



Published in final edited form as:

Biometrics. 2019 June ; 75(2): 414–427. doi:10.1111/biom.13012.

A Modified Partial Likelihood Score Method for Cox Regression with Covariate Error Under the Internal Validation Design

David M. Zucker^{1,*}, Xin Zhou², Xiaomei Liao³, Yi Li⁴, Donna Spiegelman⁵

¹Department of Statistics and Data Science, The Hebrew University of Jerusalem, Mount Scopus, Jerusalem 91905, ISRAEL

²Department of Biostatistics, Harvard T.H. Chan School of Public Health, 677 Huntington Avenue, Boston, MA 02115, U.S.A.,

³Department of Biostatistics, Harvard T.H. Chan School of Public Health, 677 Huntington Avenue, Boston, MA 02115, U.S.A. Currently employed at AbbVie Incorporated, North Chicago, IL, U.S.A.

⁴Department of Biostatistics, University of Michigan School of Public Health, 1415 Washington Heights, Ann Arbor, MI 48109-2029, U.S.A.,

⁵Department of Biostatistics, Yale School of Public Health and Department of Statistics, Yale University, New Haven, CT 06520 U.S.A, and Departments of Epidemiology, Biostatistics, Nutrition and Global Health, Harvard T.H. Chan School of Public Health, Boston, MA 02115, U.S.A

Summary:

We develop a new method for covariate error correction in the Cox survival regression model, given a modest sample of internal validation data. Unlike most previous methods for this setting, our method can handle covariate error of arbitrary form. Asymptotic properties of the estimator are derived. In a simulation study, the method was found to perform very well in terms of bias reduction and confidence interval coverage. The method is applied to data from Health Professionals Follow-Up Study (HPFS) on the effect of diet on incidence of Type II diabetes.

Keywords

Cox model; Measurement error; modified score

1. Introduction

In the Cox (1972) regression model for survival data, the hazard function $\lambda(t|\mathbf{x})$ for an individual with covariate vector $\mathbf{x} \in \mathbb{R}^p$ is modeled semiparametrically as

* david.zucker@mail.huji.ac.il.

Current address: Center for Methods in Implementation and Prevention Science (CMIPS) and Department of Biostatistics, Yale School of Public Health, 60 College Street, New Haven, CT 06520, U.S.A

⁶Supplementary Materials

The Web Appendix, referenced in Sections 2 and 3, is available with this paper at the Biometrics website on Wiley Online Library, as is the code we used to implement the various method.

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

where $\boldsymbol{\beta} \in \mathbb{R}^p$ is a vector of regression coefficients and $\lambda_0(t)$ is an unspecified baseline hazard function $\lambda_0(t)$. Cox proposed drawing inference on $\boldsymbol{\beta}$ based on the notion of partial likelihood, which was subsequently justified rigorously by Tsiatis (1981), who used classical limit theory, and by Andersen and Gill (1982), who used a martingale theory approach.

In many applications, however, the covariate \mathbf{X} is not measured exactly, but is subject to measurement error of some degree, often substantial. Thus, instead of observing \mathbf{X} , we observe a surrogate measure \mathbf{W} . Starting from Prentice (1982), a considerable literature has been developed on inference for the Cox regression model with covariate error in various contexts; see Zucker (2005) for a brief review.

The existing methods generally involve some model assumptions on the joint distribution of the true covariate and the surrogate. Many of the methods make use of specific parametric forms for this joint distribution. Other methods, such as those of Huang and Wang (2000) and Kong and Gu (1999), avoid use of a specific parametric form but still rely on an assumption that the covariate error is of independent additive structure. Some papers, such as Zhou and Pepe (1995), Zhou and Wang (2000), and Chen (2002), present methods without this additive error assumption for the internal validation design in which there is a subsample of individuals with a measurement on both the true covariate and the surrogate. These methods, however, have challenges as well. The approach taken by Zhou and Pepe (1995) and by Zhou and Wang (2000) involves stratification or smoothing in the covariate space; when the number of covariates is moderate to large, this approach breaks down due to the “curse of dimensionality.” Chen (2002) 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 to hundreds of thousands, and the validation study sample size is, as in all applications we know of, only a few hundred. Under these circumstances, the number of events in the validation study is very small, so that a satisfactory initial estimate of the regression coefficient vector based on the validation sample alone cannot be obtained. Thus, in such situations, which often arise in practice, Chen’s approach is problematic.

This paper presents a new method for the Cox model with covariate error, which overcomes the limitations of previously proposed methods. The method involves a modified version of the classical Cox partial likelihood score function, with the internal validation data incorporated in a suitable way. Our approach is very simple in concept. It is in the spirit of Lin and Ying’s (1993) work on Cox regression with incomplete covariate data. There is also some resemblance to Huang and Wang’s (2000) method for Cox regression with covariate error, and to work of Kulich and Lin (2000, 2004). The method requires no assumptions on the form of the covariate error. It is especially designed for the internal validation design with a relatively small validation sample and a moderate to large number of covariates, which, as indicated above, is a challenging situation that often arises in epidemiological

studies. The method is easy to implement, and its practical utility is backed by large-sample theory and small-sample simulations.

The outline of the remainder of the paper is as follows. Section 2 presents the proposed method and its asymptotic properties, Section 3 a simulation study, Section 4 an application to data from the Health Professionals Follow-Up Study (HPFS), and Section 5 a brief summary. The Web Appendix provides theoretical details.

2. The Proposed Method and Its Asymptotic Properties

2.1 The Proposed Method

We assume a classical survival data setup. We have i.i.d. observations on n individuals. Associated with each individual i is a set of random variables $(T_i^\circ, C_i, \mathbf{X}_i, \mathbf{W}_i)$, with T_i° representing the time to event, C_i representing the time to censoring, \mathbf{X}_i representing a p -vector of true covariate values, and \mathbf{W}_i representing a p -vector of surrogate covariate values. We assume that the covariates are arranged so that the first p_1 covariates are the error-prone covariates and the remaining $p_2 = p - p_1$ covariates are error-free. For the error-free covariates, the relevant component of \mathbf{W}_i is identical to the corresponding component of \mathbf{X}_i . We denote the maximum follow-up time by τ . The available data on all individuals consist of $(T_i, \delta_i, \mathbf{W}_i)$, where $T_i = \min(T_i^\circ, C_i)$ is the follow-up time and $\delta_i = I(T_i^\circ \leq C_i)$, with $I(\cdot)$ being the indicator function, is the event indicator. In addition, within the main study we have a random internal validation sample of size m of individuals with both \mathbf{X}_i and \mathbf{W}_i observed. We take $m = \text{ceil}(\pi n)$, where π is a specified number in $(0, 1)$ and $\text{ceil}(u)$ denotes the smallest integer greater than or equal to u . We define ω_i to be equal to 1 if individual i is in the internal validation sample and 0 otherwise. Thus, the random vector $(\omega_1, \dots, \omega_n)$ has a uniform distribution over the finite set $\mathcal{O}(n, m)$ of vectors with m ones and $n - m$ zeros (i.e., $\mathcal{O}(n, m)$ expresses the various ways of selecting m elements from a set of n elements). We write $\hat{\pi} = m/n$. Note that $\hat{\pi}$ is not an estimate, but rather is fixed by design. Also, as usual, we define $Y_i(t) = I(T_i \geq t)$ and $N_i(t) = \delta_i I(T_i \leq t)$. Left truncation is handled by setting $Y_i(t)$ to zero until the time at which individual i comes under observation.

We assume, as usual, that T_i° and C_i are conditionally independent given \mathbf{X}_i . We assume further that the measurement error is noninformative in the sense that \mathbf{W}_i is conditionally independent of (T_i°, C_i) given \mathbf{X}_i . We make no assumptions about the form of the measurement error. Finally, we assume that the survival time T_i° follows the Cox model (1).

