



Published in final edited form as:

Stat Med. 2008 May 20; 27(11): 1911–1933. doi:10.1002/sim.3159.

Corrected score estimation in the proportional hazards model with misclassified discrete covariates

David M. Zucker^{1,*},† and Donna Spiegelman²

¹Department of Statistics, Hebrew University, Mount Scopus, 91905 Jerusalem, Israel

²Departments of Epidemiology and Biostatistics, Harvard School of Public Health, 677 Huntington Avenue, Boston, MA 02115, U.S.A.

SUMMARY

We consider Cox proportional hazards regression when the covariate vector includes error-prone discrete covariates along with error-free covariates, which may be discrete or continuous. The misclassification in the discrete error-prone covariates is allowed to be of any specified form. Building on the work of Nakamura and his colleagues, we present a corrected score method for this setting. The method can handle all three major study designs (internal validation design, external validation design, and replicate measures design), both functional and structural error models, and time-dependent covariates satisfying a certain ‘localized error’ condition. We derive the asymptotic properties of the method and indicate how to adjust the covariance matrix of the regression coefficient estimates to account for estimation of the misclassification matrix. We present the results of a finite-sample simulation study under Weibull survival with a single binary covariate having known misclassification rates. The performance of the method described here was similar to that of related methods we have examined in previous works. Specifically, our new estimator performed as well as or, in a few cases, better than the full Weibull maximum likelihood estimator. We also present simulation results for our method for the case where the misclassification probabilities are estimated from an external replicate measures study. Our method generally performed well in these simulations. The new estimator has a broader range of applicability than many other estimators proposed in the literature, including those described in our own earlier work, in that it can handle time-dependent covariates with an arbitrary misclassification structure. We illustrate the method on data from a study of the relationship between dietary calcium intake and distal colon cancer.

Keywords

errors in variables; nonlinear models; proportional hazards

1. INTRODUCTION

Many regression analyses involve explanatory variables that are measured with error. It is well known that failing to account for covariate error can lead to biased estimates of the regression coefficients. For linear models, theory for handling covariate error has been developed over the past 50 or more years; Fuller [1] provides an authoritative exposition. For nonlinear models, theory has been developing over the past 25 or so years. Carroll *et al.* [2] provide a comprehensive summary of the development to date; currently, the covariate error problem for nonlinear models remains an active research area. In particular, beginning with Prentice [3], a growing literature has developed on the Cox [4] proportional hazards survival regression model when some covariates are measured with error.

Three basic design setups are of interest. In all three designs, we have a main survival cohort for which surrogate covariate measurements and survival time data are available on all individuals. The three designs are as follows: (1) the internal validation design, where the true covariate values are available on a subset of the main survival cohort, (2) the external validation design, where the measurement error distribution is estimated from data outside the main survival study, and (3) the replicate measurements design, where replicate surrogate covariate measurements are available, either on a subset of the survival study cohort or on individuals outside the main survival study. Also, two types of models for the measurement error are of interest (see [1, p. 2; 2, Section 2]): structural models, where the true covariates are random variables, and functional models, where the true covariates are fixed values. Structural model methods generally involve estimation of some aspects of the distribution of the true covariate values; in functional model methods, this process is avoided.

We focus here on discrete covariates subject to misclassification, which are of interest in many epidemiological studies. In the case of a binary event outcome, there is extensive literature on the effects of misclassification on estimation of, and inference about, the relative risk and related parameters. Bross [5] is an early seminal paper. Detailed reviews have been given by Chen [6], Kuha and Skinner [7], and Walter and Irwig [8], while Kuha *et al.* [9] provide a concise summary of much of the development. Correction for misclassification entails estimating the classification probabilities through one of the designs listed in the preceding paragraph. Given appropriate estimates of the classification probabilities, consistent estimates of the relative risk and related inferences can proceed using ‘matrix methods’ [5, 10, 11] or maximum likelihood methods [12, 13].

The Cox survival regression model with covariate errors has been examined in a number of settings. Much of the existing work focuses on the independent additive error model, which assumes that the observed covariate value is equal to the true value plus a random error whose distribution is independent of the true value. In the case of discrete covariates subject to misclassification, this model practically never holds, and hence, the methods built upon it do not apply. Other methods are available in the literature that do apply to misclassification problems, but they are subject to substantial limitations, as we now describe.

The work of Zhou and Pepe [14] and of Zhou and Wang [15] deals with the internal validation design. Their approach involves empirically estimating the conditional mean of a certain function of the true covariate vector conditional on the observed covariate vector. This process entails stratification or smoothing with respect to the observed covariate vector. When the covariate vector is of moderate to high dimension, the ‘curse of dimensionality’ causes this approach to break down, even if only one of the covariates is error prone. Chen [16] presents an alternate method for the internal validation design. His method combines the regression coefficient estimate based on the validation sample only with information gleaned from the rest of the main study cohort. Chen’s approach assumes that it is possible to form a satisfactory initial estimate of the regression coefficient vector based on the validation sample alone. This is not the case, however, for studies where the event rate is low to moderate, the main study sample size is in the thousands, and the validation study sample size is in the low hundreds. Thus, in such cases, which often arise in practice, Chen’s approach breaks down. In addition, the methods of Zhou and Pepe, Zhou and Wang, and Chen do not cover the external validation or replicate measures setups.

Spiegelman *et al.* [34] and Wang *et al.* [18] discuss the simple and well-known regression calibration method, which applies to all three design setups. This method, however, is only an approximate method and does not yield a consistent estimator. Zucker and Spiegelman [19] and Zucker [20] present general methods suitable for all three designs, but their methods cover only time-independent covariates and their approaches do not seem generalizable to time-dependent covariates. Hu *et al.* [21] present a method that can handle time-dependent covariates under all three designs in a more general setting, but their approach is complex and its asymptotic properties were not formally examined. Most of the methods cited above apply only to structural models.

Another option is to apply the SIMEX approach, which is a general approach for covariate error problems [2, Chapter 5]. Recently, Küchenhoff *et al.* [22] developed a SIMEX method for the misclassification setting. Küchenhoff *et al.* treat generalized linear models; their method could be extended to survival models. The SIMEX approach, however, has some disadvantages. It requires multiple runs of the model-fitting process, and it relies on an extrapolation scheme that is uncertain and does not necessarily yield a consistent estimator.

Finally, it is possible in principle to apply methods that have been developed in the missing covariate literature, such as that of Herring and Ibrahim [23]. Most of this work deals with the internal validation setup, though the approach could be extendable to the other two design setups. These methods, however, are highly complex, and they cannot effectively handle the functional model setting or time-dependent covariates.

Thus, the currently existing methods are subject to substantial limitations, even for the internal validation design, and all the more for the external validation and replicate measures designs. There is a need for a new method that overcomes these limitations. In particular, there is a need for a convenient method for all three study designs that can handle general measurement error structures, both functional and structural models, and time-dependent covariates.

The aim of this paper is to present such a method for the case where the error-prone covariates are discrete. The misclassification may take any specified form desired, including an unstructured form. The error-free covariates are allowed to be either discrete or continuous. The time-dependent covariates are required to satisfy a certain ‘localized error’ condition, which we describe later. We present basic asymptotic properties of the method and examine its finite-sample performance in a simulation study.

In the case where the classification probabilities are estimated from replicate measurements data, it is necessary in this estimation process to regard the true covariate value as a random variable, as in a structural model (see the end of Section 4). However, when external replicate data are used, the marginal distribution of the true covariate need not be the same in the replicate sample as in the main study; we need only portability of the conditional distribution of the observed covariate given the true covariate.

Our proposed method follows the corrected score approach. Nakamura [24, 25] described the basic idea behind the approach. He then developed it by some detail under the independent additive error model. However, as noted above, this error model is not appropriate for discrete covariates. Accordingly, we build instead on the work of Akazawa *et al.* [26], which dealt with logistic regression with discrete covariates subject to misclassification. We extend the work of Akazawa *et al.* in a number of directions. First, we extend their approach from logistic regression to the Cox survival regression model; this extension involves substantial new technical development. Second, we allow for the case where, in addition to the error-prone discrete covariates, there may be a large number of other covariates measured without error; these other covariates can be either discrete or continuous. Finally, while Akazawa *et al.* assume the classification probabilities to be known, we allow them to be estimated, and we derive corrections to the estimated covariance matrix of the parameter estimates that reflect the error in estimating the classification probabilities. In the absence of misclassification, our method reduces to the classical Cox partial likelihood method.

The paper is organized as follows. Section 2 presents our proposed method for the case where the misclassification probabilities are assumed known. Section 3 discusses the case where the misclassification probabilities are estimated. Section 4 presents a simulation study of the method in the case of a single binary error-prone covariate. In the simulations, we consider both the case where the misclassification probabilities are known and the case where these probabilities are estimated from external replicate measurement data. Section 5 presents an application to data from the Nurses Health Study on the relationship between dietary calcium intake and distal colon cancer [27]. Section 6 provides a summary and discussion.

2. THE PROPOSED METHOD

2.1. The setup

We assume a standard survival analysis setup. We have observations on n independent individuals. We denote, for individual i , the survival time by T_i^0 and the time of right censoring by C_i . The observed survival data consist of the observed follow-up time

$T_i = \min(T_i^0, C_i)$ and the event indicator $\delta_i = I(T_i^0 \leq C_i)$. We let $Y_i(t) = I(T_i \geq t)$ denote the at-risk indicator. As usual, we assume that the covariate processes are left continuous with right limits, and that the failure process and the censoring process are conditionally independent given the covariate process in the sense described by Kalbfleisch and Prentice [28, Section 6.3.2]. Left truncation can be handled by setting $Y_i(t)$ to zero until the time at which individual i comes under observation.

The covariate structure is as follows. We denote the true covariate vector by $\mathbf{X}_i(t)$, and its dimension by p . We partition the vector $\mathbf{X}_i(t)$ into subvectors $\mathbf{W}_i(t)$ and $\mathbf{Z}_i(t)$, where $\mathbf{W}_i(t)$ is a p_1 -vector of error-prone covariates and $\mathbf{Z}_i(t)$ is a p_2 -vector of error-free covariates. We denote the observed value of $\mathbf{W}_i(t)$ by $\tilde{\mathbf{W}}_i(t)$. The vectors $\mathbf{W}_i(t)$ and $\tilde{\mathbf{W}}_i(t)$ are assumed to be discrete. The possible values of $\mathbf{W}_i(t)$ (each one a p_1 -vector) are denoted by $\mathbf{w}_1, \dots, \mathbf{w}_K$. The range of values of $\tilde{\mathbf{W}}_i(t)$ is assumed to be the same as that for $\mathbf{W}_i(t)$. For example, we could have a scalar binary covariate representing the presence or the absence of a given condition. Or the covariate might be the number of servings of a certain food that a person consumes per day. We denote by $k(i, t)$ the value of k such that $\tilde{\mathbf{W}}_i(t) = \mathbf{w}_k$. The vector $\mathbf{Z}_i(t)$ of error-free covariates is allowed to be either discrete or continuous. The case where the model involves interaction terms between the error-prone and error-free covariates can be accommodated with suitable minor notational changes.

