

*University of Michigan School of Public
Health*

The University of Michigan Department of Biostatistics Working
Paper Series

Year 2004

Paper 37

Model checking techniques for regression
models in cancer screening

Debashis Ghosh*

*University of Michigan, debashis.ghosh@ucdenver.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/umichbiostat/paper37>

Copyright ©2004 by the author.

Model checking techniques for regression models in cancer screening

Debashis Ghosh

Abstract

There has been much work on developing statistical procedures for associating tumor size with the probability of detecting a metastasis. Recently, Ghosh (2004) developed a unified statistical framework in which equivalences with censored data structures and models for tumor size and metastasis were examined. Based on this framework, we consider model checking techniques for semiparametric regression models in this paper. The procedures are for checking the additive hazards model. Goodness of fit methods are described for assessing functional form of covariates as well as the additive hazards assumption. The finite-sample properties of the methods are assessed using simulation studies.

Model checking techniques for regression models in cancer screening

Running headline: Regression residual diagnostics

Author: Debashis Ghosh

Author's affiliation:

Department of Biostatistics

University of Michigan

1420 Washington Heights

Ann Arbor, Michigan 48109-2029, USA



ABSTRACT

There has been much work on developing statistical procedures for associating tumor size with the probability of detecting a metastasis. Recently, Ghosh (2004) developed a unified statistical framework in which equivalences with censored data structures and models for tumor size and metastasis were examined. Based on this framework, we consider model checking techniques for semiparametric regression models in this paper. The procedures are for checking the additive hazards model. Goodness of fit methods are described for assessing functional form of covariates as well as the additive hazards assumption. The finite-sample properties of the methods are assessed using simulation studies.

Key words: Additive risk; Empirical process; Interval censoring; Regression diagnostic; Right censoring.



1. Introduction

Given the morbidity and mortality and associated costs of treating people with cancer, it is of interest to clinical researchers to determine optimal screening schedules for early detection of cancer. There has been much work done on developing mathematical models of screening (Yakovlev and Tsodikov, 1996, Ch. 5). In this area, the natural history of the disease has been broken down into multiple stages: the disease-free stage, the preclinical stage and the clinical stage. There has been much development of parametric statistical procedures using this multi-stage framework (Zelen and Feinleib, 1969; Albert, Gertman and Louis, 1978; Day and Walter, 1984; Shen and Zelen, 1999).

An alternative approach is to use statistical methods to better understand the relationship between various aspects of tumor biology and progression. One example of this is the development of procedures to associate the probability of detecting a metastatic tumor with tumor size. Solid cancers develop through a process in which tumors originate as a progenitor cell, which grows to a local lesion that shed cancer cells into the lymphatic system and/or blood stream (Foulds, 1969). Some of these cells are transported to distant organs and lead to the development of metastases. In most oncology settings, cancers where metastases have developed are more likely to be associated with worse clinical prognosis. There have been proposals for correlating size of tumor with probability of detecting a metastasis (Kimmel and Flehinger, 1991; Xu and Prorok, 1997, 1998). These authors developed nonparametric estimation procedures for the distribution of tumor size at which metastatic transitions occur based on data from a screening trial. A limitation of the methods proposed in the previous paragraph is that they do not allow for adjustment of covariates. In the cancer setting, covariates such as the tissue of origin of the tumor or age of the patient can affect the relationship between tumor size and probability of metastatic spread. Recently, we have proposed a general hypothesis testing and semiparametric regression framework for analyzing screening trial data in which tumor size is treated as a failure time variable (Ghosh, 2004).

An important issue in the fitting of any probabilistic model is assessing model adequacy.

While there has been goodness of fit testing developed for mechanistic models of carcinogenesis (Gregori et al., 2002), such methods are unavailable for the empirical methods described in the previous paragraph. In this paper, we will develop goodness of fit methods for assessing the adequacy of semiparametric regression models. The course of this paper is as follows. In Section 2, we review the results of Ghosh (2004) and describe regression estimation procedures for the additive hazards model under two sets of assumptions. Goodness of fit methods for assessing functional form of covariates and the additive hazards assumption are proposed in Section 3. In Section 4, we report the results of some simulation studies. Finally, we conclude with some brief remarks in Section 5.

2. Statistical Models for Tumor Size Progression

2.1 Notation and model assumptions

Let V denote the size of the tumor, \mathbf{Z} a p -dimensional vector of covariates and δ be an indicator of tumor metastasis (i.e., $\delta = 1$ if metastases are present, $\delta = 0$ otherwise). We observe the data $(V_i, \delta_i, \mathbf{Z}_i)$, $i = 1, \dots, n$, a random sample from (V, δ, \mathbf{Z}) . For the sake of completeness, we will now state the model assumptions utilized by Kimmel and Flehinger (1991):

1. Primary cancers grow monotonically, and metastases are irreversible.
2. The cancer samples are characterized by the primary tumor sizes at which metastatic transitions take place. We will denote Y as the random variable for this quantity. Let the cdf of Y be denoted by F^Y .
3. Let $\lambda_1(x)$ denote the hazard function for detecting a cancer with metastasis when the tumor size is x . Let $\lambda_0(x)$ denote the hazard function for detecting a cancer with no metastases when the tumor size is x . Assume that $\lambda_1(x) \geq \lambda_0(x)$.

In general, F^Y is nonidentifiable. If we assume that cancers are detected immediately when the metastasis occurs, then F^Y becomes identifiable. Alternatively, if we assume that detection of the cancer is not affected by the presence of metastases, then F^Y is identifiable. We refer to these two situations as Case I and Case II, respectively.

The effect of \mathbf{Z} on Y is formulated through the additive risk model:

$$\lambda(y|\mathbf{Z}) = \lambda_0(y) - \beta_0^T \mathbf{Z}, \tag{1}$$

where $\lambda(y|\mathbf{Z})$ is the hazard for Y conditional on covariates, $\lambda_0(\cdot)$ is an unspecified baseline hazard function and β_0 is a $p \times 1$ vector of unknown regression coefficients. In the case of rare events, β_0 in (1) has an interpretation as risk differences associated with a unit change in the covariates, adjusting for other variables. The use of this model has been argued by Breslow and Day (1980, pp. 53 – 57).

