_2 is negative, it follows that for large enough C , $\tilde{Q}_n(\beta_{n0} + \alpha_n u_n) - \tilde{Q}_n(\beta_{n0})$ is negative with probability tending to one as $n \rightarrow \infty$. This completes the proof of Theorem 1.

Proof of Theorem 2. The assertion that $\hat{\beta}_{n,I} = 0$ with probability tending to one as $n \rightarrow \infty$ follows directly from Lemma A7. To prove the second assertion, we first show that

$$n^{1/2} u_n^T \Gamma_{n11}^{-1/2} [(I_{n11} + \Sigma_n)(\hat{\beta}_{n,I} - \beta_{n0,I}) \{1 + o_p(1)\} + B_n] = n^{-1/2} u_n^T \Gamma_{n11}^{-1/2} \tilde{\ell}'_{n1}(\beta_{n0}) + o_p(1), \tag{A1}$$

where $\tilde{\ell}'_{n1}(\beta_{n0})$ consists of the first k_n components of $\tilde{\ell}'_n(\beta_{n0})$. Since $\hat{\beta}_{n,I}$ is the maximum penalized pseudo-partial likelihood estimator, $\partial \tilde{Q}_n(\hat{\beta}_n) / \partial \beta_{n,I} = 0$. By Taylor expansion of $\partial \tilde{Q}_n(\hat{\beta}_n) / \partial \beta_{n,I}$ at

$\beta_{n0,I}$ and the fact that $\hat{\beta}_{n,II} - \beta_{n0,II} = 0$ with probability tending to one, we have

$$\begin{aligned} & \tilde{\ell}'_{n1}(\beta_{n0}) + \tilde{\ell}''_{n1}(\beta_{n0})(\hat{\beta}_{n,I} - \beta_{n0,I}) + (\hat{\beta}_{n,I} - \beta_{n0,I})^T \tilde{\ell}'''_{n1}(\beta_n^*)(\hat{\beta}_{n,I} - \beta_{n0,I})/2 \\ & - nB_n - n\Sigma_n^{**}(\hat{\beta}_{n,I} - \beta_{n0,I}) = 0 \end{aligned} \quad (\text{A2})$$

with probability tending to one, where $\tilde{\ell}''_{n1}(\beta_{n0})$ consists of the first $k_n \times k_n$ components of $\tilde{\ell}''_n(\beta_{n0})$, $\tilde{\ell}'''_{n1}(\beta_n^*)$ consists of the first $k_n \times k_n \times k_n$ components of $\tilde{\ell}'''_n(\beta_n^*)$, β_n^* lies between $\hat{\beta}_n$ and β_{n0} , $\Sigma_n^{**} = \Sigma_n(\beta_n^{**})$, β_n^{**} lies between $\hat{\beta}_n$ and β_{n0} . Rearranging (A2) we have

$$\begin{aligned} & \{\tilde{\ell}''_{n1}(\beta_{n0}) - n\Sigma_n^{**}\}(\hat{\beta}_{n,I} - \beta_{n0,I}) - nB_n \\ & = -\tilde{\ell}'_{n1}(\beta_{n0}) - \frac{1}{2}(\hat{\beta}_{n,I} - \beta_{n0,I})^T \tilde{\ell}'''_{n1}(\beta_n^*)(\hat{\beta}_{n,I} - \beta_{n0,I}). \end{aligned} \quad (\text{A3})$$