We assume that the measurement error process is ‘localized’ in the sense that it depends only on the current true covariate value. More precisely, the assumption is that, conditional on the value of $\mathbf{X}_i(t)$, the value of $\tilde{\mathbf{W}}_i(t)$ is independent of the survival and censoring processes and of the values of $\mathbf{X}_i(s)$ for $s < t$. This assumption is plausible in many circumstances, such as situations in which the main source of error is technical or a laboratory error, or a reading/coding error, as with diagnostic X-rays and dietary intake assessments. The assumption will not be directly satisfied for covariates that represent cumulative exposure, though it may be possible to adapt our approach to cumulative exposure variables by working with the successive increments in observed exposure. For time-independent covariates, the assumption reduces to an assumption that the measurement error is independent of the survival and censoring processes. Under the localized error assumption, $Y_i(t)$ and $\tilde{\mathbf{W}}_i(t)$ are conditionally independent given $\mathbf{X}_i(t)$.

We denote $A_{kl}^{(i,t)} = Pr(\tilde{\mathbf{W}}_i(t) = \mathbf{w}_l | \mathbf{W}_i(t) = \mathbf{w}_k, \mathbf{Z}_i(t))$, which defines a square matrix $\mathbf{A}^{(i,t)}$ of classification probabilities. Note that the formulation here differs from that of Zucker and Spiegelman [19]. In addition, we allow here for the possibility that the classification probabilities depend on individual-specific factors, including the error-free covariates. Note also that the classification probabilities $A_{kl}^{(i,t)}$ are allowed to depend on t , either directly or through time-dependent individual-specific factors. Under this formulation, we can account for improvements in measurement techniques over time. In addition, if internal validation data are available, we can dispense with the localized error assumption. The assumption can be avoided by using classification rate estimates based only on the internal validation sample units that are still at risk at each given point in time. For now, we assume that $\mathbf{A}^{(i,t)}$ is known. In Section 3, we will consider the case where $\mathbf{A}^{(i,t)}$ is estimated.

We work under the Cox proportional hazards model, where the hazard function is of form

$$\lambda(t|\mathbf{x})=\lambda_0(t)e^{\boldsymbol{\beta}^T\mathbf{x}} \quad (1)$$

with $\lambda_0(t)$ being a baseline hazard function of unspecified form and $\boldsymbol{\beta}$ being a p -vector of unknown regression parameters that we wish to estimate. It is possible to extend the methodology to the case where $e^{\boldsymbol{\beta}^T\mathbf{x}}$ is replaced by a general relative risk function $\psi(\mathbf{x}; \boldsymbol{\beta})$, as in Thomas [29] and Breslow and Day [30, Section 5.1(c)]. This extension is described in a version of this paper available at the following website: <http://www.hsph.harvard.edu/faculty/spiegelman/manuscripts.html>

2.2. The key idea

The key idea behind our method is as follows. The Cox partial likelihood score function involves terms of the form $G(\mathbf{X}_i(t))=G(\mathbf{W}_i(t), \mathbf{Z}_i(t))$, where G is some function. Since $\mathbf{W}_i(t)$ is not directly observed, $G(\mathbf{W}_i(t), \mathbf{Z}_i(t))$ cannot be directly evaluated. Instead, we seek an observable function $G_i^*(\tilde{\mathbf{W}}_i(t), \mathbf{Z}_i(t))$ such that

$$E[G_i^*(\tilde{\mathbf{W}}_i(t), \mathbf{Z}_i(t))|\mathbf{W}_i(t), \mathbf{Z}_i(t)]=G(\mathbf{W}_i(t), \mathbf{Z}_i(t)) \quad (2)$$

In the case of discrete error-prone covariates, a G_i^* satisfying (2) may be constructed by a simple device we define

$$G_i^*(\tilde{\mathbf{W}}_i(t), \mathbf{Z}_i(t))=\sum_{l=1}^K B_{k(i,t)l}^{(i,t)} G(\mathbf{w}_l, \mathbf{Z}_i(t)) \quad (3)$$

where $k(i, t)$ is as defined previously and $B_{kl}^{(i,t)}$ is the (k, l) element of the matrix $\mathbf{B}^{(i,t)}=[\mathbf{A}^{(i,t)}]^{-1}$. We then have

$$\begin{aligned} E[G_i^*(\tilde{\mathbf{W}}_i(t), \mathbf{Z}_i(t))|\mathbf{W}_i(t)=\mathbf{w}_m, \mathbf{Z}_i(t)] & \\ &= \sum_{k=1}^K A_{mk}^{(i,t)} \sum_{l=1}^K B_{kl}^{(i,t)} G(\mathbf{w}_l, \mathbf{Z}_i(t)) \\ &= \sum_{l=1}^k \left(\sum_{k=1}^K A_{mk}^{(i,t)} B_{kl}^{(i,t)} \right) G(\mathbf{w}_l, \mathbf{Z}_i(t)) \\ &= G(\mathbf{w}_m, \mathbf{Z}_i(t)) \end{aligned} \quad (4)$$

so that (2) is indeed satisfied. Akazawa *et al.* [26] introduced this device for the case of logistic regression with covariate error.

2.3. The method

We now present our method. The classical Cox [4, 31] partial likelihood score function in the case with no measurement error is given by

$$U_r(\boldsymbol{\beta})=\frac{1}{n} \sum_{i=1}^n \delta_i \left(X_{ir}(T_i) - \frac{e_{1r}(T_i)}{e_0(T_i)} \right) \quad (5)$$

where

$$e_0(t) = \frac{1}{n} \sum_{j=1}^n Y_j(t) e^{\boldsymbol{\beta}^T \mathbf{X}_i(t)}$$

$$e_{1r}(t) = \frac{1}{n} \sum_{j=1}^n Y_j(t) X_{ir}(t) e^{\boldsymbol{\beta}^T \mathbf{X}_i(t)}$$

Now define

$$\xi_{ir}(t) = \sum_{l=1}^k B_{k(i,t)l}^{(i,t)} \mathbf{x}_r(\mathbf{w}_l, \mathbf{z}_i(t)) \quad (6)$$

$$\psi_i(t, \boldsymbol{\beta}) = \sum_{l=1}^k B_{k(i,t)l}^{(i,t)} \exp(\boldsymbol{\beta}^T \mathbf{x}(\mathbf{w}_l, \mathbf{z}_i(t))) \quad (7)$$

$$\eta_{ir}(t, \boldsymbol{\beta}) = \sum_{l=1}^K B_{k(i,t)l}^{(i,t)} \mathbf{x}_r(\mathbf{w}_l, \mathbf{Z}_i(t)) \exp(\boldsymbol{\beta}^T \mathbf{x}(\mathbf{w}_l, \mathbf{Z}_i(t))) \quad (8)$$

$$e_0^*(t) = \frac{1}{n} \sum_{j=1}^n Y_j(t) \psi_j(t, \boldsymbol{\beta}) \quad (9)$$

$$e_{1r}^*(t) = \frac{1}{n} \sum_{j=1}^n Y_j(t) \eta_{jr}(t, \boldsymbol{\beta}) \quad (10)$$

where $\mathbf{x}(\mathbf{w}, \mathbf{z})$ denotes the \mathbf{x} vector formed by the subvectors \mathbf{w} and \mathbf{z} . Our proposed corrected score function is then given by the following obvious analogue of (5):

$$U_r^*(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \delta_i \left(\xi_{ir}(T_i) \frac{e_{1r}^*(T_i)}{e_0^*(T_i)} \right) \quad (11)$$

The proposed corrected score estimator is the solution to $\mathbf{U}^*(\boldsymbol{\beta}) = \mathbf{0}$, where \mathbf{U}^* denotes the vector whose components are U_r^* .

We have

$$E[Y_i(t) \xi_{ir}(t) | \mathbf{X}_i(t)] = E[Y_i(t) E[\xi_{ir}(t) | \mathbf{X}_i(t), Y_i(t)] | \mathbf{X}_i(t)] = E[Y_i(t) E[\xi_{ir}(t) | \mathbf{X}_i(t)] | \mathbf{X}_i(t)] = E[Y_i(t) X_{ir}(t) | \mathbf{X}_i(t)] \quad (12)$$

where the second equality follows from the conditional independence of $Y_i(t)$ and $\tilde{\mathbf{W}}_i(t)$ given $\mathbf{X}_i(t)$, and the third equality follows from the argument of Section 2.2. Similarly,

$$E[Y_i(t)\psi_i(t, \boldsymbol{\beta})|\mathbf{X}_i(t)] = E[Y_i(t)e^{\boldsymbol{\beta}^T \mathbf{X}_i(t)}|\mathbf{X}_i(t)] \quad (13)$$

$$E[Y_i(t)\eta_{ir}(t, \boldsymbol{\beta})|\mathbf{X}_i(t)] = E[Y_i(t)X_{ir}(t)e^{\boldsymbol{\beta}^T \mathbf{X}_i(t)}|\mathbf{X}_i(t)] \quad (14)$$

Thus, referring to the quantity in parentheses in (11), the first term and the numerator and denominator of the second term all have the correct expectation. Now, $\mathbf{U}^*(\boldsymbol{\beta})$ is not an exactly unbiased estimating function, because the expectation of a ratio is not equal to the ratio of the expectations. However, as indicated in Appendix A.1, it follows from the law of large numbers that $\mathbf{U}^*(\boldsymbol{\beta})$ is an asymptotically unbiased score function.