The major result demonstrated by Ghosh (2004) is the equivalence of (V, δ, \mathbf{Z}) with censored data structures from survival analysis. In the Case I situation, $V = Y \wedge C$ and $\delta = I(Y \leq C)$, where C is a random monitoring time, $I(A)$ is the indicator function for the set A , and $a \wedge b$ is the minimum of two numbers a and b . Thus, V can be treated as a right-censored version of Y . For the case II situation, $V = C$ and $\delta = I(Y \leq C)$. What this implies is that Y can be treated as interval-censored data subject to monitoring size V under the Case II assumptions. Based on this equivalence, Ghosh (2004) developed a comprehensive hypothesis testing and regression framework. Before describing the proposed goodness of fit procedures, we briefly outline the regression estimation procedures for the two situations.

2.2 Case I Estimation

Here, the data structure consists of $(V_i, \delta_i, \mathbf{Z}_i)$ ($i = 1, \dots, n$), a random sample from (V, δ, \mathbf{Z}) , where $V = Y \wedge C$ and $\delta = I(Y \leq C)$. Estimation in model (1) has been previously developed by Lin & Ying (1994). Define the following processes: $N_i(t) = I(V_i \leq t, \delta_i = 1)$

and $R_i(t) = I(V_i \geq t)$. The following estimating function can be used for estimation of β in (2.1):

$$\mathbf{U}(\beta) = \sum_{i=1}^n \int_0^\tau \{\mathbf{Z}_i - \bar{\mathbf{Z}}(t)\} \{dN_i(t) + R_i(t)\beta^T \mathbf{Z}_i dt\}, \quad (2)$$

where $\bar{\mathbf{Z}}(t) = \sum_{j=1}^n R_j(t)\mathbf{Z}_j / \sum_{j=1}^n R_j(t)$ and $\tau > 0$ is a constant chosen to satisfy certain technical conditions. In practice, we can take τ to be the largest size with an observed metastasis. Setting (2) equal to zero yields the following estimator for β_0 :

$$\hat{\beta} = - \left[\sum_{i=1}^n \int_0^\tau R_i(t) \{\mathbf{Z}_i - \bar{\mathbf{Z}}(t)\}^{\otimes 2} dt \right]^{-1} \left[\sum_{i=1}^n \int_0^\tau \{\mathbf{Z}_i - \bar{\mathbf{Z}}(t)\} dN_i(t) \right],$$

where $a^{\otimes 2} = aa^T$. As in Lin & Ying (1994), we can apply standard martingale arguments to show that the limiting distribution of $n^{1/2}(\hat{\beta} - \beta_0)$ is a p -dimensional normal random vector with mean zero and variance $\mathbf{A}^{-1}\mathbf{B}\mathbf{A}^{-1}$, where

$$\mathbf{A} = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n \int_0^\tau R_i(t) \{\mathbf{Z}_i - \bar{\mathbf{Z}}(t)\}^{\otimes 2} dt$$

and

$$\mathbf{B} = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n \int_0^\tau \{\mathbf{Z}_i(t) - \bar{\mathbf{Z}}(t)\}^{\otimes 2} dN_i(t).$$

Note that it is easy to consistently estimate \mathbf{A} and \mathbf{B} based on sample quantities. Denote these estimators as $\hat{\mathbf{A}}$ and $\hat{\mathbf{B}}$, respectively.

2.3 Case II Estimation

We now consider the situation where Y is treated as a random variable subject to interval censoring by V . We now describe the approach of Lin et al. (1998) for estimating β_0 in model (1) using the counting process $\tilde{N}_i(t) \equiv (1 - \delta_i)I(V_i \leq t)$ for the i th individual, $i = 1, \dots, n$. Note that $\tilde{N}_i(t)$ can potentially take one jump so that its corresponding hazard function is well-defined. Let $dH_i(t)$ denote the hazard and let $d\tilde{N}_i(t)$ be the increment corresponding to $\tilde{N}_i(t)$. For $d\tilde{N}_i(t)$ to equal one, V_i must equal t and the subject must be metastasis free before t . Denote the hazard of the former event be denoted by $d\Lambda^C(t)$. Under model (1), the second event has probability

$$\Pr(Y_i \geq t | \mathbf{Z}_i) = \exp \left[- \int_0^t \{\lambda_0(u) - \beta_0^T \mathbf{Z}_i\} du \right] = \exp \{-\Lambda_0(t) + \beta_0^T \mathbf{Z}_i^*(t)\},$$

where $\Lambda_0(t) = \int_0^t \lambda_0(u) du$ and $\mathbf{Z}_i^*(t) = t\mathbf{Z}$ ($i = 1, \dots, n$). Multiplying the probabilities for these two events yields

$$dH_i(t) = dH_0(t) \exp\{\beta_0^T \mathbf{Z}_i^*(t)\}, \quad (3)$$

where $dH_0(t) = \exp\{-\Lambda_0(t)\} d\Lambda^C(t)$. The model in (3) has a form identical to that of the proportional hazards model (Cox, 1972). Consequently, estimation of β_0 can be done based on the partial likelihood. The partial likelihood score function is given by

$$\tilde{\mathbf{U}}(\beta) = \sum_{i=1}^n \int_0^{\tau^*} \left\{ \mathbf{Z}_i^*(t) - \frac{\mathbf{S}^{(1)}(\beta, t)}{S^{(0)}(\beta, t)} \right\} d\tilde{N}_i(t), \quad (4)$$

where $\mathbf{S}^{(k)}(\beta, t) = n^{-1} \sum_{j=1}^n I(V_j \geq t) \mathbf{Z}_j^*(t)^{\otimes k} \exp\{\beta^T \mathbf{Z}_j^*(t)\}$, $k = 0, 1, 2$, $\mathbf{a}^{\otimes 0} = 1$, $\mathbf{a}^{\otimes 1} = \mathbf{a}$, $\mathbf{a}^{\otimes 2} = \mathbf{a}\mathbf{a}^T$. The constant $\tau^* > 0$ is a truncation time chosen to satisfy certain technical conditions. Let $\tilde{\beta}$ be the solution from setting (4) equal to $\mathbf{0}$. Using martingale theory, $n^{1/2}(\tilde{\beta} - \beta_0)$ converges in distribution to a normal random vector with mean zero and variance $\mathcal{I}(\beta_0) \equiv \lim_{n \rightarrow \infty} n^{-1} \mathbf{I}(\beta_0)$, where

$$\mathbf{I}(\beta) = \sum_{i=1}^n \int_0^{\tau^*} \left\{ \frac{\mathbf{S}^{(2)}(\beta, t)}{S^{(0)}(\beta, t)} - \frac{\mathbf{S}^{(1)}(\beta, t)^{\otimes 2}}{S^{(0)}(\beta, t)^2} \right\} d\tilde{N}_i(t).$$

This variance can be consistently estimated by $\tilde{\mathbf{I}} \equiv n^{-1} \mathbf{I}(\tilde{\beta})$.

