

The Effects of Misspecifying Cox's
Regression Model on Randomized Treatment
Group Comparisons

Greg DiRienzo*

*Harvard University, dirienzo@sdac.harvard.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/harvardbiostat/paper3>

Copyright ©2003 by the author.

The Effects of Misspecifying Cox's Regression Model on Randomized Treatment Group Comparisons

A.G. DiRienzo and S.W. Lagakos

1. Introduction

Hypothesis tests arising from Cox's proportional hazards model (Cox, 1972) are often used to compare randomized treatment groups with respect to the distribution of a failure time outcome. Some of these tests adjust for covariates that may be predictive of outcome, while others, and most notably, the log-rank test, do not. In addition to adjusting for any imbalances that may arise between treatment groups, covariate-adjusted test may enjoy greater efficiency than that of the log-rank test. Tsiatis et al. (1985) demonstrated the gain in efficiency of covariate-adjusted tests relative to the log-rank test when the working proportional hazards model is properly specified. Slud (1991) provided asymptotic relative efficiency formulae of the log-rank test to the optimal score test that arises from a properly specified model for covariates when the effect of treatment is multiplicative on the survival time hazard function. Lagakos and Schoenfeld (1984) studied the effects of various types of model misspecification on the power of tests based on Cox's model.

An important consideration in the application of these tests is their validity when the underlying proportional hazards working model is misspecified. Recent work has shown that the impact of model misspecification on the validity of resulting tests hinges on whether the distribution of the potential censoring time either (i) is conditionally independent of treatment group given covariates or conditionally independent of covariates given treatment group, or (ii) depends on both treatment group and covariates. In the first case, resulting test statistics have an asymptotic normal distribution with mean zero under the null hypothesis and that consistent variance estimates are readily obtainable (see Kong and Slud, 1997 and DiRienzo and Lagakos, 2001a). In the second case, the asymptotic mean of the test statistic is not necessarily equal to zero under the null hypothesis when the proportional hazards working model is misspecified. In such cases, the bias of tests can be large, as was demonstrated in DiRienzo and Lagakos (2001a, 2001b).

In this chapter we summarize the properties of hypothesis tests derived from proportional hazards regression models. We introduce notation and define uncorrected statistics in Section 2. In Section 3 we describe conditions necessary for the asymptotic

1 validity of these test statistics, and also discuss efficiency considerations and effects of 1
2 model misspecification on the power of uncorrected test statistics. We describe a class 2
3 of corrected test statistics for use when censoring depends on both treatment group and 3
4 covariates in [Section 4](#), and also examine estimation procedures and the efficiency of 4
5 such bias-corrected tests. We provide some recommendations for the use of these tests 5
6 in [Section 5](#), and give MATLAB code for the computation of the various test statistics 6
7 in [Appendix A](#). 7

10 2. Notation and statistics 10

11
12
13 Let the continuous random variable T denote time from randomization to failure and let 13
14 C denote a potential censoring time. Assume that we observe $T^* = \min(T, C)$ and the 14
15 indicator $\delta = \mathbb{1}(T \leq C)$ of whether T is observed ($\delta = 1$) or right-censored ($\delta = 0$). Let 15
16 the binary random variable X denote treatment group and let W denote a $q \times 1$ vector of 16
17 bounded baseline covariates. Throughout this paper we assume that censoring acts non- 17
18 informatively, that is, $T \perp C \mid (X, W)$, and also that $X \perp W$, as is the case in most ran- 18
19 domized clinical trials. The true conditional hazard functions of T and C given (X, W) 19
20 are denoted by $\kappa(t \mid X, W)$ and $\kappa_C(t \mid X, W)$, respectively, and are not necessarily of a 20
21 proportional hazards form. The observed data is assumed to consist of n independent 21
22 and identically distributed realizations of (T^*, δ, X, Z^*) , denoted $(T_i^*, \delta_i, X_i, Z_i^*)$ for 22
23 $i = 1, \dots, n$, where Z^* is a $p \times 1$ vector whose components are bounded functions of 23
24 W . 24

25 The null hypothesis of interest is $H_0: X \perp T \mid W$; that is, that the failure time distri- 25
26 bution does not depend on treatment group, under which we will denote $\kappa(t \mid X, W)$ by 26
27 $\kappa(t \mid W)$. Consider tests of H_0 that are based on statistics of the form 27

$$28 \quad n^{-1/2}U_n = \sum_{i=1}^n \int_0^\infty n^{-1/2}G_n(t)\{X_i - \mathcal{E}_n(t)\}dN_i(t), \quad (1) \quad 29$$

30 where 30

$$31 \quad \mathcal{E}_n(t) = \sum_{j=1}^n Y_j(t)\psi_n(Z_j)X_j / \sum_{j=1}^n Y_j(t)\psi_n(Z_j), \quad 32$$

33
34
35
36
37 $Y_i(t) = \mathbb{1}(T_i^* \geq t)$, $N_i(t) = \delta_i \mathbb{1}(T_i^* \leq t)$ and $\psi_n(\cdot)$ is a nonrandom bounded function 37
38 whose form is known but whose parameters can be estimated from the data. The co- 38
39 variates Z_i are some bounded function of Z_i^* , $i = 1, \dots, n$. The bounded predictable 39
40 process $G_n(\cdot)$ is also assumed to be nonrandom, converging uniformly in probability 40
41 to a bounded function $G(\cdot)$. It may be the case that one would want to consider time- 41
42 dependent covariates, for example an external ancillary covariate process ([Kalbfleisch](#) 42
43 [and Prentice, 1980, p. 123](#)). Although results hold when the components of W_i and Z_i 43
44 are uniformly bounded and predictable functions of time, for ease of notation we only 44
45 consider fixed covariates. 45

1 Statistics of the form in (1) arise as the numerator of partial likelihood score tests of
2 $\alpha = 0$ based on working proportional hazards models for $\kappa(t | X, W)$ that take the form

$$3 \exp(\alpha X_i) \psi(\beta; Z_i) h(t). \tag{2} 4$$

5 The working model (2) is misspecified when it is not equivalent to $\kappa(t | X, W)$, in
6 which case the parameters α , β , and $h(\cdot)$ have no simple interpretation. The statistic
7 $n^{-1/2}U_n$ should then generally be viewed simply as statistic from which tests of H_0
8 may be derived.