Accordingly, under standard conditions like those of Andersen and Gill [32], our corrected score estimator will be consistent and asymptotically normal. Appendix A.1 presents an outline of the asymptotic arguments. See Huang and Wang [33] for a related discussion in a similar context. Denoting the true value of $\boldsymbol{\beta}$ by $\boldsymbol{\beta}_0$, the asymptotic covariance matrix of $n^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$ may be estimated by the sandwich formula

$$\hat{\mathbf{V}} = \mathbf{D}(\hat{\boldsymbol{\beta}})^{-1} \mathbf{H}(\hat{\boldsymbol{\beta}}) \mathbf{D}(\hat{\boldsymbol{\beta}})^{-1} \quad (15)$$

Here $\mathbf{D}(\boldsymbol{\beta})$ is -1 times the matrix of derivatives of $\mathbf{U}^*(\boldsymbol{\beta})$ with respect to the components of $\boldsymbol{\beta}$ and $\mathbf{H}(\boldsymbol{\beta})$ is an empirical estimate of the covariance matrix of $n^{1/2}\mathbf{U}^*(\boldsymbol{\beta})$. To define these matrices, some additional notations are needed. We define

$$\hat{\mathbf{Y}}_{ir}(\boldsymbol{\beta}) = \delta_i \left[\xi_{ir}(T_i) - \frac{e_{1r}^*(T_i)}{e_0^*(T_i)} \right] - \frac{1}{n} \sum_{j: T_j \leq T_i} \delta_j \left[\frac{\eta_{ir}(T_j, \boldsymbol{\beta})}{e_0^*(T_j)} - \frac{e_{1r}^*(T_i)}{e_0^*(T_i)} \frac{\psi_i(T_j, \boldsymbol{\beta})}{e_0^*(T_j)} \right] \quad (16)$$

In this definition, the first term tends to be the dominant term, especially if the event is rare. We further define

$$e_{2rs}^*(t) = \frac{1}{n} \sum_{j=1}^n Y_j(t) \sum_{l=1}^K B_{k(j)l}^{(i,t)} \mathbf{x}_r(\mathbf{w}_l, \mathbf{Z}_i(t)) \mathbf{x}_s(\mathbf{w}_l, \mathbf{Z}_i(t)) \exp(\mathbf{x}(\mathbf{w}_l, \mathbf{Z}_i(t))^T \boldsymbol{\beta})$$

With these definitions, we have

$$H_{rs}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{Y}}_{ir}(\boldsymbol{\beta}) \hat{\mathbf{Y}}_{is}(\boldsymbol{\beta}) \quad (17)$$

$$D_{rs}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \delta_i \left[\frac{e_{2rs}^*(T_i)}{e_0^*(T_i)} - \left(\frac{e_{1r}^*(T_i)}{e_0^*(T_i)} \right) \left(\frac{e_{1s}^*(T_i)}{e_0^*(T_i)} \right) \right] \quad (18)$$

The expression for $D_{rs}(\boldsymbol{\beta})$ is derived by straightforward differentiation. The derivation of the expression for $H_{rs}(\boldsymbol{\beta})$ is given in Appendix A.1.

In the internal validation design, for each individual i in the internal validation sample we can carry out the estimation with $\tilde{\mathbf{W}}_i$ replaced by \mathbf{W}_i and $\mathbf{A}^{(i,t)}$ replaced by the identity matrix. Alternatively, we can employ the hybrid scheme of Zucker and Spiegelman [19, Section 5], where a separate estimator of $\boldsymbol{\beta}$ is computed for the validation sample and for the main study excluding the validation sample, and the two estimators are then combined. The hybrid scheme is likely to be more efficient when the validation sample is sizable.

The case where there are replicate measurements, $\tilde{\mathbf{W}}_{ij}$ of $\tilde{\mathbf{W}}$ on the individuals in the main study can be handled in various ways. A simple approach is to redefine the quantities given in (6)–(8) by replacing $B_{k(i,t)}$ with the mean of $B_{k(i,j,t)l}$ over the replicates for individual i , with $k(i, j, t)$ defined as the value of k such that $\tilde{\mathbf{W}}_{ij}(t) = \mathbf{w}_k$. The development then proceeds as before.

2.4. Estimation of the cumulative hazard

The cumulative hazard $\Lambda_0(t) = \int_0^t \lambda_0(u) du$ can be estimated using the Breslow-type estimator

$$\hat{\Lambda}_0(t) = \sum_{i=1}^n \frac{\delta_i I(T_i \leq t)}{e_0^*(T_i, \hat{\boldsymbol{\beta}})} \quad (19)$$

As discussed at the end of Appendix A.1, the quantity $n^{1/2}(\hat{\Lambda}_0(t) - \Lambda_0(t))$, for a given t , is asymptotically mean-zero normal. In Appendix A.1, we present an estimator of the variance of this estimator.

3. ESTIMATED CLASSIFICATION PROBABILITIES

We now indicate the changes needed to handle the case where $A_{kl}^{(i,t)}$ are estimated. The relevant estimates may be obtained in several ways. In some cases, estimates are obtained from an external validation study, that is, a separate study with measurements of both \mathbf{W}_i and $\tilde{\mathbf{W}}_i$. Alternately, an internal validation design is used, with some individuals in the main survival study having measurements on both \mathbf{W}_i and $\tilde{\mathbf{W}}_i$. Another possibility is to base the estimates on internal or external replicate measures data; we discuss this in more detail at the end of the section. The theory developed in this section represents a step beyond the work of Akazawa *et al.* [26], who considered only the case where the classification probabilities are known. This theory is applied in the example presented in Section 5.

The main issue is how to adjust the covariance matrix of the estimates to account for the estimation error in $A_{kl}^{(i,t)}$. Following Zucker and Spiegelman [19], we express $\mathbf{A}^{(i,t)}$ as $\mathbf{A}^{(i,t)}(\boldsymbol{\omega})$ for some q' -vector of parameters $\boldsymbol{\omega}$. The nature of the function $\mathbf{A}^{(i,t)}(\boldsymbol{\omega})$ is dictated by the measurement error model employed. As an illustration, consider the simplest case: a single binary covariate with a common classification matrix $\mathbf{A}(\boldsymbol{\omega})$ for all individuals. In this case, with the false-positive and false-negative rates allowed to be different, $\mathbf{A}(\boldsymbol{\omega})$ takes the following form (where we assume that the sum of the off-diagonal elements is less than 1):

$$\mathbf{A}(\boldsymbol{\omega}) = \begin{bmatrix} \omega_1 & (1 - \omega_1) \\ (1 - \omega_2) & \omega_2 \end{bmatrix} \quad (20)$$

We presume that the parameter vector $\boldsymbol{\omega}$ is estimated from a study of one of the types described above, with m independent units. We presume further that the study yields an estimator $\hat{\boldsymbol{\omega}}$ having an approximate normal distribution with mean $\boldsymbol{\omega}$ and covariance matrix $m^{-1}\boldsymbol{\Gamma}$, along with an estimator $\hat{\boldsymbol{\Gamma}}$ of the matrix $\boldsymbol{\Gamma}$. This setup is a typical one in practical applications. For example, for the case of a single 0–1 binary covariate with internal or external validation data, the estimates of $\omega_k = \Pr(\tilde{W} = k-1 | W = k-1)$, $k = 1, 2$, are given by the obvious sample proportions, and $\boldsymbol{\Gamma}$ is a 2×2 diagonal matrix with $\Gamma_{kk} = \omega_k(1 - \omega_k)/\vartheta_k$, where ϑ_k is the fraction of individuals with $W = k-1$ in the validation study. The procedure for a replicate measures study is discussed at the end of this section. For the asymptotics we assume that m and n are of the same order of magnitude, i.e. $m/n \rightarrow \zeta$ for some constant ζ as $n \rightarrow \infty$. Otherwise, the error in $\mathbf{A}^{(i,t)}(\boldsymbol{\omega})$ will either be dominated by or will dominate the error in $\hat{\boldsymbol{\beta}}$ due to the variation in the survival data. Typically, ζ will be between 0 and 1.

Let us now write the corrected score function as $\mathbf{U}^*(\boldsymbol{\beta}, \boldsymbol{\omega})$ to indicate explicitly the dependence on $\boldsymbol{\omega}$. Also, let us denote the true value of $\boldsymbol{\beta}$ by $\boldsymbol{\beta}_0$ (as before) and the true value of $\boldsymbol{\omega}$ by $\boldsymbol{\omega}_0$. Since we are now estimating $\boldsymbol{\omega}_0$ by $\hat{\boldsymbol{\omega}}$, our estimating equation for $\boldsymbol{\beta}$ is now $\mathbf{U}^*(\boldsymbol{\beta}, \hat{\boldsymbol{\omega}}) = \mathbf{0}$. Using Taylor's theorem, we can write

$$\mathbf{0} = \mathbf{U}^*(\boldsymbol{\beta}, \hat{\boldsymbol{\omega}}) \doteq \mathbf{U}^*(\boldsymbol{\beta}_0, \boldsymbol{\omega}_0) - \mathbf{D}(\boldsymbol{\beta}_0)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + \dot{\mathbf{U}}^*(\boldsymbol{\beta}_0, \boldsymbol{\omega}_0)(\hat{\boldsymbol{\omega}} - \boldsymbol{\omega}_0)$$

where $-D_{rs}$ is the partial derivative of $U_r^*(\boldsymbol{\beta}, \boldsymbol{\omega})$ with respect to β_s evaluated at $\boldsymbol{\omega}_0$, and $\dot{\mathbf{U}}(\boldsymbol{\beta}, \boldsymbol{\omega})$ is a matrix whose (r, v) element is the partial derivative of $U_r^*(\boldsymbol{\beta}, \boldsymbol{\omega})$ with respect to ω_v . Hence, we have

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 \doteq \mathbf{D}(\boldsymbol{\beta}_0)^{-1} [\mathbf{U}^*(\boldsymbol{\beta}_0, \boldsymbol{\omega}_0) + \dot{\mathbf{U}}^*(\boldsymbol{\beta}_0, \boldsymbol{\omega}_0)(\hat{\boldsymbol{\omega}} - \boldsymbol{\omega}_0)] \quad (21)$$

If $\boldsymbol{\omega}$ is estimated from external data, then $\mathbf{U}^*(\boldsymbol{\beta}_0, \boldsymbol{\omega}_0)$ and $\hat{\boldsymbol{\omega}}$ are obviously independent. When $\boldsymbol{\omega}$ is estimated from an internal validation sample, the following argument can be put forth to show that these two quantities are asymptotically independent. The contribution to $\mathbf{U}^*(\boldsymbol{\beta}_0, \boldsymbol{\omega}_0)$ made by each individual in the study has asymptotically expectation zero conditional on the true covariate values. This is true in particular of the individuals in the internal validation sample. It hence follows from an iterated expectation argument that $\mathbf{U}^*(\boldsymbol{\beta}_0, \boldsymbol{\omega}_0)$ and $\hat{\boldsymbol{\omega}}$ are asymptotically uncorrelated (cf. [17, Appendix A]). Since we have asymptotic normality as well, this implies asymptotic independence. As a result, by Slutsky's theorem, the two terms in brackets in (21) are asymptotically independent, since $\dot{\mathbf{U}}^*(\boldsymbol{\beta}_0, \boldsymbol{\omega}_0)$ converges to a deterministic limit.