We denote the true value of $\boldsymbol{\beta}$ by $\boldsymbol{\beta}^*$. We present our development for the case of the classical Cox relative risk function $e^{\boldsymbol{\beta}^T \mathbf{X}_i}$, but it is straightforward to extend the development to more general relative risk functions, as in Thomas (1981) and in Breslow and Day, 1993, Sec.5.1(c).

We construct our procedure as follows. Let \hat{E} denote empirical expectation, so that, for example,

$$\hat{E}\{Y(t)\exp(\boldsymbol{\beta}^T \mathbf{X})\} = \frac{1}{n} \sum_{j=1}^n Y_{j(t)} \exp(\boldsymbol{\beta}^T \mathbf{X}_j),$$

$$\hat{E}\{Y(t)\mathbf{X}\exp(\boldsymbol{\beta}^T \mathbf{X})\} = \frac{1}{n} \sum_{j=1}^n Y_{j(t)} \mathbf{X}_j \exp(\boldsymbol{\beta}^T \mathbf{X}_j).$$

In the absence of measurement error, the Cox partial likelihood score function is given by

$$\mathbf{U}_{COX}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \delta_i \left(\mathbf{X}_i - \left[\frac{\hat{E}\{Y(t)\mathbf{X}\exp(\boldsymbol{\beta}^T \mathbf{X})\}}{\hat{E}\{Y(t)\exp(\boldsymbol{\beta}^T \mathbf{X})\}} \right]_{t=T_i} \right). \quad (2)$$

When \mathbf{X} is measured only for a sample of the individuals and only \mathbf{W} is available for the others, a naive Cox analysis involves simply substituting \mathbf{W} in place of \mathbf{X} for the individuals without a measurement of \mathbf{X} . In other words, defining $\mathbf{W}_i^\circ = \omega_i \mathbf{X}_i + (1 - \omega_i) \mathbf{W}_i$, the naive Cox analysis is based on the score function

$$\mathbf{U}_{NAI}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \delta_i \left(\mathbf{W}_i^\circ - \left[\frac{\hat{E}\{Y(t)\mathbf{W}^\circ \exp(\boldsymbol{\beta}^T \mathbf{W}^\circ)\}}{\hat{E}\{Y(t)\exp(\boldsymbol{\beta}^T \mathbf{W}^\circ)\}} \right]_{t=T_i} \right), \quad (3)$$

with

$$\hat{E}\{Y(t)\exp(\boldsymbol{\beta}^T \mathbf{W}^\circ)\} = S_0^\circ(t, \boldsymbol{\beta}) = \frac{1}{n} \sum_{j=1}^n Y_{j(t)} \exp(\boldsymbol{\beta}^T \mathbf{W}_j^\circ),$$

$$\hat{E}\{Y(t)\mathbf{W}^\circ \exp(\boldsymbol{\beta}^T \mathbf{W}^\circ)\} = \mathbf{S}_1^\circ(t, \boldsymbol{\beta}) = \frac{1}{n} \sum_{j=1}^n Y_{j(t)} \mathbf{W}_j^\circ \exp(\boldsymbol{\beta}^T \mathbf{W}_j^\circ)$$

We denote the corresponding estimator by $\hat{\boldsymbol{\beta}}_{NAI}$. The terms in $\mathbf{U}_{NAI}(\boldsymbol{\beta})$ are of “observed – expected” form, but the “expected” term is incorrect. Consequently, the naive score function does not have zero asymptotic expectation under $\boldsymbol{\beta}^*$, and therefore $\hat{\boldsymbol{\beta}}_{NAI}$ is biased.

An improved estimator can be obtained using regression calibration, which is an established technique for measurement error problems; see, for example, Carroll et al. (2006, Chapter 4). In regression calibration, we redefine \mathbf{W}_i° to be $\mathbf{W}_i^\circ = \omega_i \mathbf{X}_i + (1 - \omega_i) \hat{\mathbf{X}}_{ir}$, with \hat{X}_{ir} ($r = 1, \dots, p_1$) defined as

$$\hat{X}_{ir} = \hat{\alpha}_{r0} + \sum_{s=1}^p \hat{\alpha}_{rs} W_{is} \quad (4)$$

where $\hat{\alpha}_{r0}, \dots, \hat{\alpha}_{rp}$ are the ordinary least squares estimates of the regression of X_{ir} on \mathbf{W}_i based on the internal validation sample. Having redefined \mathbf{W}_i^o , we redefine $S_0^o(t, \boldsymbol{\beta})$ and $S_1^o(t, \boldsymbol{\beta})$ correspondingly. We denote the resulting estimator by $\hat{\boldsymbol{\beta}}_{RC}$. In (4), for the sake of generality, we have included all of the components of \mathbf{W}_i in the regression, but in typical applications of regression calibration the regression model for X_{ir} includes only W_{ir} and perhaps one or two additional components of \mathbf{W}_i . Substantial improvement is often achieved with regression calibration approach, but the “expected” term is still not exactly correct, and therefore the resulting estimator is not exactly consistent. The regression calibration approximation is good when the degree of measurement error is small or the regression coefficients of the error-prone covariates are small, but otherwise the approximation can be unsatisfactory (Spiegelman, Rosner, and Logan, 2000).

We present an estimator that builds on the regression calibration estimator but is exactly consistent. As in regression calibration, we use the regression model (4). However, we use this model only as a working model, and it is not necessary for the model to be correct for our estimator to be consistent. As with standard regression calibration, it is possible in principle, as written in (4), to include all the components of \mathbf{W}_i in the model, but in practice we recommend using only W_{ir} and perhaps one or two additional components.

The idea of our approach is to replace the incorrect “expected” term with a correct one. Let $\boldsymbol{\alpha}^{(r)}$ be the column vector with components $\alpha_{r0}, \alpha_{r1}, \dots, \alpha_{rp}$, let $\boldsymbol{\alpha}$ denote the vector formed by stacking the vectors $\boldsymbol{\alpha}^{(r)}$ one on top of the other, and let $\boldsymbol{\alpha}^*$ denote the true value of $\boldsymbol{\alpha}$. To emphasize the dependence of \hat{X}_{ir} on $\boldsymbol{\alpha}$, we denote the vector of X_{ir} 's by $\hat{\mathbf{X}}_i(\boldsymbol{\alpha})$.

Define

$$S_{0a}(t, \boldsymbol{\beta}) = \frac{1}{n} \sum_{j=1}^n \omega_j Y_j(t) \exp(\boldsymbol{\beta}^T \mathbf{X}_j) \quad (5)$$

$$S_{0b}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{n} \sum_{j=1}^n (1 - \omega_j) Y_j(t) \exp\{\boldsymbol{\beta}^T \hat{\mathbf{X}}_j(\boldsymbol{\alpha})\} \quad (6)$$

$$S_{0c}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{n} \sum_{j=1}^n \omega_j Y_j(t) \exp\{\boldsymbol{\beta}^T \hat{\mathbf{X}}_j(\boldsymbol{\alpha})\} \quad (7)$$

$$\mathbf{S}_{1a}(t, \boldsymbol{\beta}) = \frac{1}{n} \sum_{j=1}^n \omega_j Y_j(t) \mathbf{X}_j \exp(\boldsymbol{\beta}^T \mathbf{X}_j) \quad (8)$$

$$\mathbf{S}_{1b}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{n} \sum_{j=1}^n (1 - \omega_j) Y_j(t) \widehat{\mathbf{X}}_j(\boldsymbol{\alpha}) \exp\{\boldsymbol{\beta}^T \widehat{\mathbf{X}}_j(\boldsymbol{\alpha})\} \quad (9)$$