9 Popular choices for $G(\cdot)$ and $\psi(\cdot)$ are $G_n(t) = 1$ and $\psi(\beta; Z_i) = \exp(\beta^\top Z_i)$, re-
10 sulting in $\psi_n(Z_i) = \exp(\hat{\beta}^\top Z_i)$, where $\hat{\beta}$ is the restricted maximum partial likelihood
11 estimator of β obtained by fitting the model with $\alpha = 0$ (Cox, 1972). Here the prob-
12 ability limit of $\psi_n(Z)$ is $\exp(\tilde{\beta}^\top Z)$, where $\tilde{\beta}$ is the probability limit of $\hat{\beta}$ (Lin and
13 Wei, 1989). Note that $\tilde{\beta} = \beta$ when the model (2) is properly specified. Another special
14 case of (1) is the class of weighted log-rank statistics (Cox and Oakes, 1984, p. 124),
15 where $\psi(\beta; Z_i) = \psi_n(Z_i) = 1$, and where the most commonly used choice for $G_n(\cdot)$ is
16 the identity function, yielding the ordinary log-rank test. In general, $\psi(\beta; Z)$ can also
17 depend on t , so long as it is uniformly bounded.

20 3. Conditions for valid tests

21 Suppose that either $C \perp X | W$ or $C \perp W | X$. Then the test statistic $n^{-1/2}U_n$ has an
22 asymptotic normal distribution with mean 0 under H_0 , regardless of whether or not the
23 model (2) is misspecified (DiRienzo and Lagakos, 2001a). Furthermore, when either of
24 these conditions hold, consistent estimates of the variance of $n^{-1/2}U_n$ are easily derived,
25 yielding asymptotically valid inference whether or not the relationship between T and
26 (X, W) is properly specified.

27 The condition $C \perp X | W$ is usually satisfied in a randomized clinical trial when
28 the only form of censoring is administrative or end-of-study censoring; that is, when C
29 represents the time from enrollment of a subject into the study until the time the data
30 are analyzed. However, when censoring can arise from premature study discontinuation
31 or loss-to-follow-up, it is well known that this condition may not hold. The condition
32 $C \perp W | X$ holds when there is a dependency of censoring on treatment group which
33 does depend on the covariates. To provide some insight into why either of these condi-
34 tions are necessary for valid inference, note that at baseline (that is, when $t = 0$), the dis-
35 tribution of W is independent of X because of randomization; when either $C \perp X | W$
36 or $C \perp W | X$ holds and H_0 is true, it is implied that $X \perp W | Y(t) = 1, t > 0$, which is
37 necessary for $n^{-1/2}U_n$ to have mean 0 asymptotically. For a proof of these results, see
38 Appendix A of DiRienzo and Lagakos (2001a) or Kong and Slud (1997).

39 We now provide test statistics for use when either $C \perp X | W$ or $C \perp W | X$ holds. It
40 follows from Kong and Slud (1997) that under H_0 , $n^{-1/2}U_n$ can be expressed as

$$41 n^{-1/2}U_n = n^{-1/2} \sum_{i=1}^n Q_i + o_p(1), 42$$

1 where

$$2 \quad Q_i = \int_0^\infty \{X_i - \mu(t)\} \{dN_i(t) - \rho(t)Y_i(t)\psi(\tilde{\beta}; Z_i) dt\},$$

$$3 \quad \mu(t) = E\{Y(t)\psi(\tilde{\beta}; Z)X\} / E\{Y(t)\psi(\tilde{\beta}; Z)\},$$

$$4 \quad \rho(t) = E\{Y(t)\kappa(t | W)\} / E\{Y(t)\psi(\tilde{\beta}; Z)\}.$$

5 It is easily verified (cf. Kong and Slud, 1997 or using Lemmas 1 and 2 in DiRienzo
6 and Lagakos, 2001a) that under H_0 , Q_i has mean 0 when $C \perp X | W$ or $C \perp W | X$.
7 This implies that $n^{-1/2}U_n$ is asymptotically normal with mean zero and variance equal
8 to the variance of Q_i . As shown in Kong and Slud (1997), a consistent estimate of the
9 variance of $n^{-1/2}U_n$ is

$$10 \quad \frac{1}{n} \sum_{i=1}^n (\hat{Q}_i - \bar{Q})^2,$$

11 where

$$12 \quad \hat{Q}_i = \int_0^\infty \{X_i - \bar{X}(t)\} \left\{ dN_i(t) - \frac{Y_i(t)\psi(\hat{\beta}; Z_i)}{\sum_{j=1}^n Y_j(t)\psi(\hat{\beta}; Z_j)} d\bar{N}(t) \right\},$$

$$13 \quad \bar{X}(t) = \frac{\sum_{i=1}^n Y_i(t)X_i}{\sum_{i=1}^n Y_i(t)},$$

$$14 \quad \bar{N}(t) = \sum_{i=1}^n N_i(t) \text{ and } \bar{Q} = (1/n) \sum_{i=1}^n \hat{Q}_i.$$