Denote $\nu_n = (\hat{\beta}_{n,I} - \beta_{n0,I})^T \tilde{\ell}'''_{n1}(\beta_n^*)(\hat{\beta}_{n,I} - \beta_{n0,I})$. Multiply both sides of (A3) by $n^{-1/2}u_n^T \Gamma_{n11}^{-1/2}$,

$$\begin{aligned} & n^{1/2}u_n^T \Gamma_{n11}^{-1/2} \left\{ \frac{1}{n} \tilde{\ell}''_{n1}(\beta_{n0}) - \Sigma_n^{**} \right\} (\hat{\beta}_{n,I} - \beta_{n0,I}) - n^{1/2}u_n^T \Gamma_{n11}^{-1/2} B_n \\ & = -n^{-1/2}u_n^T \Gamma_{n11}^{-1/2} \tilde{\ell}'_{n1}(\beta_{n0}) - n^{-1/2}u_n^T \Gamma_{n11}^{-1/2} \nu_n/2. \end{aligned} \quad (\text{A4})$$

By the Cauchy–Schwarz inequality, $\|\nu_n\| \leq \|\hat{\beta}_{n,I} - \beta_{n0,I}\|^2 \sum_{i=1}^n \{\sum_{j,k,l=1}^{k_n} \tilde{\ell}'''_{i1}(\beta_n^*)_{jkl}^2\}^{1/2}$. As shown in the proof of Theorem 1, $\tilde{\ell}'''_{i1}(\beta_n^*)_{jkl} = O_p(1)$, so $\|\nu_n\| = O_p\{(d_n/n)nk_n^{3/2}\} = O_p(d_n^{5/2})$. By spectral decomposition of $\Gamma_{n11}^{-1/2}$, $d_n^5/n \rightarrow 0$, and Condition 4,

$$\frac{1}{2}n^{-1/2}u_n^T \Gamma_{n11}^{-1/2} \nu_n \leq \frac{\|u_n\| \|\nu_n\|}{2} n^{-1/2} \lambda_{\max}(\Gamma_n^{-1/2}) = O_p(d_n^{5/2}n^{-1/2}) = o_p(1). \quad (\text{A5})$$

The inequality in (A5) holds by the Cauchy–Schwarz inequality and the Cauchy interlacing inequality of symmetric matrices. Moreover, $u_n^T \Gamma_{n11}^{-1/2} n^{-1} \tilde{\ell}''_{n1}(\beta_{n0})(\hat{\beta}_{n,I} - \beta_{n0,I}) = u_n^T \Gamma_{n11}^{-1/2} \{n^{-1} \tilde{\ell}''_{n1}(\beta_{n0}) + I_{n11}(\beta_{n0})\}(\hat{\beta}_{n,I} - \beta_{n0,I}) - u_n^T \Gamma_{n11}^{-1/2} I_{n11}(\beta_{n0})(\hat{\beta}_{n,I} - \beta_{n0,I}) = J_1 - J_2$. By the Cauchy–Schwarz inequality and Lemma A6, we have $|J_1| \leq \|u_n^T \Gamma_{n11}^{-1/2}\| \|n^{-1} \tilde{\ell}''_{n1}(\beta_{n0}) + I_{n11}(\beta_{n0})\| \|\hat{\beta}_{n,I} - \beta_{n0,I}\| = \|u_n^T \Gamma_{n11}^{-1/2}\| \|\hat{\beta}_{n,I} - \beta_{n0,I}\| O_p(d_n n^{-1/2})$. Also, we have $|J_2| \geq \|u_n^T \Gamma_{n11}^{-1/2}\| \|\hat{\beta}_{n,I} - \beta_{n0,I}\| \lambda_{\min}(I_{n11}) \geq \|u_n^T \Gamma_{n11}^{-1/2}\| \|\hat{\beta}_{n,I} - \beta_{n0,I}\| \lambda_{\min}(I_n)$. Then, by Condition 4 we have

$$\left| \frac{J_1}{J_2} \right| \leq \frac{\|u_n^T \Gamma_{n11}^{-1/2}\| \|\hat{\beta}_{n,I} - \beta_{n0,I}\| O_p(d_n n^{-1/2})}{\|u_n^T \Gamma_{n11}^{-1/2}\| \|\hat{\beta}_{n,I} - \beta_{n0,I}\| \lambda_{\min}(I_n)} = O_p(d_n n^{-1/2}) = o_p(1).$$