$$\mathbf{S}_{1c}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{n} \sum_{j=1}^n \omega_j Y_j(t) \widehat{\mathbf{X}}_j(\boldsymbol{\alpha}) \exp\{\boldsymbol{\beta}^T \widehat{\mathbf{X}}_j(\boldsymbol{\alpha})\} \quad (10)$$

$$\widetilde{\mathbf{S}}_{1a}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{n} \sum_{j=1}^n \omega_j Y_j(t) \widehat{\mathbf{X}}_j(\boldsymbol{\alpha}) \exp(\boldsymbol{\beta}^T \mathbf{X}_j) \quad (11)$$

$$S_0(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = S_{0a}(t, \boldsymbol{\beta}) + \left\{ \frac{S_{0a}(t, \boldsymbol{\beta})}{S_{0c}(t, \boldsymbol{\beta}, \boldsymbol{\alpha})} \right\} S_{0b}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) \quad (12)$$

$$\widetilde{\phi}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = S_{0b}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) / S_{0c}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) \quad (13)$$

$$\mathbf{S}_1(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \mathbf{S}_{1a}(t, \boldsymbol{\beta}) + \mathbf{S}_{1b}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) + \widetilde{\phi}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) \{ \widetilde{\mathbf{S}}_{1a}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) - \mathbf{S}_{1c}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) \} \quad (14)$$

We then take the score function to be

$$\mathbf{U}_{MS}(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{n} \sum_{i=1}^n \delta_i \left\{ \mathbf{W}_i^\circ - \frac{\mathbf{S}_1(T_i, \boldsymbol{\beta}, \boldsymbol{\alpha})}{S_0(T_i, \boldsymbol{\beta}, \boldsymbol{\alpha})} \right\}. \quad (15)$$

The estimator $\widehat{\boldsymbol{\beta}}_{MS}$ is defined to be the solution to the score equation $\mathbf{U}_{MS}(\boldsymbol{\beta}, \widehat{\boldsymbol{\alpha}}) = 0$. We could have used $\widehat{\phi} = (1 - \widehat{\pi})/\widehat{\pi} = (n - m)/m$ in place of $\widetilde{\phi}(t, \boldsymbol{\beta}, \boldsymbol{\alpha})$, but we found that better finite-sample performance is achieved with $\widetilde{\phi}(t, \boldsymbol{\beta}, \boldsymbol{\alpha})$.

The motivation behind $\mathbf{U}_{MS}(\boldsymbol{\beta}, \boldsymbol{\alpha})$ is as follows. The regression calibration function $\mathbf{U}_{RC}(\boldsymbol{\beta})$ can be written in counting process notation as

$$\mathbf{U}_{RC}(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{n} \sum_{i=1}^n \int_0^{\tau} \{\mathbf{W}_i^{\circ} - E(t, \boldsymbol{\beta}, \boldsymbol{\alpha})\} dN_i(t)$$

with

$$E(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{\mathbf{S}_1^{\circ}(t, \boldsymbol{\beta}, \boldsymbol{\alpha})}{S_0^{\circ}(t, \boldsymbol{\beta}, \boldsymbol{\alpha})}.$$

Let us now define $dM_i(t) = dN_i(t) - Y_i(t)e^{\boldsymbol{\beta}^* T \mathbf{X}_i} \lambda_0(t) dt$. We can then write

$$\begin{aligned} \mathbf{U}_{RC}(\boldsymbol{\beta}, \boldsymbol{\alpha}) &= \frac{1}{n} \sum_{i=1}^n \int_0^{\tau} \{\mathbf{W}_i^{\circ} - E(t, \boldsymbol{\beta}, \boldsymbol{\alpha})\} Y_i(t) e^{\boldsymbol{\beta}^* T \mathbf{X}_i} \lambda_0(t) dt \quad (16) \\ &+ \frac{1}{n} \sum_{i=1}^n \int_0^{\tau} \{\mathbf{W}_i^{\circ} - E(t, \boldsymbol{\beta}, \boldsymbol{\alpha})\} dM_i(t) \\ &= \int_0^{\tau} \{\mathbf{S}_1^{\circ\circ}(t, \boldsymbol{\beta}^*, \boldsymbol{\alpha}) - E(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) S_{0d}(t, \boldsymbol{\beta}^*)\} \lambda_0(t) dt \\ &+ \frac{1}{n} \sum_{i=1}^n \int_0^{\tau} \{\mathbf{W}_i^{\circ} - E(t, \boldsymbol{\beta}, \boldsymbol{\alpha})\} dM_i(t). \end{aligned}$$

where

$$S_{0d}(t, \boldsymbol{\beta}) = \frac{1}{n} \sum_{j=1}^n Y_j(t) e^{\boldsymbol{\beta}^T \mathbf{X}_j}$$

$$\mathbf{S}_1^{\circ\circ}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \frac{1}{n} \sum_{j=1}^n Y_j(t) \mathbf{W}_j^{\circ} e^{\boldsymbol{\beta}^T \mathbf{X}_j} = \mathbf{S}_{1d}(\boldsymbol{\beta}, \boldsymbol{\alpha}) + \frac{1}{n} \sum_{j=1}^n (1 - \omega_j) Y_j(t) \widehat{\mathbf{X}}_j(\boldsymbol{\alpha}) e^{\boldsymbol{\beta}^T \mathbf{X}_j}$$

Using counting process theory (Gill, 1984), it can be seen that the second term of (16) has expectation zero. In the absence of measurement error, the value at $\boldsymbol{\beta}^*$ of the quantity in brackets in the first term of (16) is zero, so that the score function is unbiased. In the presence of measurement error, the value at $\boldsymbol{\beta}^*$ of this quantity is in general nonzero. We need to redefine $E(t, \boldsymbol{\beta}, \boldsymbol{\alpha})$ so that the limiting value of this quantity at $\boldsymbol{\beta}^*$, $\boldsymbol{\alpha}^*$ is zero. Define

$$\mathcal{E}_X(t, \boldsymbol{\beta}) = E\{Y(t) \exp(\boldsymbol{\beta}^T \mathbf{X})\},$$

$$\mathcal{E}_{XX}(t, \boldsymbol{\beta}) = E\{Y(t) \mathbf{X} \exp(\boldsymbol{\beta}^T \mathbf{X})\},$$

$$\mathcal{E}_{WX}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = E\left\{Y(t)\widehat{\mathbf{X}}(\boldsymbol{\alpha})\exp(\boldsymbol{\beta}^T\mathbf{X})\right\}.$$

The limiting value of $S_{0d}(t, \boldsymbol{\beta})$ is then $\mathcal{E}_X(t, \boldsymbol{\beta}, \boldsymbol{\alpha})$ and the limiting value of $S_1^{\circ}(t, \boldsymbol{\beta}, \boldsymbol{\alpha})$ is $s_1(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = \pi\mathcal{E}_{XX}(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) + (1 - \pi)\mathcal{E}_{WX}(t, \boldsymbol{\beta}, \boldsymbol{\alpha})$. We thus have to redefine $E(t, \boldsymbol{\beta}, \boldsymbol{\alpha})$ so that its limiting value is equal to $s_1(t, \boldsymbol{\beta}, \boldsymbol{\alpha})/\mathcal{E}_X(t, \boldsymbol{\beta}, \boldsymbol{\alpha})$. Taking $E(t, \boldsymbol{\beta}, \boldsymbol{\alpha}) = S_1(t, \boldsymbol{\beta}, \boldsymbol{\alpha})/S_0(t, \boldsymbol{\beta}, \boldsymbol{\alpha})$ achieves this objective. At the same time, our estimator reduces to the usual Cox estimator under zero measurement error. We regard this reducibility property to be important for measurement error correction methods.