3. Goodness of fit methods

We now develop numerical and graphical methods for model checking corresponding to the estimation procedures in Sections 2.1 and 2.2 by extending the simulation-based procedure of Lin, Wei and Ying (1993). We will primarily be interested in checking the functional form for the covariates \mathbf{Z} and assessing the additive hazards assumption. While such methods have been considered by Kim and Lee (1998) and Ghosh (2002), the methods are different. In particular, Kim and Lee (1998) considered two-sample goodness of fit tests to test time-invariant regression effects in the additive risk with right-censored data, which corresponds to the Case I situation. However, they did not consider methods for assessing functional form of covariates. Ghosh (2002) developed model checking techniques in the situation where the monitoring times themselves depend on covariates. For the case II scenario, the monitoring sizes do not depend on covariates.

3.1 Case 1 Methodology

Here we assume that metastases are detected immediately after they occur so that Y is a right-censored random variable. We first consider assessing the functional form of covariates.

Since

$$M_i(t) \equiv N_i(t) - \int_0^t R_i(u) \{d\Lambda_0(u) - \beta_0^T \mathbf{Z}_i du\} \quad (i = 1, \dots, n)$$

are martingales, goodness of fit methods for assessing functional form can be based on $\widehat{M}_i(t)$, where

$$\widehat{M}_i(t) \equiv N_i(t) - \int_0^t R_i(u) \{d\widehat{\Lambda}_0(u; \hat{\beta}) - \hat{\beta}^T \mathbf{Z}_i du\}$$

and $\widehat{\Lambda}_0(t; \beta) = \sum_{i=1}^n \int_0^t \{dN_i(u) + R_i(u) \beta^T \mathbf{Z}_i du\} / \sum_{j=1}^n R_j(u)$. Let Z_{ji} ($i = 1, \dots, n; j = 1, \dots, p$) denote the j th component of \mathbf{Z}_i and $\widehat{M}_i = \widehat{M}_i(\tau)$. The functional form of the j th component of \mathbf{Z} can be graphically assessed by plotting \widehat{M}_i versus Z_{ji} . If there are systematic deviations from 0, then this indicates misspecification of the functional form for the j th component of \mathbf{Z} . We can construct more formal tests for the functional form by considering cumulative sums of the \widehat{M}_i against values of the covariates. For $j = 1, \dots, p$, define $W_j(x)$ by

$$W_j(x) \equiv n^{-1/2} \sum_{i=1}^n I(Z_{ji} \leq x) \widehat{M}_i.$$

Let (G_1, \dots, G_n) be n i.i.d. realizations from a $N(0, 1)$ distribution,

$$l(t, x) = \sum_{i=1}^n \int_0^t \frac{R_i(u) I(Z_{ji} \leq x)}{\sum_{j=1}^n R_j(u)},$$

and

$$\mathbf{Q}(x) = n^{-1} \sum_{i=1}^n \int_0^\tau \{I(Z_{ji} \leq x) - l(u, x)\} R_i(u) \mathbf{Z}_i du.$$

Under the null hypothesis that the j th component of \mathbf{Z} is correctly specified in (2.1), we show in Appendix A.1. that the distribution of $W_j(x)$ can be approximated by the zero-mean Gaussian process

$$\begin{aligned} \widehat{W}_j(x) &= n^{-1/2} \sum_{i=1}^n \int_0^\tau \{I(Z_{ji} \leq x) - l(u, x)\} dN_i(u) G_i \\ &\quad - \mathbf{Q}^T(x) \widehat{\mathbf{A}}^{-1} n^{-1/2} \sum_{i=1}^n \int_0^\tau \{\mathbf{Z}_i - \bar{\mathbf{Z}}(t)\} dN_i(u) G_i. \end{aligned}$$

The distribution of $\widehat{W}_j(\cdot)$ can be easily simulated by repeatedly generating random samples $\{G_i\}$. One can plot $W_j(\cdot)$ for the observed data, along with a few realizations from $\widehat{W}_j(\cdot)$, to assess how unusual the observed residual pattern is. A more objective test can be constructed with the Kolmogorov-Smirnov type statistic $s_j \equiv \sup_x |W_j(x)|$; a p-value can be derived using the simulation method of Lin et al. (1993). By arguments of Appendix 3 of that paper, it can be shown that this test is consistent against departures from the null hypothesis of correct functional specification for the j th covariate.

To check of the assumption of additive hazards with respect to the j th covariate in (1), we consider the standardized score process $U_j^*(t)$:

$$U_j^*(t) = \widehat{\mathbf{B}}_{jj}^{-1/2} n^{-1/2} U_j(\tilde{\beta}, t),$$

where $U_j(\beta, t)$ is the j th component of

$$\mathbf{U}(\beta, t) = \sum_{i=1}^n \int_0^t \{\mathbf{Z}_i - \bar{\mathbf{Z}}(u)\} \{dN_i(u) + R_i(u)\beta^T \mathbf{Z}_i du\},$$

and $\widehat{\mathbf{B}}_{jj}^{-1}$ is the j th diagonal element of $\widehat{\mathbf{B}}^{-1}$, $j = 1, \dots, p$. Let

$$\mathbf{L}(\beta, t) = n^{-1} \sum_{i=1}^n \int_0^t R_i(u) \{\mathbf{Z}_i(u) - \bar{\mathbf{Z}}(u)\} d\widehat{\Lambda}_0(u; \beta).$$