It therefore follows that, when external data or an internal validation sample is used to estimate $\boldsymbol{\omega}$, the necessary covariance adjustment may be accomplished by replacing the matrix $\mathbf{H}(\boldsymbol{\beta})$ with the following corrected version:

$$\mathbf{H}^{(\text{corr})}(\boldsymbol{\beta}, \boldsymbol{\omega}) = \mathbf{H}(\boldsymbol{\beta}, \boldsymbol{\omega}) + \zeta^{-1} \dot{\mathbf{U}}^*(\boldsymbol{\beta}, \boldsymbol{\omega}) \hat{\boldsymbol{\Gamma}} \dot{\mathbf{U}}^*(\boldsymbol{\beta}, \boldsymbol{\omega})^T \quad (22)$$

To present the formulas for $\dot{U}_{rv}^*(\boldsymbol{\beta})$, we define $\dot{\mathbf{A}}_{\nu}^{(i,t)}$ and $\dot{\mathbf{B}}_{\nu}^{(i,t)}$ to be, respectively, the partial derivative of the matrices $\mathbf{A}^{(i,t)}$ and $\mathbf{B}^{(i,t)}$ with respect to ω_{ν} . By the rule for differentiating an inverse matrix, we have $\dot{\mathbf{B}}_{\nu}^{(i,t)} = -\mathbf{B}^{(i,t)} \dot{\mathbf{A}}_{\nu}^{(i,t)} \mathbf{B}^{(i,t)}$. We then obtain

$$\dot{U}_{rv}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \delta_i \left[\dot{\xi}_{ir:v}(T_i, \boldsymbol{\beta}) - \frac{\dot{e}_{1r:v}^*(T_i)}{e_{0}^*(T_i)} + \left(\frac{e_{1r}^*(T_i)}{e_{0}^*(T_i)} \right) \left(\frac{\dot{e}_{0:v}^*(T_i)}{e_{0}^*(T_i)} \right) \right] \quad (23)$$

where $\dot{\xi}_{ir:v}$, $\dot{e}_{0:v}^*$, and $\dot{e}_{1r:v}^*$ are defined analogously to ξ_{ir} , e_0^* , and e_{1r}^* with $B_{k(i,t)l}^{(i,t)}$ replaced by $\dot{B}_{k(i,t)l:\nu}^{(i,t)}$.

When $\boldsymbol{\omega}$ is estimated from an internal replicate measures sample, $\mathbf{U}^*(\boldsymbol{\beta}_0, \boldsymbol{\omega}_0)$ and $\hat{\boldsymbol{\omega}}$ are no longer asymptotically independent, and hence we must work out the covariance between them. A typical scenario is the case of an i.i.d. setup where $\boldsymbol{\omega}$ is estimated by maximum likelihood. In Appendix A.2, we present the appropriate corrected version of \mathbf{H} for this setting.

The estimate for $\text{var}(\hat{\Lambda}_0(t))$ can be corrected in a similar manner. Appendix A.3 provides the details.

We now discuss in more depth the estimation of classification probabilities from a replicate measures study. There is a large literature on this problem (and extended versions thereof), especially for the case of a binary error-prone covariate, where the problem is generally known as the problem of evaluating a diagnostic marker without a gold standard. See [35] for a review of this area. The simplest case is that of a single binary (0–1) error-prone covariate W , with subject i in the replicate measures study having R_i replicate measurements \tilde{W}_{ij} on the surrogate measure \tilde{W} , where the replicates are conditionally i.i.d. given W . This case is relevant to many risk factors of interest in epidemiological studies, such as high blood pressure and high estradiol level, where the risk factor is assessed using direct

physiological measurements. In this setting, the subject totals $\tilde{W}_i^{(\text{tot})} = \sum_j \tilde{W}_{ij}$ are sufficient statistics. Denote $\alpha_1 = \Pr(\tilde{W} = 1 | W = 0)$ and $\alpha_2 = \Pr(\tilde{W} = 1 | W = 1)$. Then, conditional on $W = k-1$ ($k = 1, 2$), $\tilde{W}_i^{(\text{tot})}$ has a $\text{Bin}(R_i, \alpha_k)$ distribution. Defining π to be $\Pr(W = 1)$ within the replicate measures sample, the marginal distribution of $\tilde{W}_i^{(\text{tot})}$ is a mixture of the $\text{Bin}(R_i, \alpha_1)$ and $\text{Bin}(R_i, \alpha_2)$ distributions, with respective mixture probabilities $1-\pi$ and π :

$$\Pr(\tilde{W}_i^{(\text{tot})} = j) = (1-\pi) \binom{R_i}{j} \alpha_1^j (1-\alpha_1)^{R_i-j} + \pi \binom{R_i}{j} \alpha_2^j (1-\alpha_2)^{R_i-j} \quad (24)$$

The model is identifiable provided that some positive proportion of subjects have $R_i \geq 3$ and the correlation between W and \tilde{W} is positive. The latter condition is equivalent to the condition $\Pr(\tilde{W} = 1 | W = 0) + \Pr(\tilde{W} = 0 | W = 1) < 1$. The likelihood function can be expressed

directly from (24), and the parameters π , α_1 , and α_2 can then be estimated by maximum likelihood.

In more general settings, the basic ideas are similar, but of course the details are more complex. Several papers have discussed methods for using replicate data to estimate classification probabilities for more general settings, including parsimonious models for polychotomous W and models that allow the conditional distribution of the replicates given W to involve some dependence [8, 36–39]. The case of correlated replicates is of particular relevance to risk factors that are measured through self-report, such as dietary or physical activity variables.

4. SIMULATION STUDY

To investigate how our method performs, we carried out a simulation study in the setting of a single 0–1 binary covariate. The design of our simulation study was patterned after Zucker and Spiegelman [19, Section 6] and Zucker [20, Section 3.1]. The assumed study duration was 5 years. The baseline survival distribution was taken to be Weibull, with baseline hazard function $\lambda_0(t) = \alpha \mu(\mu t)^{\alpha-1}$. The power parameter α was taken equal to 5, which is typical of many types of cancer [30, Section 6.3; 40]. The scale parameter μ was chosen so as to yield a 25 per cent 5-year cumulative incidence rate for the unexposed population. Censoring was taken to be exponential with a rate of 1 per cent per year. For brevity of presentation, the false-positive rate $\Pr(\tilde{W} = 1|W = 0)$ and the false-negative rate $\Pr(\tilde{W} = 0|W = 1)$ were taken to be equal to a common classification error rate. A range of values was explored for the prevalence of the risk factor (5, 25, 40 per cent), the classification error rate (1, 5, 10, 20 per cent), and the true relative risk (1.5, 2.0). The number of simulation replications was 5000.

In our first simulation scenario, we assumed that the classification probabilities are known. We took the sample size to be 2000, leading to approximately 500 events in total. Table I shows the results. For comparison, we also present the simulation results given in [19] for the naive Cox partial likelihood estimator ignoring the measurement error, and for the parametric log relative risk estimator obtained by maximizing the full Weibull log likelihood under the relevant measurement error model.

The naive Cox estimator was typically badly biased except under 1 per cent misclassification with exposure prevalence of 25 or 40 per cent. By contrast, our method exhibited excellent performance, comparable to that of the fully parametric Weibull estimator. Under an exposure prevalence of 25 or 40 per cent, our method yielded nearly zero bias in the estimated log relative risk, nearly unbiased standard deviation estimates, and accurate confidence interval coverage. With an exposure prevalence of 5 per cent, the performance of all three estimators under consideration was degraded. This finding is not surprising, because the 5 per cent exposure situation presents two difficulties: (1) the expected number of events in the exposed group is only of the order of 25–50, (2) with a misclassification rate of 5 per cent or more, the predictive value of an observed positive exposure is low. The naive Cox estimator was drastically biased. Our estimator and the Weibull estimator were dramatically less biased, but still exhibited some bias. This bias was

due in part to outlying values; for both our estimator and the Weibull estimator, the deviation between the median value of the estimates and the true log relative risk was noticeably lower than the deviation between the mean estimated value and the true value. Overall, in terms of mean square error, the performance of our estimator was found to be nearly identical to that of Weibull estimator. In a few cases, our estimator was better; this reflects the fact that, for a given finite sample size, the asymptotically optimal parametric maximum likelihood expectation can be outperformed by an alternate estimator. The performance of the method proposed here essentially matches that of the methods of Zucker and Spiegelman [19] and Zucker [20], except that the method here was better for the problematic cases with 5 per cent exposure prevalence. The same pattern is seen under a sample size of 10 000 with a cumulative incidence rate of 5 per cent for the unexposed (results not shown).

In our second simulation scenario, we assumed that the classification probabilities are estimated from an external replicate measures study with 250 subjects and three replicate measurements per subject. The procedure for estimating the classification probabilities is described at the end of the preceding section. To get around low cell counts when the

misclassification rate is very small (viz. 0.01), we added $\frac{1}{2}$ to all the cell counts. We ran two sets of simulations, one for a sample size of 2000 (about 500 events in total) and the other for a sample size of 1000 (about 250 events in total). Table II presents the results.

Our method performed very well when the exposure prevalence was 25 or 40 per cent. Across the board, for both $n=2000$ and 1000, the bias in the estimated log relative risk was minimal, the standard deviation estimate was on target, and the confidence interval coverage was accurate. With $n=2000$, in most cases there was minimal change in the standard deviation of the log relative risk estimate due to the estimation of the misclassification rates, as compared with the standard deviation under known misclassification rates (shown in Table I). The one exception to this was the case of 20 per cent misclassification, where there was a 5–20 per cent increase in the standard deviation due to the estimation of the misclassification rates. The standard deviations for $n=1000$ were greater than those for $n=2000$ by about the expected factor of $\sqrt{2}$.

The method performed somewhat less well when the exposure prevalence was 5 per cent. The bias was higher in this situation, in some cases reaching the 10–20 per cent level. Still, this is much better than the bias of the naive Cox estimate (shown for $n=2000$ in Table I). The standard deviation estimates and confidence interval coverage were noticeably inaccurate in some cases. Also, under true misclassification rates of 20 per cent, the misclassification rates could not be successfully estimated in around 10–15 per cent of the simulation replications.

In summary, our method generally performed very well. Good performance was maintained even with estimated misclassification probabilities, except for some problems when the exposure prevalence rate was very low.