15 Thus, provided that $C \perp X | W$ or that $C \perp W | X$, the test statistic $U_n /$
16 $\sqrt{\{\sum_{i=1}^n (\hat{Q}_i - \bar{Q})^2\}}$ is asymptotically standard normal under H_0 when $C \perp X | W$
17 or $C \perp W | X$, regardless of whether the working model (2) is properly specified. The
18 motivation for replacing $\mu(t)$ with $\bar{X}(t)$ above is that $\mu(t) = E\{X | Y(t) = 1\}$ under H_0
19 when $C \perp X | W$ or $C \perp W | X$. We note that for the special case of the log-rank test,
20 use of the model-based variance estimator of $n^{-1/2}U_n$ also results in a valid asymptotic
21 test, and appears to usually provide nominal finite-sample type I errors (see DiRienzo
22 and Lagakos, 2001a), so that use of a robust variance estimator is not needed.

23 3.1. Efficiency considerations

24 Lagakos and Schoenfeld (1984) investigated the effects of various types of misspeci-
25 fication of the working model (2) on the power of $n^{-1/2}U_n$. When covariates have a
26 multiplicative effect on the true hazard $\kappa(t | X, W)$, but the ratio $\kappa(t | X = 1, W) / \kappa(t |$
27 $X = 0, W)$, is non-constant but either greater or less than one for all $t > 0$, i.e., the
28 hazards do not cross, there is often only a small loss in power. One exception to this is
29 when the ratio $\kappa(t | X = 1, W) / \kappa(t | X = 0, W)$ departs from one only after the ma-
30 jority of failures have occurred; in this case, the loss in power can be great. In contrast,
31 when the ratio $\kappa(t | X = 1, W) / \kappa(t | X = 0, W)$ crosses one, the loss in power is often
32 substantial.

1 Suppose that the effect of covariates in the true model $\kappa(t | X, W)$ is not multiplica- 1
2 tive, that is the ratio $\kappa(t | X = 1, W)/\kappa(t | X = 0, W)$ is a function of W , but that the 2
3 interaction is qualitative, in the sense that $\kappa(t | X = 1, W)/\kappa(t | X = 0, W)$ is either 3
4 greater or less than one for all W . In this case, the loss in the power of $n^{-1/2}U_n$ is not 4
5 in general large unless the discrepancy in the ratio $\kappa(t | X = 1, W)/\kappa(t | X = 0, W)$ 5
6 between levels of W is substantial, especially if larger ratios tend to occur within levels 6
7 of W that are less prevalent. 7

8 More generally, the loss in power of $n^{-1/2}U_n$ can be large when a component of W 8
9 that has a strong effect on the hazard of T is either omitted or misspecified in such a way 9
10 that the direction of its effect is not maintained. Further details on all of these situations 10
11 can be found in [Lagakos and Schoenfeld \(1984\)](#). [Morgan \(1986\)](#) provides a correction 11
12 to [Lagakos and Schoenfeld's \(1984\)](#) asymptotic relative efficiency formula of the log- 12
13 rank test to the score test arising from a properly specified model for covariates. See 13
14 also [Lagakos \(1988\)](#), who derived asymptotic relative efficiency formulae in the one- 14
15 sample problem when evaluating the effect of a misspecified form of a time-dependent 15
16 covariate. 16

19 4. Bias correction 19

21 When the distribution of the censoring variable depends on both treatment group and 21
22 covariates, that is when the conditions $C \perp X | W$ and $C \perp W | X$ both fail to hold, the 22
23 statistic $n^{-1/2}U_n$ in general has a non-zero asymptotic mean under H_0 . One exception 23
24 is when the model (2) is equal to $\kappa(t | X, W)$, i.e., the working proportional hazards 24
25 model is properly specified. [DiRienzo and Lagakos \(2001a, 2001b\)](#) present simulation 25
26 results which demonstrate that the bias of tests based on $n^{-1/2}U_n$ can be severe when 26
27 in this setting and the working proportional hazards model is misspecified. 27

28 In an attempt to correct for this bias, [DiRienzo and Lagakos \(2001b\)](#) present a class 28
29 of tests that are asymptotically standard normal under H_0 regardless of the joint distri- 29
30 bution between C and (X, W) , provided that either the conditional distribution of T 30
31 given (X, W) or the conditional distribution of C given (X, W) is properly modeled. 31
32 Consequently, these tests are more robust than those arising from $n^{-1/2}U_n$ when the 32
33 working model is misspecified, and do not appear to lose much efficiency when the 33
34 working model is correctly specified and bias correction is unnecessary. 34

35 Consider the generalization of (1) given by 35

$$36 \quad n^{-1/2}U_n^* = \sum_{i=1}^n \int_0^\infty n^{-1/2}G_n(t)\varphi(t; X_i, W_i)\{X_i - \mathcal{E}_n^*(t)\}dN_i(t), \quad (3) \quad 36$$

37 where 37

$$38 \quad \mathcal{E}_n^*(t) = \sum_{j=1}^n Y_j^*(t)\psi_n(Z_j)X_j / \sum_{j=1}^n Y_j^*(t)\psi_n(Z_j), \quad 38$$

$$39 \quad Y_i^*(t) = Y_i(t)\varphi(t; X_i, W_i), \quad 39$$

$$\begin{aligned} \varphi(t; X_i, W_i) &= \min\{\text{pr}(C \geq t \mid X_i = 0, W_i), \\ &\text{pr}(C \geq t \mid X_i = 1, W_i)\} / \text{pr}(C \geq t \mid X_i, W_i), \end{aligned} \quad (4)$$

for $i = 1, \dots, n$. Unlike the binary indicator variable Y normally used in Cox's model, $Y_i^*(t)$ can assume any value in the unit interval. Also, note that $\varphi(t; X_i, W_i)$ is only defined when $Z_i = W_i, i = 1, \dots, n$.