($i = 1, \dots, n; j = 1, \dots, p$). The null distribution of $U_j^*(t)$ can be approximated by that of $\widehat{U}_j^*(t)$, where

$$\widehat{U}_j^*(t) = \widehat{\mathbf{B}}_{jj}^{-1/2} \left[n^{-1/2} \sum_{i=1}^n \int_0^t \{\mathbf{Z}_i - \bar{\mathbf{Z}}_j(u)\} dN_i(u) G_i - \mathbf{Q}^T(x) \widehat{\mathbf{A}}^{-1} n^{-1/2} \sum_{i=1}^n \int_0^\tau \{\mathbf{Z}_i - \bar{\mathbf{Z}}(t)\} dN_i(u) G_i \right].$$

and $\bar{\mathbf{Z}}_j(t)$ is the j th component of $\bar{\mathbf{Z}}_j(t)$ ($j = 1, \dots, p$). We prove this result in the Appendix.

As with $\widehat{W}_j(\cdot)$, the process $\widehat{U}_j^*(\cdot)$ is very easy to simulate from. Graphical assessments of the additive hazards assumption follows from a plot of $U_j^*(\cdot)$ and several realizations from $\widehat{U}_j^*(\cdot)$.

We can conduct a test of the null hypothesis of no violation of additive hazards for the j th covariate based on $h_j \equiv \sup_t |U_j^*(t)|$. A p-value for this test statistic can be computed by

simulation, similar to that for s_j . This test will be consistent against any deviations from the null hypothesis of additive hazards for the j th covariate. An overall test of additive hazards for model (1) can be based on $\sup_t \|\mathbf{U}(\hat{\beta}, t)\|$ or $\sup_t \sum_{j=1}^p |U_j^*(t)|$. Similar to h_j , these tests are consistent against any nonadditive hazards alternative.

3.2 Case 2 Methodology

Let us first consider the problem of assessing the functional form of the covariates in (1).

Since

$$X_i(t) \equiv \tilde{N}_i(t) - \int_0^t R_i(u) \exp\{\beta_0^T \mathbf{Z}_i^*(u)\} dH_0(u) \quad (i = 1, \dots, n)$$

are martingales, goodness of fit methods for assessing functional form can be based on

$$\tilde{X}_i(t) \equiv \tilde{N}_i(t) - \int_0^t R_i(u) \exp\{\tilde{\beta}^T \mathbf{Z}_i^*(u)\} d\tilde{H}_0(u; \tilde{\beta}),$$

and $\tilde{H}_0(t; \beta) = \sum_{i=1}^n \int_0^t d\tilde{N}_i(u) / S^{(0)}(\beta, u)$. The functional form of the j th component of \mathbf{Z} can be graphically assessed by plotting $\tilde{X}_i \equiv \tilde{X}_i(\tau^*)$ versus Z_{ji} . We can construct more formal tests for the functional form by considering cumulative sums of the \tilde{X}_i against values of the covariates. For $j = 1, \dots, p$, define $\tilde{W}_j(x)$ by

$$\tilde{W}_j(x) \equiv n^{-1/2} \sum_{i=1}^n I(Z_{ji} \leq x) \tilde{X}_i.$$

Define the following:

$$S(\beta, t, x) = n^{-1} \sum_{j=1}^n R_j(t) \exp\{\beta^T \mathbf{Z}_j^*(t)\} I(Z_{ji} \leq x),$$

and

$$\mathbf{B}(\beta, x) = n^{-1} \sum_{i=1}^n \int_0^{\tau^*} R_i(u) \exp\{\beta^T \mathbf{Z}_i^*(u)\} I(Z_{ji} \leq x) \left\{ \mathbf{Z}_i^*(u) - \frac{\mathbf{S}^{(1)}(\beta, u)}{S^{(0)}(\beta, u)} \right\} d\tilde{H}_0(u; \beta).$$

Under the null hypothesis that the j th component of \mathbf{Z} is correctly specified in (1), it can be shown that the distribution of $\tilde{W}_j(x)$ can be approximated by the zero-mean Gaussian process

$$\begin{aligned} \tilde{W}_j^*(x) &= n^{-1/2} \sum_{i=1}^n \int_0^{\tau^*} \left\{ I(Z_{ji} \leq x) - \frac{S(\tilde{\beta}, u, x)}{S^{(0)}(\tilde{\beta}, u)} \right\} d\tilde{N}_i(u) G_i \\ &\quad - \mathbf{B}^T(\tilde{\beta}, x) \tilde{\mathbf{I}}^{-1} n^{-1/2} \sum_{i=1}^n \int_0^{\tau} \left\{ \mathbf{Z}_i^*(u) - \frac{\mathbf{S}^{(1)}(\tilde{\beta}, u)}{S^{(0)}(\tilde{\beta}, u)} \right\} d\tilde{N}_i(u) G_i. \end{aligned}$$

A proof of this result can be found in Appendix A.2. One can plot $W_j(\cdot)$ for the observed data, along with a few realizations from $\tilde{W}_j^*(\cdot)$, to assess how unusual the observed residual pattern is. A more formal test can be constructed with $\tilde{s}_j \equiv \sup_x |\tilde{W}_j(x)|$. Using arguments similar to those in Appendix 3 of Lin et al. (1993), it can be shown that this test is consistent against departures from the null hypothesis of proper functional specification for the j th covariate. The p-value for this statistic can be approximated by $\Pr(\tilde{S}_j \geq \tilde{s}_j)$, where $\tilde{S}_j = \sup_x |\tilde{W}_j^*(x)|$; it is done using the simulation-based method of Lin et al. (1993).