We reiterate that our method makes no assumptions about the form of the covariate error, and that the model (4) is only a working model, with our estimator still being consistent even if the working model is misspecified. In addition, our method requires only estimation of unconditional means involving Y , \mathbf{W} , and \mathbf{X} , and therefore does not require use of smoothing methods. For this reason, a modestly-sized internal validation sample is sufficient. By contrast, the approaches taken by Zhou and Pepe (1995) and by Zhou and Wang (2000) require consistent estimates of conditional means, which involve stratification or smoothing in the covariate space, and thus require a larger validation sample. In addition, since our method is based on separate empirical averages for each risk set, a rare disease approximation is not needed.

We have worked in the setting of time-independent covariates, but it is possible to consider extension to the case of time-dependent covariates. When the covariate processes are measured on an approximately continuous basis ($\mathbf{W}(t)$ for the full cohort and $\mathbf{X}(t)$ for the internal validation sample), the method and its asymptotic theory carries over with notational changes only. Since the method is based on separate empirical averages for each risk set, changes over time in the measurement error distribution are handled automatically. The method and the asymptotic theory also carry over to the case where the covariate processes are measured only intermittently, as commonly occurs in practice, but the processes vary slowly, so that carrying forward the last observed covariate value is a reasonable approximation. In the case where the the covariate processes are measured only intermittently and vary more rapidly, the extension to the case of time-dependent covariates is more complex and is beyond the scope of this paper.

2.2 Asymptotic Properties

The asymptotic properties of the estimator are presented in the following theorem.

Theorem 1: Under the regularity conditions stated in the Web Appendix, $\widehat{\boldsymbol{\beta}}_{MS}$ converges almost surely to $\boldsymbol{\beta}^$, and $\sqrt{n}(\widehat{\boldsymbol{\beta}}_{MS} - \boldsymbol{\beta}^*)$ is asymptotically mean-zero multivariate normal with covariance matrix that can be estimated consistently by the sandwich-type estimator described below.*

We present here a sketch of the proof of this result. The details are presented in the Web Appendix.

The consistency proof hinges on the fact that, as explained above, $\mathbf{U}_{MS}(\boldsymbol{\beta}, \boldsymbol{\alpha}_*)$ is constructed so that it converges to a limit $\mathbf{u}(\boldsymbol{\beta}, \boldsymbol{\alpha})$ for which $\mathbf{u}(\boldsymbol{\beta}^*, \boldsymbol{\alpha}^*) = \mathbf{0}$. We can then appeal to arguments of Foutz (1977) to obtain the consistency result.

The asymptotic normality proof is based on estimating equations theory, and uses an argument along the lines of Lin and Wei (1989). Setting $\boldsymbol{\theta} = (\boldsymbol{\beta}, \boldsymbol{\alpha})$, we can define the estimator $\hat{\boldsymbol{\theta}}$ of $\boldsymbol{\theta}$ to be the solution $\hat{\boldsymbol{\theta}}$ to $\mathcal{U}(\boldsymbol{\theta}) = 0$ with $\mathcal{U}(\boldsymbol{\theta}) = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)})$, where $\mathbf{U}^{(1)}(\boldsymbol{\theta})$ is the $\mathbf{U}_{MS}(\boldsymbol{\beta}, \boldsymbol{\alpha})$ defined in (15) and $\mathbf{U}^{(2)}(\boldsymbol{\theta})$ is given by stacking the vectors

$$\mathbf{U}_r^{(2)}(\boldsymbol{\alpha}) = \frac{1}{n} \sum_{i=1}^n \omega_i \left(X_{ir} - \alpha_{r0} - \sum_{s=1}^p \alpha_{rs} W_{is} \right) \begin{bmatrix} 1 \\ \mathbf{W}_i \end{bmatrix}$$

where we include $\mathbf{U}_r^{(2)}$ only for covariates that are subject to measurement error. We can write

$$\mathbf{U}^{(2)}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \omega_i \mathbf{Z}_i^{(12)}(\boldsymbol{\theta})$$

with

$$\mathbf{Z}_i^{(12)}(\boldsymbol{\theta}) = (\mathbf{x}_i \otimes \bar{\mathbf{w}}_i) - \left\{ \mathbf{I}_{p_1} \otimes (\bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T) \right\} \boldsymbol{\alpha}$$

where \mathbf{x}_i consists of X_{i1}, \dots, X_{ip_1} , $\bar{\mathbf{w}}_i$ consists of a 1 followed by the components of \mathbf{W}_i . \otimes denotes the Kronecker product, and \mathbf{I}_b denotes the $b \times b$ identity matrix. The vector $\mathbf{U}^{(2)}(\boldsymbol{\theta})$ is of length $(p+1)p_1$. When the model for a given X_{ir} includes only some of the W_{is} 's, we delete the superfluous elements of $\boldsymbol{\alpha}$ and $\mathbf{Z}_i^{(12)}(\boldsymbol{\theta})$.

In the Web Appendix we show that $\mathbf{U}^{(1)}(\boldsymbol{\theta}^*)$ is asymptotically equivalent to the quantity

$$\mathbf{U}^{(1)*}(\boldsymbol{\theta}^*) = \hat{\pi} \left\{ \frac{1}{m} \sum_{i=1}^n \omega_i \mathbf{Z}_i^{(11)}(\boldsymbol{\theta}^*) \right\} + (1 - \hat{\pi}) \left\{ \frac{1}{n-m} \sum_{i=1}^n (1 - \omega_i) \mathbf{Z}_i^{(21)}(\boldsymbol{\theta}^*) \right\}$$

where $\{\mathbf{Z}_i^{(11)}(\boldsymbol{\theta}) : \omega_i = 1\}$ and $\{\mathbf{Z}_i^{(21)}(\boldsymbol{\theta}) : \omega_i = 0\}$ are each sets of i.i.d. vectors with mean zero under $\boldsymbol{\theta} = \boldsymbol{\theta}^*$, the expressions for which are presented in the Web Appendix. Thus, the solution to $\mathcal{U}(\boldsymbol{\theta}^*) = \mathbf{0}$ is asymptotically equivalent to the solution to $\mathcal{U}^*(\boldsymbol{\theta}^*) = 0$, with $\mathcal{U}^* = (\mathbf{U}^{(1)*}, \mathbf{U}^{(2)})$. Let $\mathbf{Z}_i^{(1)}$ denote the stacked vector formed by $\mathbf{Z}_i^{(11)}$ and $\mathbf{Z}_i^{(12)}$ and let $\mathbf{Z}_i^{(2)}$ denote the stacked vector formed by $\mathbf{Z}_i^{(21)}$ and the zero vector of length $(p+1)p_1$. We can then write

$$\mathcal{U}^*(\boldsymbol{\theta}) = \hat{\pi} \left\{ \frac{1}{m} \sum_{i=1}^n \omega_i \mathbf{Z}_i^{(1)}(\boldsymbol{\theta}^*) \right\} + (1 - \hat{\pi}) \left\{ \frac{1}{n-m} \sum_{i=1}^n (1 - \omega_i) \mathbf{Z}_i^{(2)}(\boldsymbol{\theta}^*) \right\}.$$