At each point in study time when a survival event occurs, this correction strives to remove any imbalances between treatment groups in the distribution of covariates that are caused solely by censoring. Mechanically, at study time t , the correction down-weights, $Y_i^*(t) < 1$, those subjects in the risk set whose risk of censoring is higher in their opposite treatment group; those subjects whose risk of censoring is lower in their opposite treatment group are unweighted, $Y_i^*(t) = Y_i(t) = 1$. To see this analytically, note that under H_0 , the conditional expectation of $Y^*(t)$ given (X, W) is

$$\begin{aligned} \varphi(t; X, W) \text{pr}\{Y(t) = 1 \mid X, W\} \\ &= \varphi(t; X, W) \text{pr}(C \geq t \mid X, W) \text{pr}(T \geq t \mid W) \\ &= \min\{\text{pr}(C \geq t \mid X = 0, W), \text{pr}(C \geq t \mid X = 1, W)\} \text{pr}(T \geq t \mid W), \end{aligned}$$

which is independent of X . The probability limit of $\mathcal{E}_n^*(t)$ under H_0 is thus

$$\frac{E\{Y^*(t)\psi(Z)X\}}{E\{Y^*(t)\psi(Z)\}} = \frac{E[X\psi(Z)E\{Y^*(t) \mid W\}]}{E[\psi(Z)E\{Y^*(t) \mid W\}]} = \pi,$$

where $\pi = E(X)$.

As shown in DiRienzo and Lagakos (2001b), $n^{-1/2}U_n^*$ can be expressed under H_0 as

$$n^{-1/2}U_n^* = n^{-1/2} \sum_{i=1}^n A_i + o_p(1),$$

where

$$\begin{aligned} A_i &= \int_0^\infty G(t)\varphi(t; X_i, W_i)(X_i - \pi) \\ &\times \left\{ dN_i(t) - Y_i(t)\psi(Z_i) \frac{E\{Y^*(t)\kappa(t \mid W)\}}{E\{Y^*(t)\psi(Z)\}} dt \right\}, \end{aligned}$$

and the A_i are independent and identically distributed with mean zero.

A consistent estimator of the variance of $n^{-1/2}U_n^*$ is

$$V_n = \frac{1}{n} \sum_{i=1}^n (A_i^{(n)} - \bar{A}^{(n)})^2, \quad (5)$$

where

$$A_i^{(n)} = \int_0^\infty G_n(t) \varphi(t; X_i, W_i) (X_i - \bar{X}) \times \left\{ dN_i(t) - \frac{Y_i(t) \psi_n(Z_i)}{\sum_{j=1}^n Y_j^*(t) \psi_n(Z_j)} \sum_{j=1}^n \varphi(t; X_j, W_j) dN_j(t) \right\},$$

\bar{X} is the mean of $\{X_1, \dots, X_n\}$ and $\bar{A}^{(n)}$ is the mean of $\{A_1^{(n)}, \dots, A_n^{(n)}\}$. Hence, regardless of the joint distribution between C and (X, W) , $n^{-1/2} U_n^* / \sqrt{V_n}$ asymptotically has the standard normal distribution under H_0 whether or not the working model is properly specified. It follows that if the working model (2) is properly specified, $n^{-1/2} U_n^* / \sqrt{V_n}$ is asymptotically standard normal under H_0 regardless of whether $\text{pr}(C \geq t | X, W)$ is properly specified and of the dependency between C and (X, W) .

In practice, $\varphi(\cdot)$ will often be unknown. Let $\hat{\varphi}(t; X_i, W_i)$ denote an estimator of $\varphi(t; X_i, W_i)$. One would then calculate

$$\hat{U}_n^* = \sum_{i=1}^n \int_0^\infty G_n(t) \hat{\varphi}(t; X_i, W_i) \{X_i - \hat{\mathcal{E}}_n^*(t)\} dN_i(t)$$

instead of (3) and

$$\hat{A}_i^{(n)} = \int_0^\infty G_n(t) \hat{\varphi}(t; X_i, W_i) (X_i - \bar{X}) \times \left\{ dN_i(t) - \frac{Y_i(t) \psi_n(Z_i)}{\sum_{j=1}^n \hat{Y}_j^*(t) \psi_n(Z_j)} \sum_{j=1}^n \hat{\varphi}(t; X_j, W_j) dN_j(t) \right\}$$

instead of $A_i^{(n)}$ in (5), where $\hat{Y}_i^*(t) = Y_i(t) \hat{\varphi}(t; X_i, W_i)$, and $\hat{\mathcal{E}}_n^*(\cdot)$ is obtained by substituting $\hat{Y}_i^*(\cdot)$ for $Y_i^*(\cdot)$ in $\mathcal{E}_n^*(\cdot)$, $i = 1, \dots, n$. Denote this variance estimate by \hat{V}_n .

Some methods for estimating $\varphi(t; X, W)$ are given in DiRienzo and Lagakos (2001b). These include the nonparametric regression methods of McKeague and Utikal (1990) as well as Cox's (1972) proportional hazards regression models. If the covariates are discrete with relatively few levels, then a stratified, left-continuous Kaplan–Meier estimator (Kaplan and Meier, 1958) of censoring can be calculated for each treatment group within each level of the covariate space.