To check of the assumption of additive hazards with respect to the j th covariate in (1), we consider the standardized score process $\tilde{U}_j(t)$:

$$\tilde{U}_j(t) = (\tilde{\mathbf{I}}_{jj}^{-1})^{1/2} n^{-1/2} \tilde{U}_j(\tilde{\beta}, t),$$

where $\tilde{U}_j(\beta, t)$ is the j th component of

$$\tilde{\mathbf{U}}(\beta, t) = \sum_{i=1}^n \int_0^t \left\{ \mathbf{Z}_i^*(u) - \frac{\mathbf{S}^{(1)}(\beta, u)}{S^{(0)}(\beta, u)} \right\} d\tilde{N}_i(u),$$

and $\tilde{\mathbf{I}}_{jj}^{-1}$ is the j th diagonal element of $\tilde{\mathbf{I}}^{-1}$, $j = 1, \dots, p$. Let

$$\mathbf{J}(\beta, t) = n^{-1} \sum_{i=1}^n \int_0^t R_i(u) e^{\beta^T \mathbf{Z}_i^*(u)} Z_{ji}^*(u) \left\{ \mathbf{Z}_i^*(u) - \mathbf{S}^{(1)}(\beta, u) / S^{(0)}(\beta, u) \right\} d\tilde{H}_0(u; \beta).$$

The null distribution of $U_j^*(t)$ can be approximated by that of $\tilde{U}_j^*(t)$, where

$$\begin{aligned} \tilde{U}_j^*(t) = & (\tilde{\mathbf{I}}_{jj}^{-1})^{1/2} \left[n^{-1/2} \sum_{i=1}^n \int_0^t \left\{ \mathbf{Z}_{ji}^*(u) - \frac{S_j^{(1)}(\tilde{\beta}, u)}{S^{(0)}(\tilde{\beta}, u)} \right\} d\tilde{N}_i(u) G_i \right. \\ & \left. - \mathbf{J}^T(\tilde{\beta}, t) \tilde{\mathbf{I}}^{-1} n^{-1/2} \sum_{i=1}^n \int_0^{\tau^*} \left\{ \mathbf{Z}_i^*(u) - \frac{\mathbf{S}^{(1)}(\tilde{\beta}, u)}{S^{(0)}(\tilde{\beta}, u)} \right\} d\tilde{N}_i(u) G_i \right], \end{aligned}$$

and $S_j^{(1)}(\beta, t)$ is the j th component of $\mathbf{S}^{(1)}(\beta, t)$. Graphical assessment of $W_j(\cdot)$ is easy. We can conduct a test of the null hypothesis of no violation of additive hazards for the j th covariate based on $\tilde{h}_j \equiv \sup_t |\tilde{U}_j(t)|$. An overall test of additive hazards for model (1) can be based on $\sup_t \|\tilde{\mathbf{U}}(\tilde{\beta}, t)\|$ or $\sup_t \sum_{j=1}^p |\tilde{U}_j^*(t)|$. As with $\tilde{W}(x)$, p-values are straightforward to calculate. Consistency of these tests from departures against the null hypotheses follows from Lin et al. (1993).

4. Simulation Studies

To assess the small-sample properties of the proposed methods, extensive simulation studies were conducted. In the ones reported here, we consider only a single covariate Z . We considered sample sizes $n = 50, 100$ and 200 . For each simulation setting, 1000 samples were generated, and 1000 resamplings were used to calculate p -values for each sample. Censoring was generated using three scenarios: an independent uniform $(0,2)$, $(0,3)$ and $(0,5)$ random variable.

The methods for assessing functional form were considered first. For the simulation studies reported here, $n/5$ subjects were assigned to one of five dose groups. The covariate Z takes values $0, 1, 2, 3$ and 4 . Failure times were generated from (1) with $\lambda_0(t) \equiv 1.0$ and $\beta = 0.05$. The sizes of s_1 and \tilde{s}_1 were assessed by using Z in the estimation procedures. To examine the powers of these statistics, the covariate Z^* was used in estimation, where $Z^* = 1$ if $Z > 2$ and 0 otherwise. The results of these numerical studies are summarized in Tables 1 and 2.

Based on these results, we find that the proposed methods perform reasonably well, at least for moderately larger sample sizes. For smaller sample sizes, the statistics tend to be somewhat unstable in terms of achieving the proper level of significance. This behavior diminishes in larger sample sizes. The method has good power. The reason s_1 has better power than \tilde{s}_1 because there is inherently more information available in the Case I scenario relative to the Case II situation.

Next, the procedures for assessing additive hazards were studied. In the ones reported here, Z is a 0–1 treatment indicator. Failure times were generated from the following model:

$$\lambda(t|Z) = \lambda_0(t) - \beta(t)Z. \quad (5)$$

To assess the size of the h_1 and \tilde{h}_1 , we set $\lambda_0(t) \equiv 1.0$ and $\beta(t) \equiv 0.5$. Power of these statistics was examined by setting $\lambda_0(t) = 2t$ and $\beta(t) = -6t$. The results are given in Table 2. We find that the proposed test statistics yield reasonable sizes, although the tests tend

to be anticonservative in smaller samples. This behavior diminishes for larger sample sizes. The proposed methods have good power as well.

5. Discussion

In this article, we have proposed model checking procedures for semiparametric regression models in cancer screening procedures. The model is based on a framework proposed by Kimmel and Flehinger (1991). Ghosh (2004) shows that this framework corresponds to treating tumor size as either a right-censored or as an interval-censored random variable. Based on this result, he develops semiparametric regression modelling procedures. It is important to have methods to assess goodness of fit for the resulting model fits, which was the goal of the work here.

Gregori et al. (2002) develop goodness of fit procedures using the method of Hjort (1990). While the goodness of fit idea is similar in spirit to that proposed here, the model being fit is substantially different. The model of Gregori et al. (2002) arises from a mechanistic model for tumor progression (Tsodikov and Yakovlev, 1996). Our model involves different assumptions from theirs.

References

- Albert, A., Gertman, P. M., and Louis, T. A. (1978). Screening for the early detection of cancer - I. The temporal history of a progressive disease state. *Math. Biosci.* **40**, 1–59.
- Andersen, P. K. and Gill, R. D. (1982). Cox's regression model for counting processes: a large sample study. *Ann. Statist.*, **10**, 1100–20.
- Breslow, N. and Day, N. E. (1980). *Statistical methods in cancer research (Vol. 1): the analysis of case-control studies*. Lyon: World Health Organization.
- Day, N. E. and Walter, S. D. (1984). Simplified models of screening for chronic disease. Estimation procedure from mass screening programs. *Biometrics* **40**, 1 – 14.