Define $\mathbf{C}_1 = \text{Cov}(\mathbf{Z}_i^{(1)})$, $\mathbf{C}_2 = \text{Cov}(\mathbf{Z}_i^{(2)})$, and $\mathbf{C} = \hat{\pi}\mathbf{C}_1 + (1 - \hat{\pi})\mathbf{C}_2$. We see that the asymptotic distribution of $\sqrt{n}\mathcal{U}^*(\boldsymbol{\theta}^*)$ is mean-zero normal with covariance matrix \mathbf{C} . Consequently $\sqrt{n}(\hat{\boldsymbol{\beta}}_{MS} - \boldsymbol{\beta}^*)$ is asymptotically mean-zero normal with covariance matrix $\mathbf{V} = \mathcal{R}\mathbf{C}\mathcal{R}^T$, where \mathcal{R} is the matrix consisting of the first p rows of $\mathbf{d}(\boldsymbol{\theta})^{-1}$, where $\mathbf{d}(\boldsymbol{\theta})$ is the limiting value of the matrix $\mathbf{D}(\boldsymbol{\theta})$ given by -1 times the Jacobian of $\mathcal{U}(\boldsymbol{\theta})$. In principle, we can estimate \mathbf{V} by $\hat{\mathbf{V}} = \hat{\mathcal{R}}\hat{\mathbf{C}}\hat{\mathcal{R}}^T$, where $\hat{\mathcal{R}}$ consists of the first p rows of $\mathbf{D}(\hat{\boldsymbol{\theta}})^{-1}$ and $\hat{\mathbf{C}} = \hat{\pi}\hat{\mathbf{C}}_1 + (1 - \hat{\pi})\hat{\mathbf{C}}_2$, where $\hat{\mathbf{C}}_s$ is the sample covariance of $\mathbf{Z}_i^{(s)}(\boldsymbol{\theta})$, i.e.

$$\hat{\mathbf{C}}_s = \frac{1}{n} \sum_{i=1}^n \mathbf{Z}_i^{(s)}(\boldsymbol{\theta}) \mathbf{Z}_i^{(s)}(\boldsymbol{\theta})^T. \quad (17)$$

In actuality, the terms of $\mathbf{U}^{(1)*}$ involve additional unknown quantities, so we compute $\hat{\mathbf{C}}_s$ using the sample covariance of the vectors $\hat{\mathbf{Z}}_i^{(s)}(\hat{\boldsymbol{\theta}})$ defined by replacing these quantities with consistent estimates. The detailed derivations of the expressions for $\mathbf{Z}_i^{(11)}(\boldsymbol{\theta})$, $\mathbf{Z}_i^{(21)}(\boldsymbol{\theta})$, and $\mathbf{D}(\boldsymbol{\theta})$ are presented in the Web Appendix.

3. Simulation Study

We examined the performance of the proposed method in a simulation study. We constructed the simulation setup so as to be representative of a typical epidemiological cohort study. We considered a setup where the time metameter is age, with age at entry to the study being uniformly distributed over the interval 30 to 50 years. The study horizon was 12 years. We took the censoring distribution to be exponential with a rate of 1% per year. We took the baseline survival function to be Weibull with shape parameter 5, as in Zucker and Spiegelman (2004, 2008). In terms of the sample size and the event rate (determined by the Weibull scale parameter), we considered two scenarios: a rare event scenario with $n = 10,000$ and a cumulative event rate of about 5% (so that the number of events is about 500), and a common event scenario with $n = 500$ and a cumulative event rate of about 25% (so that the number of events is about 125). The internal validation sample size was 200. Thus, in the rare event case, the internal validation sample size included a mere handful of events, which may hamper the use of Chen's (2002) approach.

We carried out two sets of simulations. In the first set, we worked with a single covariate X , generated from a standard normal distribution. We considered two measurement error models, as follows:

Independent Measurement Error Model: $W = X + \epsilon$ with $\epsilon \sim \mathcal{N}(0, a)$ independently of X

Dependent Measurement Error Model: $W = X + \epsilon$ with $\epsilon | X \sim N(0, a(1 + |X|))$

We chose a range of a values corresponding to the following range of values for the correlation between X and W : 0.9, 0.7, 0.5. Finally, we took $e^\beta = 1.5, 2.5, \text{ or } 4$. We compared our proposed estimator (MS) against Chen's (2002) estimator (CH), the regression calibration estimator obtained by replacing X by \hat{X} in the Cox score function (RC), the "complete case" (CC) estimator based only on the data with a measurement of X . In the second set of simulations, we worked with five covariates X_1, \dots, X_5 , with X_1 error-prone and the other covariates error-free. We took the five covariates to be $N(0, 1)$ random variables, either independent or equally-correlated with a correlation of 0.2. We took the hazard function to be $\lambda(t) = \lambda_0(t) \exp(\beta_1 \times_1 + \beta_2 \times_2 + \beta_3 \times_3 + \beta_4 \times_4 + \beta_5 \times_5)$ with $\beta_2 = \beta_3 = \beta_4 = \beta_5 = \log(1.5)$, where, as before, we took $\lambda_0(t)$ to be Weibull with shape parameter 5 and $e^\beta = 1.5, 2.5, \text{ or } 4$. The other settings were as in the first set of simulations. The simulation results were based on 10,000 replications. If the zero-finding procedure with our method failed to converge, we used the RC estimate. In the univariate simulations this usually occurred in less than 1% of the replications, and the worst instance it occurred in 6% of the replications. In the multivariate simulations, convergence failure usually occurred in less than 5% of the replications, and in the worst instance it occurred in 10% of the replications. In both the univariate and multivariate simulation, the worst instance was with highest value of β_1 and highest degree of measurement error. The results for the rare event scenario are presented in Tables 1–6. The corresponding results for the common event scenario are presented in the Supplementary Web Materials in Tables S1–S6.

The naive estimator was seriously biased in all cases studied, often dramatically. In the single covariate setup, the MS method exhibited low bias across the board, while the RC method often exhibited appreciable bias, especially under the dependent error model, with the bias increasing as the true β increases and as the degree of measurement error increases. In the rare disease case, as expected, the CC method had very high variance, while the variance of the MS method was usually considerably lower. In the common disease case, the MS method had lower variance than the CC method in most configurations, although there are some configurations in which the CC method had lower variance. As expected, Chen's method performed very well in the common disease setup, where the MS method and Chen's method are comparable in terms of bias, variance and coverage probability. In the rare disease setup, Chen's estimator had low bias in some cases and considerable bias in other cases. In addition, the standard deviation of Chen's estimator was substantially greater than that of the MS estimator, in some cases around 3 times greater. Also, the estimate of the standard deviation tended to underestimate, leading to considerably lower than nominal confidence interval coverage rates.

In the multiple-covariate setup, MS method exhibited noticeable bias in some configurations, but the bias with the MS method was typically lower than with the RC method, often considerably so. The patterns were similar across the disease incidence levels (common/rare) and the measurement error models (independent/dependent). The performance of the MS method with dependent covariates was similar to that with independent covariates, and no systematic trends emerged between the dependent covariate case and the independent covariate case in the relative performance of the MS method as

compared with the other methods. Chen's method had a noticeably lower standard deviation than the MS method in the multivariate common disease setting with for large β and moderate correlation between the surrogate and the true exposure (Tables S3–S6 in the Web Appendix, bottom panel). To explore the relative performance of the two methods further, we conducted additional simulations with $e^\beta = 4$ under an “intermediate event rate” scenario with $n = 500$, validation sample size of 200, and a cumulative event rate of about 15% (Table S7 in the Web Appendix) In these simulations, Chen's method again had a noticeably lower standard deviation than the MS method; at the same time, Chen's estimate of the standard deviation of the estimate was noticeably lower than the empirical standard deviation. As a rough practical guide, we suggest that the MS estimator is to be preferred when the number of events in the validation study is very low, while Chen's estimator is to be preferred when the number of events in the validation study is 30 or more, with some caution needed with Chen's estimate of the standard deviation of the estimator.

In both the single-covariate and the multiple-covariate setups, the empirical coverage rate of the asymptotic confidence interval based on the MS method is generally close to the nominal level of 95%, while for the RC method the coverage rate tended to be considerably below nominal for $e^\beta = 4$.