For example, an estimate for $\varphi(t; X, W)$ can be obtained via the stratified proportional hazards model for $\kappa_C(t | X, W)$,

$$\lambda^{(X)}(t) \exp(\gamma^{(X)} Z^C),$$

where Z_i^C is some bounded function of Z_i^* , $i = 1, \dots, n$. The maximum partial-likelihood estimator, $\hat{\gamma}_t^{(X)}$, and the Breslow (1972, 1974) estimator of the baseline cumulative hazard function of censoring, $\hat{\Lambda}^{(X)}(t)$, may then be calculated within each treatment group at each censoring time, using data accumulated before that time, and the continuous estimator

$$\hat{\text{pr}}(C_i \geq t | X_i, Z_i^C) = \exp\{-\hat{\Lambda}^{(X_i)}(t) \exp(\hat{\gamma}_t^{(X_i)} Z_i^C)\}$$



1 obtained by linear interpolation between censoring times of $\widehat{\Lambda}_i^{(X_i)}(t)$. Here estimation 1
2 was stratified on X , but stratification may additionally be based on any covariate that 2
3 might possibly have a strong interaction with treatment. 3

4 When $\varphi(t; X, W)$ is estimated using a semiparametric or nonparametric model, 4
5 $\widehat{\varphi}(t; X_i, W_i)$ contains estimates of an infinite dimensional parameter, for which case 5
6 a consistent estimate of the variance of $n^{-1/2}\widehat{U}_n^*$ would not necessarily be given by \widehat{V}_n . 6
7 However, given the choice for an estimate of $\varphi(t; X, W)$, if it can be shown that \widehat{U}_n^* is 7
8 asymptotically linear, then the nonparametric bootstrap estimate of variance of \widehat{U}_n^* will 8
9 be consistent (Gill, 1989). DiRienzo and Lagakos (2001b) have shown via simulation 9
10 that when using a semiparametric proportional hazards model to calculate $\widehat{\varphi}(t; X_i, W_i)$, 10
11 the variance estimate \widehat{V}_n appears to be adequate. 11

12 For any given data set, there is no guarantee that it will be possible to specify and 12
13 estimate $\varphi(\cdot)$ well enough to make the correction for a misspecified model for T reli- 13
14 able. It is thus of utmost importance to check and validate the fit of both the model for 14
15 censoring and survival. Some well known techniques for checking the appropriateness 15
16 of proportional hazards regression models are given in Lin et al. (1993) and Klein and 16
17 Moeschberger (1997). 17

18 A related consideration in the use of bias-adjusted tests are the relative efficiencies. 18
19 When the working proportional hazards model (2) is properly specified, i.e., equal to 19
20 $\kappa(t | X, W)$, then the uncorrected, fully model-based test of H_0 is asymptotically valid 20
21 regardless of the dependency between C and (X, W) . In this situation, it is of interest 21
22 to examine the relative efficiency of the corrected test to that of the uncorrected test and 22
23 determine if there are situations for which unnecessary use of the corrected test could 23
24 lead to loss in power. 24

25 DiRienzo and Lagakos (2001b) provide formulae for the asymptotic mean and vari- 25
26 ance of $n^{-1/2}\widehat{U}_n$ and $n^{-1/2}\widehat{U}_n^*$ under the contiguous alternative $H_n: \alpha = c/\sqrt{n}$, for 26
27 some constant c , when the true hazard for T is given by 27

$$\kappa(t | X_i, W_i) = \exp(\alpha X_i)\psi(\beta, W_i)h(t).$$

28
29
30 That is, the working proportional hazards model is properly specified and calculation of 30
31 a corrected test is unnecessary since the uncorrected, fully model-based test of H_0 is as- 31
32 ymptotically valid. In their accompanying simulations, the empirical relative efficiency 32
33 of the corrected test to that of the uncorrected test appears to almost always be close to 33
34 one. 34

35 Other choices for the functional form of $\varphi(\cdot)$ may be of interest; one example is 35
36 $\varphi(t; X, W) = 1/\text{pr}(C \geq t | X, W)$. However, using simulations, DiRienzo and La- 36
37 gakos (2001b) have found that this choice for $\varphi(\cdot)$ can be much less efficient than the 37
38 choice (4). They present an efficacy formula for the corrected test; this may be used to 38
39 compare the efficiencies of tests using different choices for $\varphi(\cdot)$. 39
40
41

42 5. Discussion 42

43
44 Given the wide use of statistical tests based on Cox's regression model, especially in 44
45 medical applications, and considering the importance of decisions that are reached from 45

1 these analyses, an understanding of their robustness to misspecification of the model is 1
2 important. Misspecification can occur in many forms, including omitted or mismodelled 2
3 covariates, the omission of treatment by covariate interactions, or a violation of the un- 3
4 derlying proportionality assumption. While goodness-of-fit methods can be applied to 4
5 check model fit (cf. Klein and Moeschberger, 1997), their failure to signal misspecifica- 5
6 tion is no assurance that this is the case and, furthermore, their subjective and post-hoc 6
7 nature can be problematic when a new treatment is being assessed, e.g., in clinical tri- 7
8 als the standard practice is to precisely prespecify how treatment comparisons will be 8
9 made. This chapter has argued that a fundamental question in assessing such robustness 9
10 is whether treatment group and the censoring variable are conditionally independent 10
11 given the underlying covariates, or whether the underlying covariates associated with 11
12 survival are conditionally independent of the censoring variable, given treatment group. 12
13 When either of these conditions apply, then statistical tests arising from fitting a propor- 13
14 tional hazards model, including the popular log-rank test, maintain their validity under 14
15 misspecification of the model-relating treatment and these covariates to the hazard func- 15
16 tion for survival. That is, when either condition holds, the resulting test statistic, when 16
17 standardized by a robust variance estimator, has a distribution under the null hypothe- 17
18 sis of no treatment effect that is asymptotically standard normal, regardless of whether 18
19 or not the model is correctly specified. For the special case of the log-rank test, use of 19
20 the model-based variance estimator to standardize the score statistic arising from the 20
21 assumed model also leads to the desired asymptotic behavior under the null hypothesis. 21
22 Thus, establishment of either of these conditions ensures that the size, or Type I error, 22
23 associated with such tests is not distorted as a result of model misspecification. More- 23
24 over, one or both of the conditions can in practice often either be checked empirically 24
25 or concluded to hold based on the analyst's knowledge of the circumstances that lead to 25
26 censored observations. 26