- Foulds, L. (1969). Characteristics of neoplasms. In *Neoplastic Development*, Volume 1, Ed. L. Foulds, pp. 97 – 136. London: Academic Press.
- Ghosh, D. (2002). Goodness of fit methods for the additive risk model in tumorigenicity experiments. *Biometrics* **28**, 721 – 726.
- Ghosh, D. (2004). Nonparametric and semiparametric inference for models of tumor size and metastasis. Technical report, Department of Biostatistics, University of Michigan.
- Gregori, G., Hanin, L., Luebeck, G., Moolgavkar, S., and Yakovlev, A. Y. (2002). Testing goodness of fit for stochastic models of carcinogenesis. *Math. Biosci.* **175**, 13 – 29.
- Hjort, N. (1990). Goodness of fit tests in models for life history based on cumulative hazard rates. *Ann. Statist.* **18**, 1221 – 1258.
- Kim, J. and Lee, S. (1998). Two-sample goodness-of-fit tests for additive risk models with censored observations. *Biometrika* **85**, 593 – 603.
- Lin, D. Y., Oakes, D., and Ying, Z. (1998). Additive hazards regression with current status data. *Biometrika* **85**, 289–298.
- Lin, D. Y., Wei, L. J., Yang, I. and Ying, Z. (2000). Semiparametric regression for the mean and rate functions of recurrent events. *J. Roy. Statist. Soc. Ser. B* **62**, 711–730.
- Lin, D. Y., Wei, L. J. and Ying, Z. (1993). Checking the Cox model with cumulative sums of martingale-based residuals. *Biometrika* **80**, 557–572.
- Lin, D. Y. and Ying, Z. (1994). Semiparametric analysis of the additive risk model. *Biometrika* **81**, 61–71.
- Pollard, D. (1990), *Empirical processes: theory and applications*, Hayward, CA: Institute of Mathematical Statistics.

- Shen, Y. and Zelen, M. (1999). Parametric estimation procedures for screening programmes: stable and nonstable disease models for multimodality case finding. *Biometrika* **86**, 503 – 515.
- Xu, J. L., and Prorok, P. C. (1997). Nonparametric estimation of solid cancer size at metastasis and probability of presenting with metastasis at detection. *Biometrics* **53**, 579 – 591.
- Xu, J. L. and Prorok, P. C. (1998). Estimating a distribution function of the tumor size at metastasis. *Biometrics* **54**, 859 – 864.
- Yakovlev, A. Y. and Tsodikov, A. D. (1996). *Stochastic models of tumor latency and their biostatistical applications*. Singapore: World Scientific Press.
- Zelen, M. and Feinleib, M. (1969). On the theory of screening for chronic diseases. *Biometrika* **56**, 601 – 614.

D. Ghosh, Department of Biostatistics, University of Michigan, Ann Arbor, Michigan, USA, 48109-2029

Appendix

A.1. Derivations of Asymptotic Results for Case I situation

We assume the usual regularity conditions as in Lin and Ying (1994). Consider the following multiparameter process:

$$W(t, z) = n^{-1/2} \sum_{i=1}^n \int_0^t f(\mathbf{Z}_i) I(\mathbf{Z}_i \leq z) d\widehat{M}_i(u),$$

where $I(\mathbf{Z}_i \leq z)$ takes a value of 1 if all the j components of \mathbf{Z}_i are less than z and 0 otherwise. If $f(\mathbf{Z}) = 1$ and $t = \tau$, the j th component of $W(t, z)$, $W_j(t, z)$, is $\widehat{W}_j(z)$. If $z = \infty$ and $f(\mathbf{Z}) = \mathbf{Z}$, then $W_j(t, z)$, suitably normalized, yields $U_j^*(t)$. To demonstrate

the weak convergence of $W_j(z)$ and $U_j^*(t)$, it thus suffices to prove the weak convergence of $W(t, z)$ for $t \in [0, \tau]$. Taylor series expansion and standard algebraic manipulations yield

$$W(t, z) = n^{-1/2} \sum_{i=1}^n \int_0^t \{f(\mathbf{Z}_i)I(\mathbf{Z}_i \leq z) - \mathbf{H}(u, z)\} dM_i(u) - \mathbf{G}(t, z)n^{1/2}(\widehat{\beta} - \beta_0) + o_P(1),$$

where $\mathbf{H}(t, z) = n^{-1} \sum_{i=1}^n \int_0^t R_i(u)f(\mathbf{Z}_i)I(\mathbf{Z}_i \leq z) / \sum_{i=1}^n R_i(u)$,

$$\mathbf{G}(t, z) = n^{-1} \sum_{i=1}^n \int_0^t \{f(\mathbf{Z}_i)I(\mathbf{Z}_i \leq z) - \mathbf{H}(u, z)\} R_i(u)\mathbf{Z}_i du,$$

and $o_P(1)$ denotes a random variable that converges to zero in probability. By the uniform strong law of large numbers (Pollard, 1990, p. 41) and the strong consistency of $\widehat{\beta}$, \mathbf{H} and \mathbf{G} converge to deterministic functions, \mathbf{h} and \mathbf{g} say. By application of the martingale central limit theorem (Fleming and Harrington, 1991, Theorem 5.3.5) and the iid representation for $n^{1/2}(\widehat{\beta} - \beta_0)$, we have that

$$W(t, z) = n^{-1/2} \sum_{i=1}^n \Psi_i(t, z) + o_P(1),$$

where

$$\begin{aligned} \Psi_i(t, z) &= \int_0^t \{f(\mathbf{Z}_i)I(\mathbf{Z}_i \leq z) - \mathbf{h}(u)\} dM_i(u) \\ &\quad - \mathbf{g}(t, z)^T \mathbf{A}^{-1} \int_0^\tau \left\{ \mathbf{Z}_i - \frac{\mathbf{r}^{(1)}(u)}{r^{(0)}(u)} \right\} dM_i(u), \end{aligned}$$

and $\mathbf{r}^{(k)}(\beta, t)$ is the limit of $n^{-1} \sum_{j=1}^n R_j(t)\mathbf{Z}_j^{\otimes k}$, $k = 0, 1$. The multivariate central limit theorem implies that $W(t, z)$ converges in finite distribution to a Gaussian process with mean zero and covariance function $\sigma(t, t', z, z') = E\{\Psi_1(t, z)\Psi_1(t', z')^T\}$. By the arguments in the Appendix of Lin et al. (2000), $W(t, z)$ is tight. The finite dimensional convergence and tightness of $W(t, z)$ imply its weak convergence. The weak convergence of $W_j(x)$ and $U_j^*(t)$ are established.