5. EXAMPLE

We illustrate our method on data from the Nurses Health Study concerning the relationship between dietary calcium (Ca) intake and cancer of the distal colon (i.e. the furthest segment of the large intestine) [27, Table 4]. The data consist of observations on female nurses whose calcium intake was assessed through a food frequency questionnaire (FFQ) in 1984 and were followed up to 31 May 1996 for distal colon cancer occurrence. Our analysis includes data on 60 575 nurses who reported in 1984 that they had never taken calcium supplements. In this cohort, there were 150 cases of distal colon cancer during the follow-up period. The analysis focuses on the effect of baseline calcium intake after adjustment for baseline body mass index (BMI) and baseline aspirin use. BMI is defined as the person's weight in kilograms divided by the square of the person's height in meters, and is a standard measure of a person's build (low BMI means thin, high BMI means fat). As in Wu *et al.*'s Table 4, we work with a binary 'high Ca' risk factor defined as 1 if the calcium intake was greater than 700 mg/day and 0 otherwise. Note that one glass of milk contains approximately 300 mg of calcium. BMI is expressed in terms of the following categories: <22, 22 to <25, 25 to <30, and 30kg/m² or greater. Aspirin use is coded as yes (1) or no (0). Thus, our model has five explanatory variables, one for the binary risk factor (W), three dummy variables for BMI (Z_1, Z_2, Z_3), and one for aspirin use (Z_4). BMI and aspirin use status are assumed to be measured without error.

It is well known that the FFQ measures dietary intake with some degree of error and more reliable information can be obtained from a diet record (DR) [41, Chapter 6]. We thus take W to be the Ca risk factor indicator based on the DR and \tilde{W} to be the Ca risk factor indicator based on the FFQ. The classification probabilities are estimated using data from the Nurse's Health Study validation study [41, pp. 122–126]. The estimated specificity was

$\hat{P}_r(\tilde{W}=0|W=0)=0.78$, with an estimated standard error of 0.042. The estimated sensitivity was $\hat{P}_r(\tilde{W}=1|W=1)=0.72$, with an estimated standard error of 0.046.

Table III presents the results of the following analyses: (1) a naive classical Cox regression analysis ignoring measurement error, corresponding to an assumption that there is no measurement error, (2) our method with \mathbf{A} assumed known and set according to the foregoing estimated classification probabilities, ignoring the estimation error in these probabilities, and (3) our method with \mathbf{A} estimated as above with the estimation error in the probabilities taken into account (main study/external validation study design). The last of these analyses makes use of the theory developed in Section 3.

The results followed the expected pattern. Adjusting for the misclassification in calcium intake had a marked effect on the estimated relative risk for high calcium intake. Accounting for the error in estimating the classification probabilities increased (modestly) the standard error of the log relative risk estimate. The relative risk estimates for high calcium intake and the corresponding 95 per cent confidence intervals obtained in the three analyses were as follows:

Method	Estimate	95 per cent CI
Naive Cox	0.71	[0.51, 0.99]
A Known	0.49	[0.24, 1.01]
A Estimated	0.49	[0.23, 1.04]

In general, in the multivariate setting, measurement error (including misclassification) can lead to either attenuation or magnification of covariate effects. In our example, the misclassification led to attenuation of the high calcium effect, so that the corrected relative risk was further from the null value of 1 than the naive relative risk. The misclassification correction had a small effect on the estimated regression coefficients for the BMI dummy variables and essentially no effect on the estimated regression coefficient for aspirin use.

6. SUMMARY AND DISCUSSION

We have considered the Cox [4] proportional hazards model with a set of covariates that includes error-prone discrete covariates along with error-free covariates, which may be discrete or continuous. The misclassification in the discrete error-prone covariates is allowed to be of any specified form. Building on the work of Nakamura [24, 25] and Akazawa *et al.* [26], we have developed an easily implemented corrected score method for this setting. The method can handle all three major study designs (internal validation design, external validation design, and replicate measures design), both functional and structural error models, and time-dependent covariates satisfying the ‘localized error’ condition described in Section 2. Also, for the internal validation design, the ‘localized error’ condition can be eliminated by using time-dependent classification rate estimates. The method thus represents a significant advance relative to other methods in the literature for this problem. The method performed well in a simulation study, both with misclassification probabilities known and with misclassification probabilities estimated from an external replicate measures study.

In most applications, the new method developed in this paper will be easier to apply than and preferable to our previous method based on weighted transformed Kaplan–Meier curves [19]. Our previous method requires defining strata for every possible configuration of the entire covariate vector (both the error-prone and error-free part). Except when the number of configurations is small, this leads to cumbersome implementation and loss of data for strata having no events. Our current method avoids this problem. Also, our previous method cannot handle continuous error-free covariates, while the current method can. Additionally, our previous method cannot handle time-dependent covariates (nor can the method of Zucker [20]), whereas our current method can handle such covariates if the ‘localized error’ condition applies or internal validation data are available. In some applications, our previous method might be preferred on account of reduced computational burden. Also, our previous method may be more convenient when it is desired to apply a measurement error correction based on published Kaplan–Meier curves for various risk groups as presented in medical and other subject-matter journals. This point is of particular relevance to meta-analysis applications.

This work focuses on the case where the error-prone covariates are discrete. This is admittedly a limitation. However, much of the existing work on the Cox model with covariate error focuses on continuous covariates with independent additive error, and as such does not apply or generalize easily to discrete covariates with misclassification. In many epidemiological studies, the error-prone covariates of interest are in fact discrete. Thus, the method presented here fills a definite need.

Still, there are cases where it is of interest to investigate continuous error-prone risk factors. In the case of a single error-prone continuous risk factor W , the basic equation (2) for classical likelihood models takes the form

$$G(w) = \int a(\tilde{w}|w) G^*(\tilde{w}) d\tilde{w} \quad (25)$$

where $a(\tilde{w}|w)$ is the conditional density of \tilde{W} given W , the integral is over the entire range of \tilde{W} , and the arguments \mathbf{Z}_i and $\boldsymbol{\beta}$ are suppressed. Analogous equations are obtained for the case of multiple error-prone continuous covariates. Equation (25) is a Fredholm integral equation of the first kind. Such equations are discussed, for example, in Delves and Mohamed [42, Chapter 12], where numerical solution techniques are discussed. These techniques could be applied to the measurement error problem in suitable cases. However, as Delves and Mohamed indicate, such equations can sometimes be ill conditioned and do not always have a solution. Thus, for example, for the logistic regression model with additive normal covariate error, Stefanski [43] showed that a corrected score function satisfying (2) does not exist.

One way around the problem is to carry out a mild discretization of the error-prone covariate, fine enough to reduce the bias satisfactorily but not so fine as to lead to numerical problems. This approach will not produce a strictly consistent estimator, but it is reasonable to expect that the bias will be small. This supposition is supported by Cochran's [44] classic work on subclassification, which indicates that the bulk of the information in a continuous variable can often be captured in a discretized version with four to six categories. We are currently exploring this discretization approach in more depth.

Alternatively, an attempt can be made to modify the corrected score approach so that it will work for the model under consideration. Huang and Wang [45] developed such a modification for logistic regression. In their work, the terms in the likelihood score function were re-weighted to yield a new likelihood score function for which a corrected score satisfying (2) can be derived. These authors dealt only with independent additive error, which for the Cox model is already covered by existing corrected score methods [33, 46, 47]. Modification of the corrected score approach under other measurement error structures is an open problem.

Acknowledgments

This work was supported in part by a grant from the U.S. National Cancer Institute. We thank Els Goetghebeur and Malka Gorfine for their helpful comments and Ruifeng Li for assistance in the data analysis for the example. We also thank the Associate Editor and referees for comments that led to a substantially improved presentation.

APPENDIX A: TECHNICAL DETAILS

A.1. Outline proof of consistency and asymptotic normality

We give here an outline derivation of the asymptotic properties of our estimator. Our goal is to indicate the main steps of the argument without dwelling on the technical details. For simplicity, we focus on the i.i.d. structural model setting; our development will go through for the functional model setting as well provided that $\{\mathbf{X}_i(t) : i = 1, 2, \dots\}$ exhibits ergodic behavior suitably similar to that of an i.i.d. sequence. We assume the parameter space is a compact set \mathcal{B} of $\boldsymbol{\beta}$ values of which the true value $\boldsymbol{\beta}_0$ is an interior point. We further assume that regularity conditions along the lines of Andersen and Gill [32] (AG) are in force over \mathcal{B} .

We take up first the issue of consistency. We are operating under AG's 'asymptotic stability' conditions, which in the i.i.d. case follow from the functional law of large numbers in Andersen and Gill's Appendix III. Define

$$s_0(t, \boldsymbol{\beta}) = E[Y_i(t) \exp(\boldsymbol{\beta}^T \mathbf{X}_i(t))], \quad s_{1r}(t, \boldsymbol{\beta}) = E[Y_i(t) X_{ir}(t) \exp(\boldsymbol{\beta}^T \mathbf{X}_i(t))]$$

Using AG's arguments, including appeal to the asymptotic stability conditions, we find that $U_r(\boldsymbol{\beta})$ converges uniformly over \mathcal{B} to

$$u_r(\boldsymbol{\beta}) = \int \left[s_{1r}(t, \boldsymbol{\beta}_0) - \left(\frac{s_{1r}(t, \boldsymbol{\beta})}{s_0(t, \boldsymbol{\beta})} \right) s_0(t, \boldsymbol{\beta}_0) \right] \lambda_0(t) dt$$

In view of equations (12)–(14) of our Section 2.3, the same arguments yield the result that $U_r^*(\boldsymbol{\beta})$ converges uniformly over \mathcal{B} to $u_r(\boldsymbol{\beta})$. We thus have the following:

1. The function $\mathbf{U}^*(\boldsymbol{\beta})$, being continuous over \mathcal{B} , is therefore uniformly continuous over \mathcal{B} .
2. The function $\mathbf{U}^*(\boldsymbol{\beta})$ converges uniformly to $\mathbf{u}(\boldsymbol{\beta})$ over \mathcal{B} .
3. As can be seen by inspection, $\mathbf{u}(\boldsymbol{\beta}_0) = \mathbf{0}$.

Moreover, as AG show, $\boldsymbol{\beta}_0$ is the *only* zero point of $\mathbf{u}(\boldsymbol{\beta})$ in \mathcal{B} . As a result, given the compactness of the parameter space, convergence of $\hat{\boldsymbol{\beta}}$ to $\boldsymbol{\beta}_0$ follows by standard subsequence arguments.

We now discuss the asymptotic normality of $\hat{\boldsymbol{\beta}}$. By Taylor expansion we may write

$$n^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) = \mathbf{D}(\tilde{\boldsymbol{\beta}}) \mathbf{U}(\boldsymbol{\beta}_0)$$

where $\tilde{\boldsymbol{\beta}}$ lies between $\boldsymbol{\beta}_0$ and $\hat{\boldsymbol{\beta}}$, and hence converges to $\boldsymbol{\beta}_0$. Hence, as in AG, $\mathbf{D}(\tilde{\boldsymbol{\beta}})$ converges to the limiting value of $\mathbf{D}(\boldsymbol{\beta}_0)$, which exists by virtue of the asymptotic stability conditions. It now remains only to show that $n^{1/2} \mathbf{U}(\boldsymbol{\beta}_0)$ is asymptotically normal.