27 When neither condition holds, that is, when either treatment or the underlying cov- 27
28 variate is not conditionally independent of time to censoring, then tests based on fitting 28
29 a proportional hazards model can be asymptotically biased under the null hypothesis. 29
30 Since in practice the significance levels used to evaluate these tests invariably resort to 30
31 their presumed asymptotic normality, the size of such tests can be seriously biased when 31
32 the working proportional hazards model is misspecified. To avoid or minimize such bi- 32
33 ases, a class of bias-corrected tests can readily be adapted. These tests require knowl- 33
34 edge or estimation of $\phi(t; X, W)$, a function of the conditional distribution of censoring. 34
35 Based on asymptotic considerations and simulations, the corrected test works well in a 35
36 variety of settings, even when the estimated form of $\phi(t; X, W)$ is only approximately 36
37 correct. That is, misspecification of the function $\phi(t; X, W)$ appears to be far less crit- 37
38 ical for the bias-corrected test than does the misspecification of the underlying hazard 38
39 model for the uncorrected test. Furthermore, use of a bias-corrected test when one is 39
40 unnecessary – that is, when the working proportional hazards model happens to be cor- 40
41 rectly specified – does not appear to result in much loss in efficiency. Thus, when there 41
42 is any suspicion that the key conditions for robustness may be violated, use of the bias- 42
43 adjusted tests instead of or as a complement to standard methods is advised. To facilitate 43
44 the computation for the adjusted tests, Appendix A gives MATLAB code for these and 44
45 uncorrected tests. 45

1 **Acknowledgement**

2
3 This work was supported in part from grant AI24643 from the US National Institutes of
4 Health.

7 **Appendix A: MATLAB code for computing statistical tests**

9 We provide below MATLAB code for calculating the uncorrected and corrected score
10 tests presented in this paper. The version of MATLAB used is 5.3.1 (R11.1) along with
11 the Statistics (Version 2.2, R11) and Optimization (Version 2.0, R11) toolboxes.

12 The uncorrected test is calculated using the model-based variance estimator, which
13 is consistent when the working proportional hazards model for $\kappa(t | X, W)$ is prop-
14 erly specified or when the log-rank test is used as the uncorrected test. The cor-
15 rected test is calculated with $G_n(t) = 1$ and using a stratified (by treatment group)
16 proportional hazards model for the conditional distribution of C given (X, W) with
17 $\psi_n(Z) = \exp(\hat{\beta}^\top Z)$, where $\hat{\beta}$ the restricted maximum partial likelihood estimate of
18 β under H_0 . We note, however, that the code can be modified to accommodate other
19 choices for these functions as well as for more covariates that are used below to illu-
20 strate the methods.

21 The observed data consists of the five $n \times 1$ column vectors T_0, d, x, Z_1, Z_2 , where
22 T_0 corresponds to $\{T_i^*\}$, d to $\{\delta_i\}$, x to $\{X_i\}$, Z_1 is the first component of $\{Z_i^*\}$, say
23 $\{Z_{1i}\}$ and Z_2 the second, say $\{Z_{2i}\}$, $i = 1, \dots, n$. Suppose that one wanted to adjust for
24 the covariates $I(Z_1 < 0), Z_2^2$ in the model for T , and calculate a corrected test using
25 a proportional hazards model for C that was conditional on $X, |Z_1|^{-1/2}, Z_2$. Then the
26 MATLAB call would be

```
27  
28 [un1, cor1] = SC(T0, d, x, [(Z1 < 0), (Z2.^2)]),  
29 [((abs(Z1)).^(-.5)), Z2]);
```

31 where the output 1×2 row vector `un1` consists of the uncorrected score statistic and
32 score test, similarly, `cor1` consists of the corrected score statistic and score test.