Therefore, $J_1 = o_p(J_2)$ and $u_n^T \Gamma_{n11}^{-1/2} n^{-1} \tilde{\ell}''_{n1}(\beta_{n0})(\hat{\beta}_{n,I} - \beta_{n0,I}) = -u_n^T \Gamma_{n11}^{-1/2} I_{n11}(\beta_{n0})(\hat{\beta}_{n,I} - \beta_{n0,I})\{1 + o_p(1)\}$. Since $\hat{\beta}_n$ converges to β_{n0} in probability, it follows that

$$\begin{aligned} & u_n^T \Gamma_{n11}^{-1/2} \left\{ \frac{1}{n} \tilde{\ell}''_{n1}(\beta_{n0}) - \Sigma_n^{**} \right\} (\hat{\beta}_{n,I} - \beta_{n0,I}) \\ & = -u_n^T \Gamma_{n11}^{-1/2} \{I_{n11}(\beta_{n0}) + \Sigma_n\} (\hat{\beta}_{n,I} - \beta_{n0,I})\{1 + o_p(1)\}. \end{aligned} \quad (\text{A6})$$

By (A4), (A5), and (A6), we have that (A1) holds. By Lemma A5, $n^{-1/2}u_n^T \Gamma_{n11}^{-1/2} \tilde{\ell}'_{n1}(\beta_{n0})$ converges to the standard normal distribution. Therefore, $n^{1/2}u_n^T \Gamma_{n11}^{-1/2} (I_{n11} + \Sigma_n)\{\hat{\beta}_{n,I} - \beta_{n0,I} + (I_{n11} + \Sigma_n)^{-1}B_n\} \rightarrow N(0, 1)$ in distribution. This proves the second assertion of Theorem 2.

BIBLIOGRAPHY

- AKAIKE, H. (1973). Maximum likelihood identification of Gaussian autoregressive moving average models. *Biometrika* **60**, 255–265.
- BARLOW, W. E. (1994). Robust variance estimation for the case-cohort design. *Biometrics* **50**, 1064–1072.