We first recall expression (11) for the corrected score function:

$$U_r^*(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \delta_i \left(\xi_{ir}(T_i) - \frac{e_{1r}^*(T_i)}{e_0^*(T_i)} \right)$$

Since $U_r^*(\boldsymbol{\beta}_0)$ does not have exactly expectation zero, the martingale approach of AG cannot be applied to derive the asymptotic distribution of $U^*(\boldsymbol{\beta}_0)$. Instead, we follow the approach of Lin and Wei [48]. From this point forward, all quantities involving $\boldsymbol{\beta}$ (including those in which the dependence is suppressed from the notation) are evaluated at the true value $\boldsymbol{\beta}_0$, except in (A1), which presents a definition for general $\boldsymbol{\beta}$. We use counting process notation, based on the definition $N_i(t) = I(T_i \leq t, \delta_i = 1)$. We define

$$\bar{N}(t) = \frac{1}{n} \sum_{i=1}^n N_i(t), \mathcal{N}(t) = E[N_i(t)]$$

We have, from the law of large numbers along with (12) and (13), that $\bar{N}(t) \rightarrow \mathcal{N}(t)$, $e_0^*(t) \rightarrow s_0(t)$, and $e_{1r}^*(t) \rightarrow s_{1r}(t)$ as $n \rightarrow \infty$. Here, the dependence on $\boldsymbol{\beta}$ is suppressed from the notation, and, as mentioned above, evaluation is at $\boldsymbol{\beta} = \boldsymbol{\beta}_0$. As seen from Andersen and Gill [32, Appendix III], the convergence is uniform in t . In the development below, the symbol \doteq will denote equality up to negligible terms.

We can write

$$\begin{aligned} U_r^*(\boldsymbol{\beta}_0) &= \frac{1}{n} \sum_{i=1}^n \int \xi_{ir}(t) dN_i(t) \\ &\quad - \int \frac{e_{1r}^*(t)}{e_0^*(t)} d\bar{N}(t) \\ &= \frac{1}{n} \sum_{i=1}^n \int \xi_{ir}(t) dN_i(t) \\ &\quad - \int \frac{e_{1r}^*(t)}{e_0^*(t)} d\mathcal{N}(t) \\ &\quad - \int \frac{s_{1r}(t)}{s_0(t)} d(\bar{N} - \mathcal{N})(t) \\ &\quad - \int \left[\frac{e_{1r}^*(t)}{e_0^*(t)} - \frac{s_{1r}(t)}{s_0(t)} \right] d(\bar{N} - \mathcal{N})(t) \\ &\quad - \mathcal{N}(t) \doteq \frac{1}{n} \sum_{i=1}^n \int \xi_{ir}(t) dN_i(t) \\ &\quad - \int \frac{e_{1r}^*(t)}{e_0^*(t)} d\mathcal{N}(t) \\ &\quad - \int \frac{s_{1r}(t)}{s_0(t)} d(\bar{N} - \mathcal{N})(t) \end{aligned}$$

In addition,

$$\frac{e_{1r}^*(t)}{e_0^*(t)} = \frac{e_{1r}^*(t)}{s_0(t)} + e_{1r}^*(t) \left[\frac{1}{e_0^*(t)} - \frac{1}{s_0(t)} \right] = \frac{e_{1r}^*(t)}{s_0(t)} - \left[\frac{e_{1r}^*(t)}{e_0^*(t)s_0(t)} \right] (e_0^*(t) - s_0(t)) \doteq \frac{1}{s_0(t)} \left[e_{1r}^*(t) - \frac{s_{1r}(t)}{s_0(t)} e_0^*(t) - s_0(t) \right]$$

Thus, substituting and re-arranging, we obtain

$$\begin{aligned}
 U_r^*(\boldsymbol{\beta}_0) &\doteq \frac{1}{n} \sum_{i=1}^n \int \xi_{ir}(t) dN_i(t) \\
 &\quad - \int \frac{s_{1r}(t)}{s_0(t)} d(\overline{N} \\
 &\quad - \mathcal{N})(t) \\
 &\quad - \int \frac{1}{s_0(t)} \left[e_{1r}^*(t) - \frac{s_{1r}(t)}{s_0(t)} (e_0^*(t) - s_0(t)) \right] d\mathcal{N}(t) \\
 &= \frac{1}{N} \sum_{i=1}^n \int (\xi_{ir}(t) \\
 &\quad - \frac{s_{1r}(t)}{s_0(t)}) dN_i(t) \\
 &\quad - \int \left(\frac{e_{1r}^*(t)}{s_0(t)} \right. \\
 &\quad \left. - \frac{s_{1r}(t)}{s_0(t)} \frac{e_0^*(t)}{s_0(t)} \right) d\mathcal{N}(t) \\
 &= \frac{1}{N} \sum_{i=1}^n \int (\xi_{ir}(t) \\
 &\quad - \frac{s_{1r}(t)}{s_0(t)}) dN_i(t) \\
 &\quad - \frac{1}{n} \sum_{i=1}^n \int \left(\frac{Y_i(t) \eta_{ir}(t, \boldsymbol{\beta}_0)}{s_0(t)} \right. \\
 &\quad \left. - \frac{s_{1r}(t)}{s_0(t)} \frac{Y_i(t) \psi_i(t, \boldsymbol{\beta}_0)}{s_0(t)} \right) d\mathcal{N}(t) \\
 &= \frac{1}{n} \sum_{i=1}^n \Upsilon_{ir}(\boldsymbol{\beta}_0)
 \end{aligned}$$

where

$$\Upsilon_{ir}(\boldsymbol{\beta}) = \delta_i \left[\xi_{ir}(T_i) \frac{e_{1r}^*(T_i)}{e_0^*(T_i)} \right] - \int Y_i(t) \left[\frac{\eta_{ir}(t, \boldsymbol{\beta})}{s_0(t)} - \frac{s_{1r}(t)}{s_0(t)} \frac{\psi_i(t, \boldsymbol{\beta})}{s_0(t)} \right] d\mathcal{N}(t) \quad (\text{A1})$$

From this result it follows immediately from the classical central limit theorem for i.i.d. random vectors that $n^{1/2} \mathbf{U}^*(\boldsymbol{\beta}_0)$ is asymptotically mean-zero multivariate normal. It is straightforward to see that the asymptotic covariance matrix of $n^{1/2} \mathbf{U}^*(\boldsymbol{\beta}_0)$ can be estimated consistently by expression (17) evaluated at $\hat{\boldsymbol{\beta}}$.

Finally, we turn to the cumulative hazard estimator (19). This estimator can be expressed as

$$\hat{\Lambda}_0(t) = \int_0^t \frac{d\overline{N}(u)}{e_0^*(u, \hat{\boldsymbol{\beta}})}$$

Using arguments similar to the above, we find that

$$\hat{\Lambda}_0(t) - \Lambda_0(t) \doteq -\mathbf{a}^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) + \frac{1}{n} \sum_{i=1}^n \int_0^t \frac{dN_i(u)}{s_0(u, \boldsymbol{\beta}_0)} - \frac{1}{n} \sum_{i=1}^n \int_0^t \frac{Y_i(u) \psi_i(u, \boldsymbol{\beta}_0)}{s_0(u, \boldsymbol{\beta}_0)} d\mathcal{N}(u)$$

with the r th element of \mathbf{a} given by

$$a_r = \int_0^t \frac{s_{1r}(u, \boldsymbol{\beta}_0)}{s_0(u, \boldsymbol{\beta}_0)^2} d\mathcal{N}(u)$$

Using the approximation $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 \doteq \mathbf{D}(\boldsymbol{\beta}_0)^{-1} \mathbf{U}^*(\boldsymbol{\beta}_0)$, and defining $\mathbf{c}(\boldsymbol{\beta}) = -\mathbf{D}(\boldsymbol{\beta})^{-1} \mathbf{a}$, we obtain

$$\hat{\Lambda}_0(t) - \Lambda_0(t) \doteq \frac{1}{n} \sum_{i=1}^n \Upsilon_i^*(\boldsymbol{\beta}_0) \quad (\text{A2})$$

where

$$\Upsilon_i^*(\boldsymbol{\beta}) = \sum_{r=1}^p c_r(\boldsymbol{\beta}) \Upsilon_{ir}(\boldsymbol{\beta}) + \int_0^t s_0(u, \boldsymbol{\beta})^{-1} dN_i(u) - \int_0^t \frac{Y_i(u) \psi_i(u, \boldsymbol{\beta})}{s_0(u, \boldsymbol{\beta})} d\mathcal{N}(u)$$

As before, it is apparent from this representation that the estimator is asymptotically mean-zero normal with variance that can be estimated by

$$\widehat{\text{var}}(\hat{\Lambda}_0(t)) = \frac{1}{n} \sum_{i=1}^n \hat{\Upsilon}_i^*(\hat{\boldsymbol{\beta}})^2 \quad (\text{A3})$$

where

$$\hat{\Upsilon}_i^*(\boldsymbol{\beta}) = \sum_{r=1}^p \hat{c}_r(\boldsymbol{\beta}) \hat{\Upsilon}_{ir}(\boldsymbol{\beta}) + \frac{\delta_i I(T_i \leq t)}{e_0^*(T_i, \boldsymbol{\beta})} - \frac{1}{n} \sum_{j=1}^n \delta_j I(T_j \leq \min\{T_i, t\}) \frac{\psi_i(T_j, \boldsymbol{\beta})}{e_0^*(T_j, \boldsymbol{\beta})}$$

with $\hat{\mathbf{c}}(\boldsymbol{\beta}) = -\mathbf{D}(\boldsymbol{\beta})^{-1} \hat{\mathbf{a}}(\boldsymbol{\beta})$ and

$$\hat{a}_r(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \delta_i I(T_i \leq t) \frac{e_{1r}^*(T_i, \boldsymbol{\beta})}{e_0^*(T_i, \boldsymbol{\beta})^2}$$

A.2. Correction to $\widehat{\text{var}}(\hat{\boldsymbol{\beta}})$ with $\mathbf{A}^{(i,t)}$ estimated from internal replicate data