33 The code for the function `SC.m` and the two functions it calls, `rPLgh.m` and
34 `BRES.m` is given by:

```
35  
36 function [un, cor] = SC(TT, dd, xx, Z1, Z2)  
37 % computes uncorrected *un* and corrected *cor* score statistics  
38 % and tests  
39 % *TT* is the column vector of N possibly right-censored event  
40 % times  
41 % there are assumed to be no TIES in *TT*  
42 % *dd* is the column vector of N indicators I(T<=C)  
43 % *xx* is the column vector of N treatment group indicators  
44 % *Z1* is the N x p matrix of covariates for *T*  
45 % *Z2* is the N x p matrix of covariates for *C*  
-----
```

```
1 % first get the restricted MPLE for T 1
2 %----- 2
3 global T d x Z; 3
4 T=TT; d=dd; x=xx; Z=Z1; 4
5 N = length(T); 5
6 p = size(Z,2); 6
7 th = zeros(1,p); 7
8 options = optimset('GradObj','on','Display','off'); 8
9 rmple = fsolve('rPLgh',th,options); 9
10 clear global 10
11 %----- 11
12 % calculate the MPLE and Breslow estimate of the baseline 11
13 % cumulative hazard of censoring within each treatment group 12
14 %----- 13
15 global T d Z; 14
16 T=TT(xx==0); d=1-dd(xx==0); Z=Z2(xx==0,:); 15
17 N = length(T); 16
18 p = size(Z,2); 17
19 th = zeros(1,p); 18
20 options = optimset('GradObj','on','Display','off'); 19
21 mple0 = fsolve('rPLgh',th,options); 20
22 [L0w,c0] = BRES(T,d,Z,mple0); 21
23 clear global 22
24 %----- 23
25 global T d Z; 24
26 T=TT(xx==1); d=1-dd(xx==1); Z=Z2(xx==1,:); 25
27 N = length(T); 26
28 p = size(Z,2); 27
29 th = zeros(1,p); 28
30 options = optimset('GradObj','on','Display','off'); 29
31 mple1 = fsolve('rPLgh',th,options); 30
32 [L1w, c1] = BRES(T,d,Z,mple1); 31
33 clear global 32
34 %----- 33
35 T = TT; d = dd; x = xx; Z = Z1; 34
36 p = size(Z,2); 35
37 I = zeros((p+1),(p+1)); 36
38 U = 0; Ur = 0; 37
39 Ts = T.*d; 38
40 N = length(T); 39
41 K = sum(d); 40
42 eBz = exp((Z*(rmple'))'); 41
43 Xb = mean(x); 42
44 wr = zeros(N,1); Wr = zeros(N,1); 43
45 temp0 = 0; dN = 0; YseM = zeros(N,1); 44
46 %----- 45
47 % calculate the baseline cumulative hazards at each observed 43
48 % failure time for each treatment group by linear interpolation 44
49 %----- 45
```



```
1 L0 = interp1(c0,L0w, min(max(c0),T)); 1
2 L1 = interp1(c1,L1w, min(max(c1),T)); 2
3 %----- 3
4 % test statistic 4
5 %----- 5
6 for mm=1:N 6
7     if (Ts(mm)>0) 7
8         Y = (T>=Ts(mm)); 8
9         Y0 = Y.*(1-x); 9
10        Y1 = Y.*x; 10
11 %----- 11
12 %treatment-specific Survival functions of censoring 12
13 %----- 12
14 F0 = exp(-L0(mm)*exp(Z2*mple0')); 13
15 F1 = exp(-L1(mm)*exp(Z2*mple1')); 14
16 F0 = F0 + (F0==0).*eps; 15
17 F1 = F1 + (F1==0).*eps; 16
18 phi0r = ((min([F1';F0']))./F0)'; 17
19 philr = ((min([F1';F0']))./F1)'; 18
20 %----- 19
21 % uncorrected test 19
22 %----- 20
23 meBz = (eBz' .*Y)*ones(1,p+1); 21
24 s0 = (eBz*Y)/N; 22
25 s1 = (sum(meBz.*[x,Z]))/N; 23
26 s2 = ([x,Z]'*(meBz.*[x,Z]))/N; 24
27 vz = ((s2/s0) - ((s1/s0)'*(s1/s0))); 25
28 I = I + vz/N; 26
29 sc = ([x(mm),Z(mm,:)]) - (s1/s0); 27
30 U = U + sc(1); 28
31 %----- 29
32 % corrected test 29
33 %----- 30
34 Y0n=Y0.*phi0r; 31
35 Y1n=Y1.*philr; 32
36 Ys = Y0n + Y1n; 33
37 YseM = [YseM, (eBz' .*Ys)]; 34
38 meBz = (eBz' .*Ys)*ones(1,p+1); 35
39 s0 = (eBz*Ys)/N; 36
40 temp0 = [temp0, s0]; 37
41 s1 = (sum(meBz.*[x,Z]))/N; 38
42 s2 = ([x,Z]'*(meBz.*[x,Z]))/N; 39
43 E = s1/s0; 40
44 Ur = Ur + ( Ys(mm)*(x(mm)-E(1)) ); 41
45 dN = [dN, Ys(mm)]; 42
46 wr(mm) = Ys(mm)*( x(mm) - Xb ); 43
47 end 44
48 end 45
49 temp0(1)=[]; 46
```

```
1      dN(1)=[];
2      YseM(:,1)=[];
3      I = N*I;
4      %-----
5      %calculate sample version of the iid terms (A)
6      %-----
7      resr=sum(((x-Xb)*ones(1,K)).*(((YseM)./(ones(N,1)*temp0)).*
8      (ones(N,1)*(dN/N))))');
9      Wr = wr - resr';
10     %-----
11     %model-based variance estimate of uncorrected test
12     %-----
13     aa=I((2:p+1),(2:p+1));
14     iiI=inv(aa);
15     % iiI=aa\eye(size(aa)); may be more efficient
16     V = I(1,1)-(I(1,(2:p+1))*iiI*I((2:p+1),1));
17     %-----
18     %variance estimate of corrected test
19     %-----
20     Rrm = sum( (Wr-mean(Wr)).^(2) );
21     un = [U, U/sqrt(V)];
22     cor = [Ur, Ur/sqrt(Rrm)];
23     %-----
24     % NOTE: to calculate log-rank, set rmple=zeros(1,p) and
25     % V = I(1,1);
26     function [dL, ddL] = rPLgh(th)
27     % computes the gradient and Hessian of Cox's partial likelihood
28     % at *th*
29     % *th* is the (p+1) row vector of coefficients
30     % *T* is the column vector of N possibly right-censored event
31     % times
32     % *d* is the column vector of N indicators I(T<=C)
33     % *Z* is the N-by-p matrix of baseline covariates
34     global T d Z;
35     N = length(T);
36     p = size(Z,2);
37     I = zeros(p,p);
38     U = zeros(1,p);
39     %-----
40     % compute  $S^0(th,t)$ ,  $S^1(th,t)$  and  $S^2(th,t)$  at each event time
41     Bz = Z*(th');
42     eBz = exp(Bz');
43     Ts = T.*d;
44     for n=1:N
45         if (Ts(n)>0)
46             Y = (T>=Ts(n));
47             meBz = (eBz'.*Y)*ones(1,p);
48             s0 = (eBz*Y)/N;
49             s1 = (sum(meBz.*Z))/N;
```

```
1          s2 = (Z'*(meBz.*Z))/N; 1
2          vz = (s2/s0) - ((s1/s0)'*(s1/s0)); 2
3          sc = Z(n,:) - (s1/s0); 3
4          U = U + sc; 4
5          I = I + vz; 5
6          end 6
7          end 7
8          dL = U'; 8
9          ddL = -I; 9
10         function [LL, tt] = BRES(T,d,z,b) 10
11         % computes Breslow's estimate of baseline cumulative baseline 11
12         % hazard fn 11
13         % *T* is the column vector of N possibly right-censored event 12
14         % times 13
15         % *d* is the column vector of N indicators I(T<=C) 14
16         % assumes no ties in the data 15
17         % *z* is the N x p matrix of covariates 16
18         % *b* is the 1 x p vector of regression coefficients 17
19         Ts=T.*d; 18
20         Ts=Ts(Ts>0); 19
21         Ts=sort(Ts); 20
22         n=length(Ts); 21
23         L=1:n; 22
24         eb = exp(z*b'); 23
25         for mm=1:n 24
26             L(mm) = 1/sum( (T>=Ts(mm)).*eb ); 25
27         end 26
28         tt=[0,Ts']; 27
29         LL=[0,cumsum(L)]; 28
30 29
31 30
```