Note that by the martingale structure of $M_i(t)$ ($i = 1, \dots, n$), its variance is $E\{N_i(t)\}$. Let $\widehat{W}(t, z) = n^{-1/2} \sum_{i=1}^n \widehat{\Psi}_i(t, z)G_i$, where

$$\begin{aligned} \widehat{\Psi}_i(t, z) &= \int_0^t \{f(\mathbf{Z}_i)I(\mathbf{Z}_i \leq z) - \mathbf{H}(u, z)\} dN_i(u) \\ &\quad - \mathbf{G}(t, z)^T \mathbf{A}^{-1} \int_0^\tau \left\{ \mathbf{Z}_i - \frac{\sum_{j=1}^n R_j(u)\mathbf{Z}_j}{\sum_{j=1}^n R_j(u)} \right\} dN_i(u). \end{aligned}$$

Conditional on the data $\{V_i, \delta_i, \mathbf{Z}_i\}$, the random components in $\widehat{W}(t, z)$ are $\{G_i\}$. By the multivariate central limit theorem, $\widehat{W}(t, z)$, conditional on the data, converges in finite dimensions to a mean-zero Gaussian process with covariance function $\widehat{\sigma}(t, t', z, z') = n^{-1} \sum_{i=1}^n \widehat{\Psi}_i(t, z) \widehat{\Psi}_i(t', z')^T$. By the strong consistency of $\widehat{\beta}$, $\widehat{H}_0(t)$ and repeated applications of the uniform strong law of large numbers (Pollard, 1990, p. 41), $\widehat{\sigma} \rightarrow \sigma$ almost surely. If $\widehat{W}(t, z)$ is tight, then $\widehat{W}(t, z)$ has the same limiting distribution as $W(t, z)$. The components of $\widehat{\Psi}_i(t, z)$ are compositions of monotone functions in t . Since the class of monotone functions is manageable (Pollard, 1990, p. 41), it follows by the functional central limit theorem (Pollard, 1990, p. 41) that $\widehat{W}(t, z)$ is tight. This shows that the null distribution of $W(t, z)$ can be approximated by that of $\widehat{W}(t, z)$.

A.2. Derivations of Asymptotic Results for Case II situation

Regularity conditions, similar to those in Andersen and Gill (1982, Theorem 4.1) are imposed for $\{V_i, \delta_i, \mathbf{Z}_i\}$, $i = 1, \dots, n$. The strong consistency of $\tilde{\beta}$, \widehat{H}_0 , and $\tilde{\mathbf{I}}$ follows from the arguments in the Appendix of Lin et al. (2000).

As in the previous section, we consider a generalization of the goodness of fit processes utilized in §3.2:

$$\tilde{W}(t, z) = n^{-1/2} \sum_{i=1}^n \int_0^t f(\mathbf{Z}_i) I(\mathbf{Z}_i \leq z) d\tilde{M}_i(u).$$

Reductions to $\tilde{W}_j(z)$ and $\tilde{U}_j(t)$ are straightforward. It thus suffices to prove the weak convergence of $W(t, z)$ for $t \in [0, \tau^*]$. Utilizing arguments similar to those in Appendix A.1.,

$$\tilde{W}(t, z) = n^{-1/2} \sum_{i=1}^n \int_0^t \left\{ f(\mathbf{Z}_i) I(\mathbf{Z}_i \leq z) - \frac{\mathbf{S}_f(\beta_0, u, z)}{S^{(0)}(\beta_0, u)} \right\} d\tilde{M}_i(u) - \mathbf{B}_f(\beta^*, t, z) n^{1/2} (\tilde{\beta} - \beta_0),$$

where $\mathbf{S}_f(\beta, u, z) = n^{-1} \sum_{i=1}^n R_i(u) \exp\{\beta^T \mathbf{Z}_i^*(u)\} f(\mathbf{Z}_i) I(\mathbf{Z}_i \leq z)$,

$$\mathbf{B}_f(\beta, t, z) = n^{-1} \sum_{i=1}^n \int_0^t R_i(u) e^{\beta^T \mathbf{Z}_i^*(u)} f(\mathbf{Z}_i) I(\mathbf{Z}_i \leq z) \{ \mathbf{Z}_i^*(u) - \mathbf{S}^{(1)}(\beta, u) / S^{(0)}(\beta, u) \} d\widehat{H}_0^*(u; \beta),$$

and β^* is on the line segment between β_0 and $\tilde{\beta}$. The uniform strong law of large numbers (Pollard, 1990, p. 41) and the strong consistency of $\tilde{\beta}$ and $\widehat{H}_0^*(t)$ yield the convergence

of \mathbf{S}_f and \mathbf{B}_f to deterministic functions \mathbf{s}_f and \mathbf{b}_f , respectively. An iid representation of $n^{1/2}(\tilde{\beta} - \beta_0)$ follows from Rebdollo's inequality (Fleming and Harrington, 1991, Ch. 5). This fact and the martingale central limit theorem (Fleming and Harrington, 1991, Theorem 5.3.5) yield

$$\tilde{W}(t, z) = n^{-1/2} \sum_{i=1}^n \Psi_i(t, z) + o_P(1),$$

where

$$\begin{aligned} \tilde{\Psi}_i(t, z) = & \int_0^t \left\{ f(\mathbf{Z}_i) I(\mathbf{Z}_i \leq z) - \frac{\mathbf{s}_f(\beta_0, u, z)}{s^{(0)}(\beta_0, u)} \right\} dM_i(u) \\ & - \mathbf{b}_f(\beta_0, t, z)^T \mathcal{I}^{-1} \int_0^\tau \left\{ \mathbf{Z}_i^*(u) - \frac{\mathbf{s}^{(1)}(\beta_0, u)}{s^{(0)}(\beta_0, u)} \right\} dM_i(u), \end{aligned}$$

and $\mathbf{s}^{(k)}(\beta, t)$ is the limit of $\mathbf{S}^{(k)}(\beta, t)$, $k = 0, 1$. By arguments as in Appendix A.1., $\tilde{W}(t, z)$ converges weakly.