For the multiple-covariate setup, we conducted additional simulations to examine the bias of the MS method for larger sample sizes. These results are reported in the Supplementary Web Materials in Tables S8–S9. When the sample size is increased, the bias decreases, eventually to a very small level.

4. Example

We illustrate the method on data from the Health Professionals Follow-Up Study (HPFS), a prospective cohort study of 51,529 middle-aged (age 40–75 years at baseline) male health professionals. Participants were recruited in 1986 and were mailed questionnaires every other year to assess health status and lifestyle. Here, we analyze the relationship between onset of Type 2 diabetes (T2D) and a diet score relating to intake of carbohydrates, protein, and fat (de Koning et al., 2011). The diet score ranged from 0 to 30, with the score increasing under a decrease in carbohydrate intake or an increase in protein or fat intake. The analysis included the 41,616 study participants who were free of T2D, cardiovascular disease, or cancer at baseline, among whom there were 2,790 cases of incident T2D during follow-up. Diet was assessed with a 131-item semiquantitative food frequency questionnaire (FFQ), an instrument which is subject to substantial measurement error. In a subsample of 105 participants, another diet assessment was carried out using a more accurate diet record (DR). The analysis was stratified by age and adjusted for body mass index (BMI). We analyzed the data using the naive Cox method, the RC method, the complete case method, Chen's method, and our proposed MS method. There were only 6 events among the 105 individuals in the validation sample, which puts Chen's method and the complete case method at a very severe disadvantage. Table 5 presents the results for the various methods. For the regression coefficient for the diet score, the RC estimate is considerably larger than the naive estimate, and the MS and complete case estimates are noticeably larger than the RC estimate. The estimate with Chen's method was lower than that with the naive method.

The standard error with Chen's method was a bit over 1.5 times the standard error with the MS method. For the regression coefficient for BMI, the estimates were similar across all methods, and the standard error with Chen's method was 2.7 times that of the standard error with the MS method.

5. Summary and Discussion

We have developed a new method for covariate error correction in the Cox survival regression model, given internal validation data. The method can handle covariate error of arbitrary form, not just independent additive measurement error. Only a modestly-sized internal validation sample is required. The method can handle the case where the number of covariates is moderate to large. In a simulation study, the method was found to perform very well in terms of bias reduction and confidence interval coverage.

We have worked in the setting of time-independent covariates, but it is possible to consider extension to the case of time-dependent covariates. When the covariate processes are measured on an approximately continuous basis ($\mathbf{W}(t)$ for the full cohort and $\mathbf{X}(t)$ for the internal validation sample), the method and its asymptotic theory carries over with notational changes only. The same is true in the case where the covariate processes are measured only intermittently, as commonly occurs in practice, but the processes vary slowly, so that carrying forward the last observed covariate value is a reasonable approximation.

If the association between \mathbf{W} and \mathbf{X} is very weak, the proposed estimate will remain consistent and asymptotically normal, but the variance will be very high. If there is no association at all between \mathbf{W} and \mathbf{X} , then \mathbf{W} is not a suitable surrogate for \mathbf{X} and no correction method will help. If the relationship between \mathbf{W} and \mathbf{X} is highly nonlinear, the working model (4) can be modified to include nonlinear W terms. A plot of X_{ir} versus W_{ir} for the individuals in the internal validation sample can be used to examine whether nonlinear W terms are needed in the working model for X_{ir} .

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgements

We thank Yi-Hau Chen for sharing with us the code for his method. In addition, we thank the editor, associate editor, and referees for helpful comments that led to substantial improvements in the paper.

References

- Andersen PK, and Gill RD (1982). Cox's regression model for counting processes: a large sample study, *Annals of Statistics* 10: 1100–1120.
- Breslow N and Day NE (1993). *Statistical Methods in Cancer Research Vol. II: The Design and Analysis of Cohort Studies*. Oxford, England: Oxford University Press.
- Carroll RJ, Ruppert D, Stefanski LA, and Crainiceanu CM (2006) *Measurement Error in Nonlinear Models: A Modern Perspective*, 2nd ed Boca Raton: Chapman and Hall / CRC.
- Chen YH (2002). Cox regression in cohort studies with validation sampling. *Journal of the Royal Statistical Society, Series B* 64:51–62.

- Cox DR (1972). Regression models and life tables (with discussion). *Journal of the Royal Statistical Society, Series B* 34:187–200.
- de Koning L, Fung TT, Liao X, Chiuve SE, Rimm EB, Willett WC, Spiegelman D, and Hu FB (2011). Low-carbohydrate diet scores and risk of type 2 diabetes in men. *American Journal of Clinical Nutrition* 93:844–850. [PubMed: 21310828]
- Foutz RV (1977). On the unique consistent solution to the likelihood equations, *Journal of the American Statistical Association* 72:147–148.
- Gill RD (1984). Understanding Cox's regression model: a martingale approach. *Journal of the American Statistical Association* 79:441–447.
- Huang Y and Wang CY (2000). Cox regression with accurate covariates unascertainable: a nonparametric-correction approach, *Journal of the American Statistical Association*, 95:1209–1219.
- Kong FH and Gu M (1999). Consistent estimation in Cox's proportional hazards model with covariate measurement errors, *Statistica Sinica* 9:953–969.
- Kulich M and Lin DY (2000). Additive hazards regression with covariate measurement error. *Journal of the American Statistical Association* 95:238–248.
- Kulich M and Lin DY (2004). Improving the efficiency of relative-risk estimation in case-cohort studies. *Journal of the American Statistical Association* 99:832–844.
- Lin DY and Wei LJ (1989). The robust inference for the Cox proportional hazards model. *Journal of the American Statistical Association* 84:1074–1078.
- Lin DY and Ying Z (1993). Cox regression with incomplete covariate measurements. *Journal of the American Statistical Association* 88:1341–1349.
- Prentice R (1982) Covariate measurement errors and parameter estimation in a failure time regression model. *Biometrika* 69:331–342.
- Spiegelman D, Rosner B and Logan R (2000). Estimation and inference for logistic regression with covariate misclassification and measurement error in main study/validation study designs. *Journal of the American Statistical Association* 95:51–61.
- Thomas DC (1981). General relative-risk models for survival time and matched case-control analysis, *Biometrics* 37:673–686.
- Tsiatis AA (1981). A large sample study of Cox's regression model, *Annals of Statistics* 9:93–108.
- Zhou H and Pepe M (1995). Auxilliary covariate data in failure time regression, *Biometrika* 82:139–149.
- Zhou H and Wang CY (2000). Failure time regression with continuous covariates measured with error, *Journal of the Royal Statistical Society, Series B* 62:657–665.
- Zucker DM (2005). A pseudo partial likelihood method for semi-parametric survival regression with covariate errors. *Journal of the American Statistical Association* 100:1264–1277.
- Zucker DM and Spiegelman D (2004). Inference for the proportional hazards model with misclassified discrete-valued covariates, *Biometrics* 60:324–334. [PubMed: 15180657]
- Zucker DM and Spiegelman D (2008). Corrected score estimation in the proportional hazards model with misclassified discrete covariates. *Statistics in Medicine* 27:1911–1933. [PubMed: 18219700]

Table 1

Simulation results for the single-covariate rare disease case with independent measurement error. β^* is the true value of β . Bias(%) is the relative bias, i.e. $Bias(\%) = 100 \times (\hat{\beta} - \beta^*) / \beta^*$. IQR is 0.74 times the interquartile range of the $\hat{\beta}$ values. SE is the mean of the estimated standard error of $\hat{\beta}$. SD is the empirical standard deviation of the $\hat{\beta}$ values. CR is the empirical coverage rate of the asymptotic 95% confidence interval. Methods considered: MS = modified score, CH = Chen, RC = regression calibration, CC = complete case, NA = naive.