References

- 32 Breslow, N.E. (1972). Discussion of the paper by D.R. Cox. *J. Roy. Statist. Soc. B* **34**, 216–217. 32
- 33 Breslow, N.E. (1974). Covariance analysis of censored survival data. *Biometrics* **30**, 89–99. 33
- 34 Cox, D.R. (1972). Regression models and life-tables (with discussion). *J. Roy. Statist. Soc. B* **34**, 187–220. 34
- 35 Cox, D.R., Oakes, D.O. (1984). *Analysis of Survival Data*. Chapman and Hall, London. 35
- 36 DiRienzo, A.G., Lagakos, S.W. (2001a). Effects of model misspecification on tests of no randomized treat- 36
- 37 ment effect arising from Cox's proportional hazards model. *J. Roy. Statist. Soc. B* **63**, 745–757. 37
- 38 DiRienzo, A.G., Lagakos, S.W. (2001b). Bias correction for score tests arising from misspecified proportional 38
- 39 hazards regression models. *Biometrika* **88**, 421–434. 39
- 40 Gill, R.D. (1989). Non and semi-parametric maximum likelihood estimators and the von mises method 40
- 41 (Part 1). *Scand. J. Statist.* **16**, 97–128. 41
- 42 Kalbfleisch, J.D., Prentice, R.L. (1980). *The Statistical Analysis of Failure Time Data*. Wiley, New York. 42
- 43 Kaplan, E.L., Meier, P. (1958). Nonparametric estimation from incomplete observations. *J. Amer. Statist.* 43
- 44 *Assoc.* **53**, 457–481. 44
- 45 Klein, J.P., Moeschberger, M.L. (1997). *Survival Analysis – Techniques for Censored and Truncated Data*. 45
- Springer, New York.
- Kong, F.H., Slud, E. (1997). Robust covariate-adjusted log-rank tests. *Biometrika* **84**, 847–862.

The effects of misspecifying Cox's regression model

15

- 1 Lagakos, S.W. (1988). The loss in efficiency from misspecifying covariates in proportional hazards regression 1
2 models. *Biometrika* **75**, 156–160. 2
3 Lagakos, S.W., Schoenfeld, D.A. (1984). Properties of proportional-hazards score tests under misspecified 3
4 regression models. *Biometrics* **40**, 1037–1048. 4
5 Lin, D.Y., Wei, L.J. (1989). The robust inference for the Cox proportional hazards model. *J. Amer. Statist.* 5
6 *Assoc.* **84**, 1074–1078. 6
7 Lin, D.Y., Wei, L.J., Ying, Z. (1993). Checking the Cox model with cumulative sums of martingale-based 7
8 residuals. *Biometrika* **80**, 557–572. 8
9 McKeague, I.W., Utikal, K.J. (1990). Inference for a nonlinear counting process regression model. *Ann. Sta-* 9
10 *tist.* **18**, 1172–1187. 10
11 Morgan, T. (1986). Omitting covariates from the proportional hazards model. *Biometrics* **42**, 993–995. 11
12 Slud, E. (1991). Relative efficiency of the log rank test within a multiplicative intensity model. *Biometrika* **78**, 12
13 621–630. 13
14 Tsiatis, A.A., Rosner, G.L., Trichler, D.L. (1985). Group sequential tests with censored survival data adjusting 14
15 for covariates. *Biometrika* **72**, 365–373. 15
16 16
17 17
18 18
19 19
20 20
21 21
22 22
23 23
24 24
25 25
26 26
27 27
28 28
29 29
30 30
31 31
32 32
33 33
34 34
35 35
36 36
37 37
38 38
39 39
40 40
41 41
42 42
43 43
44 44
45 45



Collection of Biostatistics
Research Archive