We work in the setting where the replicate data are i.i.d. across individuals, and $\boldsymbol{\omega}$ is estimated by maximum likelihood. Let R_i denote the number of replicates on individual i , and let $g_i(\boldsymbol{\omega})$ denote the log likelihood function for $(\tilde{\mathbf{W}}_{i1}, \dots, \tilde{\mathbf{W}}_{iR_i})$. The overall normalized log likelihood is then $g(\boldsymbol{\omega}) = m^{-1} \sum_{i \in \mathcal{R}} g_i(\boldsymbol{\omega})$, where \mathcal{R} denotes the set of individuals in the internal replicate measures sample. Let $\mathbf{g}'(\boldsymbol{\omega})$ and $\mathbf{g}''(\boldsymbol{\omega})$ denote the gradient vector and Hessian matrix, respectively, of $g(\boldsymbol{\omega})$, and let $\mathbf{g}'_i(\boldsymbol{\omega})$ denote the gradient

of $g_i(\omega)$. We can then express $\hat{\omega}$ in terms of the classic asymptotic approximation $\hat{\omega} \doteq -\mathbf{g}''(\omega_0)^{-1} \mathbf{g}'(\omega_0)$. Let $\mathbf{T}_i(\omega)$ denote the vector comprising the quantities $T_{ir}(\beta, \omega)$, as defined in (A1), and let $\hat{\mathbf{T}}_i(\hat{\beta}, \hat{\omega})$ be the vector of corresponding estimated values $T_{ir}(\hat{\beta}, \hat{\omega})$ as defined by (16). Define

$$\Phi = \text{COV} \left(\left[\frac{1}{\sqrt{m}} \sum_{i \in \mathfrak{R}} \mathbf{T}_i^*(\omega) \right], \sqrt{m} \mathbf{g}'(\omega_0) \right)$$

The limiting value of Φ can then be estimated empirically by

$$\hat{\Phi} = \frac{1}{m} \sum_{i \in \mathfrak{R}} \hat{\mathbf{T}}_i(\hat{\beta}, \hat{\omega}) \mathbf{g}'_i(\hat{\omega})^T \quad (\text{A4})$$

The appropriate corrected version of \mathbf{H} is then

$$\mathbf{H}^{(\text{corr})} = \mathbf{H} + \zeta^{-1} \dot{\mathbf{U}}^*(\hat{\beta}, \hat{\omega}) \hat{\Gamma} \dot{\mathbf{U}}^*(\hat{\beta}, \hat{\omega})^T - \hat{\Phi} \mathbf{g}''(\hat{\omega})^{-1} \dot{\mathbf{U}}^*(\hat{\beta}, \hat{\omega})^T \quad (\text{A5})$$

where for the present setup we have $\hat{\Gamma} = \mathbf{g}''(\hat{\omega})^{-1}$.

A.3. Correction to $\text{var}(\hat{\Lambda}_0(t))$ with $\mathbf{A}^{(i,t)}$ estimated

When, as in Section 3, we estimate the parameters ω that determine the classification probabilities, representation (A2) becomes

$$\hat{\Lambda}_0(t) - \Lambda_0(t) \doteq \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^*(\beta_0) + \mathbf{h}^T (\hat{\omega} - \omega_0)$$

where the v th element of \mathbf{h} is given by

$$h_v = [\dot{\mathbf{U}}^*(\beta, \omega)^T \mathbf{c}(\beta)]_v - \frac{1}{n} \sum_{i=1}^n \int_0^t \frac{Y_i(u)}{s_0(u, \beta)} \sum_{l=1}^K [\dot{B}_v^{(i,t)}]_{k(i,u)l} \exp(\beta^T \mathbf{x}(\mathbf{w}_l, \mathbf{Z}_i(u))) d\mathcal{N}(u)$$

The estimate (A3) for $\text{var}(\hat{\Lambda}_0(t))$ can be corrected as follows. Define $\hat{\mathbf{h}} = \hat{\mathbf{h}}(\beta, \omega)$ by

$$\hat{h}_v = [\dot{\mathbf{U}}^*(\beta, \omega)^T \hat{\mathbf{c}}(\beta)]_v - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \delta_i I(T_j \leq \min\{T_i, t\}) s_0(T_j, \beta)^{-1} \sum_{l=1}^K [\dot{B}_v^{(i,t)}]_{k(i,T_j)l} \exp(\beta^T \mathbf{X}(\mathbf{w}_l, \mathbf{Z}_i(T_j)))$$

With external data or an internal validation study, it is necessary merely to add to (A3) the term $\zeta^{-1} \hat{\mathbf{h}}^T \hat{\Gamma} \hat{\mathbf{h}}$. For the case of internal replicate data, the additional term is

$$\zeta^{-1} \hat{\mathbf{h}}^T \hat{\Gamma} \hat{\mathbf{h}} + \left[\frac{1}{m} \sum_{i \in \mathfrak{R}} \hat{\mathbf{T}}_i^*(\hat{\beta}) \mathbf{g}'_i(\hat{\omega}) \right]^T \mathbf{g}''(\hat{\omega})^{-1} \hat{\mathbf{h}}$$

REFERENCES

1. Fuller, WA. Measurement Error Models. New York: Wiley; 1987.
2. Carroll, R.J.; Ruppert, D.; Stefanski, L.A.; Crainiceanu, CM. Measurement Error in Nonlinear Models: A Modern Perspective. 2nd edn. Boca Raton, FL, London: Chapman & Hall/CRC; 2006.
3. Prentice R. Covariate measurement errors and parameter estimation in a failure time regression model. *Biometrika*. 1982; 69:331–342.
4. Cox DR. Regression models and life-tables (with Discussion). *Journal of the Royal Statistical Society, Series B*. 1972; 34:187–220.
5. Bross IDJ. Misclassification in 2×2 tables. *Biometrics*. 1954; 10:478–486.
6. Chen TT. A review of methods for misclassified categorical data in epidemiology. *Statistics in Medicine*. 1989; 8:1095–1106. [PubMed: 2678350]
7. Kuha, J.; Skinner, C. Categorical data analysis and misclassification. In: Lyberg, L.; Biemer, P.; Collins, M.; DeLeeuw, E.; Dippo, C.; Schwarz, N.; Trewin, D., editors. *Survey Measurement and Process Quality*. New York: Wiley; 1997. p. 633–670.
8. Walter SD, Irwig LM. Estimation of test error rates, disease prevalence, and relative risk from misclassified data: a review. *Journal of Clinical Epidemiology*. 1988; 41:923–937. [PubMed: 3054000]
9. Kuha, J.; Skinner, C.; Palmgren, J. Misclassification error. In: Armitage, P.; Colton, T., editors. *Encyclopedia of Biostatistics*. 1st edn. Vol. 1. New York: Wiley; 1998. p. 2615–2621.
10. Marshall JR. Validation study methods for estimating exposure proportions and odds ratios with misclassified data. *Journal of Clinical Epidemiology*. 1990; 43:941–947. [PubMed: 2213082]
11. Morrissey MJ, Spiegelman D. Matrix methods for estimating odds ratios with misclassified exposure data: extensions and comparisons. *Biometrics*. 1999; 55:338–344. [PubMed: 11318185]
12. Espeland MA, Hui SL. A general approach to analyzing epidemiologic data that contain misclassification errors. *Biometrics*. 1987; 43:1001–1012. [PubMed: 3427157]
13. Spiegelman D, Rosner B, Logan R. Estimation and inference for logistic regression with covariate misclassification and measurement error in main study/validation study designs. *Journal of the American Statistical Association*. 2000; 95:51–61.
14. Zhou H, Pepe M. Auxiliary covariate data in failure time regression. *Biometrika*. 1995; 82:139–149.
15. Zhou H, Wang CY. Failure time regression with continuous covariates measured with error. *Journal of the Royal Statistical Society, Series B*. 2000; 62:657–665.
16. Chen YH. Cox regression in cohort studies with validation sampling. *Journal of the Royal Statistical Society, Series B*. 2002; 64:51–62.
17. Spiegelman D, McDermott A, Rosner B. The regression calibration method for correcting measurement error bias in nutritional epidemiology. *American Journal of Clinical Nutrition*. 1997; 65(Suppl.):1179S–1186S. [PubMed: 9094918]
18. Wang CY, Hsu L, Feng ZD, Prentice RL. Regression calibration in failure time regression. *Biometrics*. 1997; 53:131–145. [PubMed: 9147589]
19. Zucker DM, Spiegelman D. Inference for the proportional hazards model with misclassified discrete-valued covariates. *Biometrics*. 2004; 60:324–334. [PubMed: 15180657]
20. Zucker DM. A pseudo partial likelihood method for semi-parametric survival regression with covariate errors. *Journal of the American Statistical Association*. 2005; 100:1264–1277.
21. Hu P, Tsiatis AA, Davidian M. Estimating the parameters in the Cox model when covariate variables are measured with error. *Biometrics*. 1998; 54:1407–1419. [PubMed: 9883541]
22. Küchenhoff H, Mwalili SM, Lesaffre E. A general method for dealing with misclassification in regression: the misclassification SIMEX. *Biometrics*. 2006; 62:85–96. [PubMed: 16542233]
23. Herring AH, Ibrahim JG. Likelihood-based methods for missing covariates in the Cox proportional hazards model. *Journal of the American Statistical Association*. 2001; 96:292–302.
24. Nakamura T. Corrected score function of errors-in-variables models: methodology and application to generalized linear models. *Biometrika*. 1990; 77:127–137.