Corr(X, W)	exp(β^*)	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.90	1.5	0.4055	MS	0.4050	-0.1	0.4011	-1.1	0.0577	0.0525	0.0517	0.965
			GH	0.4230	4.3	0.4313	6.4	0.1621	0.1373	0.1743	0.879
			RG	0.4036	-0.5	0.4032	-0.6	0.0562	0.0511	0.0506	0.957
			GG	0.4188	3.3	0.4163	2.7	0.3120	0.3384	0.3684	0.945
			NA	0.3287	-18.9	0.3309	-18.4	0.0404	0.0402	0.0386	0.543
0.70	1.5	0.4055	MS	0.4088	0.8	0.4040	-0.4	0.0737	0.0738	0.0753	0.945
			GH	0.4277	5.5	0.4403	8.6	0.2545	0.2126	0.2979	0.855
			RG	0.4029	-0.6	0.4032	-0.6	0.0704	0.0688	0.0690	0.965
			GG	0.4188	3.3	0.4163	2.7	0.3120	0.3384	0.3684	0.945
			NA	0.1993	-50.9	0.2011	-50.4	0.0346	0.0313	0.0306	0.000
0.50	1.5	0.4055	MS	0.4129	1.8	0.4103	1.2	0.1081	0.1099	0.1186	0.938
			GH	0.4326	6.7	0.4459	10.0	0.3006	0.2518	0.3558	0.859
			RG	0.4030	-0.6	0.4004	-1.3	0.1052	0.0976	0.1011	0.938
			GG	0.4188	3.3	0.4163	2.7	0.3120	0.3384	0.3684	0.945
			NA	0.1022	-74.8	0.1036	-74.4	0.0237	0.0224	0.0223	0.000
0.90	2.5	0.9163	MS	0.9279	1.3	0.9221	0.6	0.0750	0.0720	0.0721	0.949
			GH	0.9401	2.6	0.9289	1.4	0.1857	0.1599	0.2120	0.875
			RG	0.9098	-0.7	0.9045	-1.3	0.0619	0.0584	0.0575	0.973
			GG	0.9380	2.4	0.9259	1.0	0.3401	0.3590	0.4109	0.941
			NA	0.7412	-19.1	0.7406	-19.2	0.0393	0.0409	0.0395	0.004
0.70	2.5	0.9163	MS	0.9449	3.1	0.9376	2.3	0.1269	0.1279	0.1324	0.957
			GH	0.9545	4.2	0.9511	3.8	0.2756	0.2394	0.3477	0.867
			RG	0.8944	-2.4	0.8910	-2.8	0.0904	0.0878	0.0845	0.949
			GG	0.9380	2.4	0.9259	1.0	0.3401	0.3590	0.4109	0.941
			NA	0.4434	-51.6	0.4438	-51.6	0.0272	0.0315	0.0307	0.000
0.50	2.5	0.9163	MS	0.9620	5.0	0.9460	3.2	0.2069	0.2263	0.2401	0.957
			GH	0.9577	4.5	0.9601	4.8	0.3152	0.2761	0.4026	0.855
			RG	0.8785	-4.1	0.8766	-4.3	0.1345	0.1263	0.1258	0.914
			GG	0.9380	2.4	0.9259	1.0	0.3401	0.3590	0.4109	0.941
			NA	0.2254	-75.4	0.2270	-75.2	0.0191	0.0224	0.0223	0.000
0.90	4.0	1.3863	MS	1.4214	2.5	1.4080	1.6	0.1166	0.1162	0.1196	0.930
			GH	1.4359	3.6	1.4159	2.1	0.2352	0.2004	0.2625	0.875

Corr(X, W)	$\exp(\beta^*)$	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.70	4.0	1.3863	RG	1.3464	-2.9	1.3476	-2.8	0.0718	0.0687	0.0651	0.906
			GG	1.4460	4.3	1.4134	2.0	0.3881	0.4063	0.4590	0.957
			NA	1.0967	-20.9	1.0997	-20.7	0.0449	0.0426	0.0422	0.000
			MS	1.4862	7.2	1.4587	5.2	0.2271	0.2405	0.2590	0.957
			GH	1.4654	5.7	1.4196	2.4	0.3281	0.2837	0.3986	0.856
			RG	1.2863	-7.2	1.2901	-6.9	0.1112	0.1079	0.1020	0.781
0.50	4.0	1.3863	GG	1.4460	4.3	1.4134	2.0	0.3881	0.4063	0.4590	0.957
			NA	0.6384	-53.9	0.6388	-53.9	0.0337	0.0319	0.0312	0.000
			MS	1.4992	8.1	1.4446	4.2	0.3384	0.3698	0.3786	0.944
			GH	1.4752	6.4	1.4155	2.1	0.3636	0.3168	0.4497	0.863
			RG	1.2358	-10.9	1.2302	-11.3	0.1650	0.1496	0.1490	0.739
			GG	1.4460	4.3	1.4134	2.0	0.3881	0.4063	0.4590	0.957
			NA	0.3206	-76.9	0.3228	-76.7	0.0227	0.0225	0.0221	0.000

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 2

Simulation results for the single-covariate rare disease case with dependent measurement error. β^* is the true value of β . Bias(%) is the relative bias, i.e. $Bias(\%) = 100 \times (\hat{\beta} - \beta^*)/\beta^*$. IQR is 0.74 times the interquartile range of the $\hat{\beta}$ values. SE is the mean of the estimated standard error of $\hat{\beta}$. SD is the empirical standard deviation of the $\hat{\beta}$ values. CR is the empirical coverage rate of the asymptotic 95% confidence interval. Methods considered: MS = modified score, CH = Chen, RC = regression calibration, CC = complete case, NA = naive.

Corr(X, W)	exp(β^*)	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.90	1.5	0.4055	MS	0.4043	-0.3	0.4012	-1.1	0.0558	0.0525	0.0516	0.957
			GH	0.4255	4.9	0.4314	6.4	0.1628	0.1345	0.1729	0.871
			RG	0.4005	-1.2	0.4009	-1.1	0.0517	0.0499	0.0497	0.949
			GG	0.4188	3.3	0.4163	2.7	0.3120	0.3384	0.3684	0.945
			NA	0.3299	-18.6	0.3315	-18.3	0.0392	0.0400	0.0382	0.531
0.70	1.5	0.4055	MS	0.4071	0.4	0.4055	0.0	0.0739	0.0730	0.0729	0.953
			GH	0.4272	5.4	0.4406	8.7	0.2622	0.2126	0.3012	0.867
			RG	0.3992	-1.5	0.3980	-1.8	0.0694	0.0669	0.0667	0.957
			GG	0.4188	3.3	0.4163	2.7	0.3120	0.3384	0.3684	0.945
			NA	0.1974	-51.3	0.1985	-51.0	0.0322	0.0308	0.0301	0.000
0.50	1.5	0.4055	MS	0.4156	2.5	0.4111	1.4	0.1020	0.1122	0.1175	0.953
			GH	0.4267	5.2	0.4287	5.7	0.3036	0.2509	0.3561	0.855
			RG	0.4016	-1.0	0.3995	-1.5	0.1007	0.0974	0.0995	0.949
			GG	0.4188	3.3	0.4163	2.7	0.3120	0.3384	0.3684	0.945
			NA	0.1023	-74.8	0.1042	-74.3	0.0224	0.0223	0.0224	0.000
0.90	2.5	0.9163	MS	0.9226	0.7	0.9091	-0.8	0.0763	0.0840	0.0801	0.949
			GH	0.9386	2.4	0.9283	1.3	0.1835	0.1623	0.2210	0.859
			RG	0.8798	-4.0	0.8778	-4.2	0.0579	0.0550	0.0552	0.887
			GG	0.9380	2.4	0.9259	1.0	0.3401	0.3590	0.4109	0.941
			NA	0.7247	-20.9	0.7229	-21.1	0.0354	0.0392	0.0376	0.000
0.70	2.5	0.9163	MS	0.9415	2.8	0.9221	0.6	0.1341	0.1433	0.1346	0.961
			GH	0.9492	3.6	0.9506	3.7	0.2732	0.2431	0.3613	0.867
			RG	0.8631	-5.8	0.8669	-5.4	0.0861	0.0807	0.0779	0.879
			GG	0.9380	2.4	0.9259	1.0	0.3401	0.3590	0.4109	0.941
			NA	0.4268	-53.4	0.4264	-53.5	0.0261	0.0297	0.0292	0.000
0.50	2.5	0.9163	MS	0.9616	4.9	0.9365	2.2	0.2144	0.2861	0.2327	0.953
			GH	0.9367	2.2	0.9367	2.2	0.3046	0.2777	0.3913	0.871
			RG	0.8675	-5.3	0.8698	-5.1	0.1360	0.1257	0.1246	0.902
			GG	0.9380	2.4	0.9259	1.0	0.3401	0.3590	0.4109	0.941
			NA	0.2230	-75.7	0.2246	-75.5	0.0202	0.0219	0.0226	0.000
0.90	4.0	1.3863	MS	1.4056	1.4	1.3595	-1.9	0.1087	0.2276	0.2203	0.928
			GH	1.4280	3.0	1.4047	1.3	0.2489	0.2089	0.2907	0.883