- BORGAN, O., LANGHOLZ, B., SAMUELSEN, S. O., GOLDSTEIN, L. & POGODA, J. (2000). Exposure stratified case-cohort designs. *Lifetime Data Analysis* **6**, 39–58.
- CAI, J., FAN, J., LI, R. & ZHOU, H. (2005). Variable selection for multivariate failure time data. *Biometrika* **92**, 303–316.
- 450 CHO, H. & QU, A. (2013). Model selection for correlated data with diverging number of parameters. *Statistica Sinica* **23**, 901–927.
- COX, D. R. (1972). Regression models and life-tables (with discussion). *Journal of the Royal Statistical Society. Series B* **34**.
- 455 CRAVEN, P. & WAHBA, G. (1979). Smoothing noisy data with spline functions: Estimating the correct degree of smoothing by the method of generalized cross-validation. *Numer. Math.* **31**, 377–403.
- CULLEN, K. J. (1972). Mass health examinations in the Busselton population, 1966 to 1970. *Australian Journal of Medicine* **2**, 714–718.
- FAN, J. & LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* **96**, 1348–1360.
- 460 FAN, J. & LI, R. (2002). Variable selection for Cox’s proportional hazards model and frailty model. *Annals of Statistics* **30**, 74–99.
- FAN, Y. & TANG, C. Y. (2013). Tuning parameter selection in high dimensional penalized likelihood. *Journal of the Royal Statistical Society. Series B* **75**, 531–552.
- 465 HASTIE, T., TIBSHIRANI, R. J. & FRIEDMAN, J. (2009). *The Elements of Statistical Learning*. Berlin: Springer, 2nd ed.
- HUBER, P. J. (1973). Robust regression: Asymptotics, conjectures, and Monte Carlo. *Annals of Statistics* **1**, 799–821.
- HUNTER, D. & LI, R. (2005). Variable selection using MM algorithms. *Annals of Statistics* **33**, 1617–1642.
- KALBFLEISCH, J. D. & LAWLESS, J. F. (1988). Likelihood analysis of multi-state models for disease incidence and mortality. *Statistics in Medicine* **7**, 149–160.
- 470 KANG, S. & CAI, J. (2009). Marginal hazards model for case-cohort studies with multiple disease outcomes. *Biometrika* **96**, 887–901.
- KIM, S., CAI, J. & LU, W. (2013). More efficient estimators for case-cohort studies. *Biometrika* **100**, 695–708.
- 475 KNUIMAN, M. W., DIVITINI, M. L., OLYNYK, J. K., CULLEN, D. J. & BARTHOLOMEW, H. C. (2003). Serum ferritin and cardiovascular disease: A 17-year follow-up study in Busselton, Western Australia. *American Journal of Epidemiology* **158**, 144–149.
- KULICH, M. & LIN, D. (2004). Improving the efficiency of relative-risk estimation in case-cohort studies. *Journal of the American Statistical Association* **99**, 832–844.
- LEEB, H. & PÖTSCHER, B. M. (2005). Model selection and inference: facts and fiction. *Econometric Theory* **21**, 21–59.
- 480 LEEB, H. & PÖTSCHER, B. M. (2006). Can one estimate the conditional distribution of post-model-selection estimators? *The Annals of Statistics* **34**, 2554–2591.
- LIN, D. (2000). On fitting Cox’s proportional hazards models to survey data. *Biometrika* **87**, 37–47.
- MEINSHAUSEN, N. & BÜHLMANN, P. (2010). Stability selection (with discussion). *Journal of the Royal Statistical Society, Series B* **72**, 417–473.
- 485 PENG, H. & FAN, J. (2004). Nonconcave penalized likelihood with a diverging number of parameters. *Annals of Statistics* **32**, 928–961.
- PORTNOY, S. (1988). Asymptotic behavior of likelihood methods for exponential families when the number of parameters tends to infinity. *Annals of Statistics* **16**, 356–366.
- 490 PÖTSCHER, B. M. & LEEB, H. (2009). On the distribution of penalized maximum likelihood estimators: The lasso, SCAD, and thresholding. *Journal of Multivariate Analysis* **100**, 2065–2082.
- PRENTICE, R. L. (1986). A case-cohort design for epidemiologic cohort studies and disease prevention trials. *Biometrika* **73**, 1–11.
- SCHWARZ, G. (1978). Estimating the dimension of a model. *Annals of Statistics* **6**, 461–464.
- 495 SELF, S. G. & PRENTICE, R. L. (1988). Asymptotic distribution theory and efficiency results for case-cohort studies. *Annals of Statistics* **16**, 64–81.
- TIBSHIRANI, R. J. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B* **58**, 267–288.
- TIBSHIRANI, R. J. (1997). The lasso method for variable selection in the Cox model. *Statistics in Medicine* **16**, 385–395.
- 500 WANG, H., LI, B. & LENG, C. (2009). Shrinkage tuning parameter selection with a diverging number of parameters. *Journal of the Royal Statistical Society Series B* **71**, 671–683.
- WANG, H., LI, R. & TSAI, C.-L. (2007). Tuning parameter selectors for the smoothly clipped absolute deviation method. *Biometrika* **94**, 553–568.
- 505 ZHANG, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *Annals of Statistics* **38**, 894–942.
- ZHANG, Y., LI, R. & TSAI, C.-L. (2010). Regularization parameter selections via generalized information criterion. *Journal of the American Statistical Association* **105**, 312–323.

- ZOU, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* **101**, 1418–1429.
- ZOU, H. & ZHANG, H. H. (2009). On the adaptive elastic-net with a diverging number of parameters. *Annals of Statistics* **37**, 1733–1751.

510