25. Nakamura T. Proportional hazards model with covariates subject to measurement error. *Biometrics*. 1992; 48:829–838. [PubMed: 1420844]
26. Akazawa K, Kinukawa N, Nakamura T. A note on the corrected score function corrected for misclassification. *Journal of the Japan Statistical Society*. 1998; 28:115–123.
27. Wu K, Willett WC, Fuchs CS, Colditz GA, Giovannucci EL. Calcium intake and risk of colon cancer in women and men. *Journal of the National Cancer Institute*. 2002; 94:437–446. [PubMed: 11904316]
28. Kalbfleisch, JD.; Prentice, RL. *The Statistical Analysis of Failure Time Data*. 2nd edn. New York: Wiley; 2002.
29. Thomas DC. General relative-risk models for survival time and matched case-control analysis. *Biometrics*. 1981; 37:673–686.
30. Breslow, N.; Day, NE. *Statistical Methods in Cancer Research, Volume 2: The Design and Analysis of Cohort Studies*. Oxford: Oxford University Press; 1993.
31. Cox DR. Partial likelihood. *Biometrika*. 1975; 62:269–276.
32. Andersen PK, Gill RD. Cox's regression model for counting processes: a large sample study. *Annals of Statistics*. 1982; 10:1100–1120.
33. Huang Y, Wang CY. Cox regression with accurate covariates unascertainable: a nonparametric-correction approach. *Journal of the American Statistical Association*. 2000; 95:1209–1219.
34. Spiegelman D, Carroll RJ, Kipnis V. Efficient regression calibration in main study/internal validation study designs with an imperfect reference instrument. *Statistics in Medicine*. 2001; 20:139–160. [PubMed: 11135353]
35. Hui SL, Zhou XH. Evaluation of diagnostic tests without gold standards. *Statistical Methods in Medical Research*. 1998; 7:354–370. [PubMed: 9871952]
36. Qu Y, Tan M, Kutner MH. Random effects models in latent class analysis for evaluating accuracy of diagnostic tests. *Biometrics*. 1996; 52:797–810. [PubMed: 8805757]
37. Torrance-Rynard VL, Walter SD. Effects of dependent errors in the assessment of diagnostic test performance. *Statistics in Medicine*. 1997; 6:2157–2175. [PubMed: 9330426]
38. Formann AK, Kohlmann T. Latent class analysis in medical research. *Statistical Methods in Medical Research*. 1996; 5:179–211. [PubMed: 8817797]
39. Albert PS, McShane LM, Shih JHUS. National Cancer Institute Bladder Tumor Marker Network. Latent class modeling approaches for assessing diagnostic error without a gold standard: with applications to p53 immunohistochemical assays in bladder tumors. *Biometrics*. 2001; 57:610–619. [PubMed: 11414591]
40. Armitage, P.; Doll, R. *Proceedings of the 4th Berkeley Symposium on Mathematical Statistics and Probability. Volume 4: Contributions to Biology and Problems of Medicine*. Berkeley, California: University of California Press; 1961. Stochastic models for carcinogenesis; p. 19-38.
41. Willett, WC. *Nutritional Epidemiology*. 2nd edn. New York: Oxford University Press; 1998.
42. Delves, LM.; Mohamed, JL. *Computational Methods for Integral Equations*. Cambridge: Cambridge University Press; 1985.
43. Stefanski L. Unbiased estimation of a nonlinear function of a normal mean with application to measurement-error models. *Communications in Statistics, Theory and Methods*. 1989; 18:4335–4358.
44. Cochran WG. The effectiveness of adjustment by subclassification in removing bias in observational studies. *Biometrics*. 1968; 24:295–313. [PubMed: 5683871]
45. Huang Y, Wang CY. Consistent function methods for logistic regression with errors in covariates. *Journal of the American Statistical Association*. 2001; 95:1209–1219.
46. Kong FH, Gu M. Consistent estimation in Cox's proportional hazards model with covariate measurement errors. *Statistica Sinica*. 1999; 9:953–969.
47. Hu C, Lin DY. Cox regression with covariate measurement error. *Scandinavian Journal of Statistics*. 2002; 29:637–655.
48. Lin DY, Wei LJ. The robust inference for the Cox proportional hazards model. *Journal of the American Statistical Association*. 1989; 84:1074–1078.

49. Liang KY, Zeger S. Longitudinal data analysis using generalized linear models. *Biometrika*. 1986; 73:13–22.

Table 1

Simulation results for a single binary covariate with known classification rates.

Per cent exposed	Error rate	True RR	Naive Cox	Per cent bias in estimated log RR			Standard deviation—CSCORE			MSE ratio	95 per cent CI coverage
				CSCORE	FWMLE	Empirical	Mean of estimates				
5	1	1.5	-15.89	-1.93	-1.92	0.196	0.194	1.00	1.00	95.12	
5	1	2.0	-14.00	-0.46	-0.43	0.176	0.173	1.00	1.00	95.10	
5	5	1.5	-48.45	-4.73	-4.77	0.264	0.255	1.00	1.00	94.90	
5	5	2.0	-46.08	-2.19	-2.16	0.231	0.218	1.00	1.00	93.90	
5	10	1.5	-66.02	-6.64	-6.99	0.349	0.338	1.03	1.03	94.76	
5	10	2.0	-64.16	-2.87	-2.90	0.313	0.281	0.99	0.99	93.62	
5	20	1.5	-82.24	-17.22	-23.81	0.591	0.620	1.22	1.22	93.07	
5	20	2.0	-80.92	-7.20	-9.87	0.531	0.482	1.04	1.04	91.87	
25	1	1.5	-3.03	-0.06	-0.02	0.096	0.096	1.00	1.00	94.92	
25	1	2.0	-3.01	-0.15	-0.10	0.090	0.089	1.00	1.00	94.46	
25	5	1.5	-14.09	0.02	0.06	0.106	0.105	1.00	1.00	94.56	
25	5	2.0	-13.89	-0.28	-0.22	0.100	0.097	1.00	1.00	94.62	
25	10	1.5	-26.61	-0.01	0.04	0.123	0.120	1.00	1.00	94.90	
25	10	2.0	-26.14	-0.33	-0.30	0.112	0.110	1.00	1.00	94.60	
25	20	1.5	-48.41	-0.54	-0.49	0.167	0.164	1.00	1.00	94.96	
25	20	2.0	-47.42	-0.18	-0.15	0.156	0.149	0.99	0.99	94.10	
40	1	1.5	-2.43	-0.36	-0.32	0.088	0.087	1.00	1.00	94.62	
40	1	2.0	-2.11	-0.03	0.01	0.081	0.082	1.00	1.00	95.22	
40	5	1.5	-10.03	0.39	0.43	0.095	0.094	1.00	1.00	94.58	
40	5	2.0	-10.35	0.02	0.07	0.091	0.089	1.00	1.00	94.90	
40	10	1.5	-20.63	0.01	0.05	0.108	0.106	1.00	1.00	94.74	
40	10	2.0	-20.40	0.28	0.30	0.102	0.101	1.00	1.00	95.04	
40	20	1.5	-41.04	-0.34	-0.33	0.140	0.142	1.00	1.00	95.58	
40	20	2.0	-40.66	0.29	0.34	0.135	0.135	1.00	1.00	94.90	

Note Sample size=2000, unexposed cumulative incidence=25 per cent.

RR, relative risk; CSCORE, corrected score estimator; FWMLE, full Weibull maximum likelihood estimator; MSE ratio, ratio of mean square error of FWMLE to that of CSCORE.

Table II

Simulation results for a single binary covariate with estimated classification rates.

Sample size	Per cent exposed	Error rate	True RR	Per cent bias in $\hat{\beta}$	Empirical Std dev.	Mean of SD estimates	95 per cent CI coverage
2000	5	1	1.5	-1.10	0.206	0.204	93.14
	5	1	2.0	0.22	0.185	0.185	93.72
	5	5	1.5	-1.96	0.293	0.292	93.30
	5	5	2.0	2.54	0.278	0.270	94.17
	5	10	1.5	-1.31	0.441	0.491	94.30
	5	10	2.0	4.48	0.412	0.439	94.87
	5	20	1.5	13.13	0.670	1.663	96.83
	5	20	2.0	11.30	0.637	1.758	95.37
	25	1	1.5	0.36	0.097	0.097	94.90
	25	1	2.0	0.23	0.091	0.090	94.70
	25	5	1.5	0.34	0.109	0.107	94.54
	25	5	2.0	0.47	0.101	0.100	95.02
	25	10	1.5	0.57	0.123	0.124	94.90
	25	10	2.0	0.83	0.117	0.117	95.44
1000	25	20	1.5	1.51	0.183	0.181	95.24
	25	20	2.0	3.39	0.187	0.179	95.56
	40	1	1.5	1.03	0.087	0.087	95.30
	40	1	2.0	0.59	0.083	0.083	94.84
	40	5	1.5	0.36	0.095	0.095	95.12
	40	5	2.0	0.54	0.091	0.091	94.78
	40	10	1.5	0.63	0.106	0.108	95.24
	40	10	2.0	0.66	0.106	0.104	94.98
	40	20	1.5	2.08	0.150	0.149	94.96
	40	20	2.0	1.79	0.148	0.148	95.72
	5	1	1.5	-4.07	0.296	0.297	91.58
	5	1	2.0	0.68	0.263	0.263	92.42
	5	5	1.5	-7.94	0.426	0.442	91.78
	5	5	2.0	0.56	0.382	0.383	92.30

Sample size	Per cent exposed	Error rate	True RR	Per cent bias in $\hat{\beta}$	Empirical Std dev.	Mean of SD estimates	95 per cent CI coverage
5	10	1.5	1.5	-12.20	0.600	0.790	91.48
5	10	2.0	2.0	-1.38	0.546	0.641	92.79
5	20	1.5	1.5	-15.98	2.805	2.130	96.76
5	20	2.0	2.0	-12.77	1.776	1.913	95.04
25	1	1.5	1.5	-0.47	0.137	0.137	94.56
25	1	2.0	2.0	0.55	0.126	0.127	95.46
25	5	1.5	1.5	0.25	0.153	0.151	94.76
25	5	2.0	2.0	0.89	0.141	0.141	95.16
25	10	1.5	1.5	0.20	0.176	0.175	94.80
25	10	2.0	2.0	0.33	0.165	0.163	94.24
25	20	1.5	1.5	0.34	0.253	0.255	94.78
25	20	2.0	2.0	2.91	0.241	0.239	95.18
40	1	1.5	1.5	0.33	0.121	0.124	95.32
40	1	2.0	2.0	0.22	0.116	0.117	95.58
40	5	1.5	1.5	0.80	0.135	0.135	94.90
40	5	2.0	2.0	0.68	0.130	0.129	94.86
40	10	1.5	1.5	0.50	0.154	0.153	95.20
40	10	2.0	2.0	0.60	0.147	0.146	94.88
40	20	1.5	1.5	2.92	0.210	0.211	95.12
40	20	2.0	2.0	1.93	0.202	0.204	95.84

Note Classification rates estimated from an external replicate measures sample of size 250. Unexposed cumulative incidence=25 per cent. RR, relative risk.

Table III

Estimated coefficients and standard errors for the Nurses Health Study of the relationship between dietary calcium intake and distal colon cancer incidence.

Method	High calcium		BMI of 22 to <25		BMI of 25 to <30		BMI of 30+		Aspirin use	
	Estimate	Std error	Estimate	Std error	Estimate	Std error	Estimate	Std error	Estimate	Std error
Cox	-0.3448	0.1694	0.6837	0.2240	0.5352	0.2395	0.5729	0.2876	-0.4941	0.1954
CS0	-0.7121	0.3690	0.7124	0.2247	0.5776	0.2419	0.6157	0.2892	-0.4994	0.1955
CS1	-0.7121	0.3832	0.7124	0.2249	0.5776	0.2423	0.6157	0.2896	-0.4994	0.1955

Note Cox, classical Cox regression analysis; CS0, corrected score method, observed classification matrix taken as known; CS1, corrected score method, accounting for uncertainty in the classification matrix.