Let $\tilde{W}^*(t, z) = n^{-1/2} \sum_{i=1}^n \tilde{\Psi}_i(t, z) G_i$, where

$$\begin{aligned} \tilde{\Psi}_i(t, z) = & \int_0^t \left\{ f(\mathbf{Z}_i) I(\mathbf{Z}_i \leq z) - \frac{\mathbf{S}_f(\tilde{\beta}, u, z)}{S^{(0)}(\tilde{\beta}, u)} \right\} dN_i(u) \\ & - \mathbf{B}_f(\tilde{\beta}, t, z)^T \tilde{\mathbf{I}}^{-1} \int_0^\tau \left\{ \mathbf{Z}_i^*(u) - \frac{\mathbf{S}^{(1)}(\tilde{\beta}, u)}{S^{(0)}(\tilde{\beta}, u)} \right\} dN_i(u). \end{aligned}$$

Note that conditional on the data $\{V_i, \delta_i, \mathbf{Z}_i\}$, the random components in $\tilde{W}(t, z)$ are $\{G_i\}$. By arguments from Appendix A.1., the null distribution of $\tilde{W}(t, z)$ can be approximated by that of $\tilde{W}^*(t, z)$.



Table 1. Empirical sizes of s_1 and \tilde{s}_1

n	α	Case I			Case II		
		$U(0, 5)$	$U(0, 3)$	$U(0, 2)$	$U(0, 5)$	$U(0, 3)$	$U(0, 2)$
50	0.05	0.05	0.05	0.04	0.04	0.03	0.03
	0.10	0.09	0.08	0.08	0.11	0.09	0.08
	0.15	0.14	0.14	0.13	0.13	0.12	0.12
	0.20	0.20	0.18	0.18	0.21	0.19	0.18
100	0.05	0.04	0.05	0.05	0.05	0.04	0.05
	0.10	0.09	0.09	0.09	0.09	0.08	0.09
	0.15	0.15	0.15	0.14	0.15	0.14	0.14
	0.20	0.19	0.19	0.19	0.19	0.18	0.19
200	0.05	0.05	0.05	0.05	0.05	0.05	0.05
	0.10	0.10	0.09	0.09	0.10	0.09	0.09
	0.15	0.14	0.14	0.15	0.15	0.14	0.15
	0.20	0.19	0.19	0.20	0.20	0.19	0.20

Note: α denotes level of significance.

Table 2 Empirical powers of s_1 and \tilde{s}_1

n	α	Case I			Case II		
		$U(0, 5)$	$U(0, 3)$	$U(0, 2)$	$U(0, 5)$	$U(0, 3)$	$U(0, 2)$
50	0.05	0.80	0.68	0.65	0.70	0.68	0.65
	0.10	0.84	0.79	0.77	0.74	0.69	0.67
	0.15	0.87	0.80	0.73	0.77	0.63	0.70
	0.20	0.91	0.88	0.85	0.81	0.78	0.75
100	0.05	0.87	0.85	0.83	0.75	0.71	0.69
	0.10	0.89	0.87	0.84	0.69	0.65	0.71
	0.15	0.92	0.88	0.87	0.82	0.77	0.74
	0.20	0.96	0.92	0.90	0.86	0.80	0.78
200	0.05	0.90	0.87	0.85	0.80	0.77	0.73
	0.10	0.92	0.90	0.85	0.82	0.80	0.75
	0.15	0.95	0.92	0.89	0.85	0.82	0.77
	0.20	0.98	0.9	0.92	0.88	0.85	0.80

Note: α denotes level of significance used in determining power.

Table 3. *Empirical sizes of h_1 and \tilde{h}_1*

n	α	Case I			Case II		
		$U(0, 5)$	$U(0, 3)$	$U(0, 2)$	$U(0, 5)$	$U(0, 3)$	$U(0, 2)$
50	0.05	0.05	0.04	0.04	0.05	0.05	0.04
	0.10	0.09	0.09	0.09	0.09	0.08	0.09
	0.15	0.13	0.13	0.13	0.16	0.16	0.15
	0.20	0.23	0.21	0.20	0.21	0.22	0.21
100	0.05	0.05	0.05	0.04	0.04	0.05	0.04
	0.10	0.08	0.09	0.09	0.10	0.09	0.10
	0.15	0.15	0.14	0.15	0.15	0.14	0.14
	0.20	0.19	0.18	0.19	0.18	0.19	0.18
200	0.05	0.04	0.05	0.04	0.05	0.05	0.05
	0.10	0.10	0.09	0.10	0.10	0.09	0.10
	0.15	0.14	0.14	0.15	0.14	0.14	0.15
	0.20	0.19	0.19	0.20	0.20	0.19	0.19

See note to Table 1.

Table 4. *Empirical powers of h_1 and \tilde{h}_1*

n	α	Case I			Case II		
		$U(0, 5)$	$U(0, 3)$	$U(0, 2)$	$U(0, 5)$	$U(0, 3)$	$U(0, 2)$
50	0.05	0.80	0.77	0.75	0.45	0.43	0.42
	0.10	0.84	0.81	0.77	0.47	0.44	0.43
	0.15	0.85	0.83	0.79	0.58	0.52	0.47
	0.20	0.89	0.85	0.81	0.61	0.55	0.52
100	0.05	0.85	0.83	0.81	0.48	0.46	0.45
	0.10	0.88	0.86	0.83	0.52	0.50	0.48
	0.15	0.90	0.87	0.86	0.62	0.58	0.55
	0.20	0.93	0.90	0.88	0.63	0.59	0.57
200	0.05	0.89	0.87	0.83	0.54	0.52	0.50
	0.10	0.92	0.90	0.86	0.57	0.55	0.53
	0.15	0.94	0.92	0.89	0.60	0.59	0.56
	0.20	0.96	0.93	0.90	0.67	0.65	0.61

See note to Table 2.