Corr(X, W)	$\exp(\beta^*)$	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.70	4.0	1.3863	RG	1.2652	-8.7	1.2619	-9.0	0.0659	0.0635	0.0608	0.508
			GG	1.4460	4.3	1.4134	2.0	0.3881	0.4063	0.4590	0.957
			NA	1.0416	-24.9	1.0429	-24.8	0.0417	0.0392	0.0386	0.000
			MS	1.4559	5.0	1.4036	1.2	0.2359	0.3021	0.2920	0.939
			GH	1.4517	4.7	1.4043	1.3	0.3466	0.2888	0.4162	0.859
			RG	1.2086	-12.8	1.2094	-12.8	0.0963	0.0962	0.0906	0.535
0.50	4.0	1.3863	GG	1.4460	4.3	1.4134	2.0	0.3881	0.4063	0.4590	0.957
			NA	0.5969	-56.9	0.5988	-56.8	0.0304	0.0288	0.0284	0.000
			MS	1.4663	5.8	1.4039	1.3	0.3186	0.3821	0.3793	0.934
			GH	1.4564	5.1	1.3922	0.4	0.3705	0.3192	0.4530	0.856
			RG	1.2105	-12.7	1.1978	-13.6	0.1572	0.1490	0.1478	0.696
			GG	1.4460	4.3	1.4134	2.0	0.3881	0.4063	0.4590	0.957
			NA	0.3135	-77.4	0.3149	-77.3	0.0239	0.0214	0.0221	0.000

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 3

Simulation results for the multiple-covariate rare disease case with independent covariates and independent measurement error. β^* is the true value of β . Bias(%) is the relative bias, i.e. $Bias(\%) = 100 \times (\hat{\beta} - \beta^*)/\beta^*$. IQR is 0.74 times the interquartile range of the $\hat{\beta}$ values. SE is the mean of the estimated standard error of $\hat{\beta}$. SD is the empirical standard deviation of the $\hat{\beta}$ values. CR is the empirical coverage rate of the asymptotic 95% confidence interval. Methods considered: MS = modified score, CH = Chen, RC = regression calibration, CC = complete case, NA = naive.

Corr(X, W)	exp(β^*)	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.90	1.5	0.4055	MS	0.4118	1.6	0.4101	1.1	0.0512	0.0684	0.0573	0.949
			GH	0.4136	2.0	0.4106	1.3	0.1962	0.1352	0.2432	0.772
			RG	0.4074	0.5	0.4063	0.2	0.0486	0.0523	0.0507	0.945
			GG	0.4717	16.3	0.4533	11.8	0.3626	0.3782	0.4538	0.961
			NA	0.3321	-18.1	0.3337	-17.7	0.0420	0.0410	0.0408	0.567
0.70	1.5	0.4055	MS	0.4175	3.0	0.4067	0.3	0.0761	0.0860	0.0814	0.957
			GH	0.4292	5.8	0.4518	11.4	0.3221	0.2026	0.4050	0.749
			RG	0.4066	0.3	0.4074	0.5	0.0700	0.0709	0.0684	0.949
			GG	0.4717	16.3	0.4533	11.8	0.3626	0.3782	0.4538	0.961
			NA	0.1997	-50.8	0.2017	-50.2	0.0333	0.0319	0.0320	0.000
0.50	1.5	0.4055	MS	0.4300	6.1	0.4148	2.3	0.1105	0.1402	0.1340	0.952
			GH	0.4274	5.4	0.4542	12.0	0.3518	0.2358	0.4852	0.749
			RG	0.4081	0.6	0.4065	0.3	0.1050	0.1016	0.0986	0.957
			GG	0.4717	16.3	0.4533	11.8	0.3626	0.3782	0.4538	0.961
			NA	0.1013	-75.0	0.1021	-74.8	0.0248	0.0228	0.0228	0.000
0.90	2.5	0.9163	MS	0.9429	2.9	0.9289	1.4	0.1076	0.1272	0.1230	0.937
			GH	0.9757	6.5	0.9480	3.5	0.2110	0.1596	0.2607	0.785
			RG	0.9109	-0.6	0.9133	-0.3	0.0536	0.0595	0.0569	0.961
			GG	1.0757	17.4	1.0305	12.5	0.3716	0.4151	0.4676	0.945
			NA	0.7424	-19.0	0.7429	-18.9	0.0441	0.0415	0.0418	0.016
0.70	2.5	0.9163	MS	0.9657	5.4	0.9398	2.6	0.1901	0.2015	0.2146	0.955
			GH	1.0211	11.4	0.9778	6.7	0.3117	0.2287	0.4124	0.754
			RG	0.8962	-2.2	0.8896	-2.9	0.0883	0.0901	0.0865	0.930
			GG	1.0757	17.4	1.0305	12.5	0.3716	0.4151	0.4676	0.945
			NA	0.4408	-51.9	0.4428	-51.7	0.0311	0.0318	0.0333	0.000
0.50	2.5	0.9163	MS	1.0264	12.0	0.9356	2.1	0.3004	0.3131	0.3227	0.938
			GH	1.0330	12.7	1.0100	10.2	0.3453	0.2602	0.4751	0.762
			RG	0.8868	-3.2	0.8738	-4.6	0.1235	0.1308	0.1308	0.930
			GG	1.0757	17.4	1.0305	12.5	0.3716	0.4151	0.4676	0.945
			NA	0.2228	-75.7	0.2244	-75.5	0.0221	0.0226	0.0236	0.000
0.90	4.0	1.3863	MS	1.4395	3.8	1.4175	2.3	0.1245	0.1333	0.1490	0.943
			GH	1.5051	8.6	1.4591	5.2	0.3151	0.2087	0.4670	0.769

Corr(X, W)	$\exp(\beta^*)$	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.70	4.0	1.3863	RG	1.3467	-2.9	1.3496	-2.6	0.0732	0.0705	0.0683	0.902
			GG	1.6563	19.5	1.5847	14.3	0.4873	0.5055	0.5971	0.945
			NA	1.0974	-20.8	1.0969	-20.9	0.0440	0.0438	0.0453	0.000
			MS	1.5253	10.0	1.4410	3.9	0.3180	0.3329	0.3486	0.936
			GH	1.5592	12.5	1.5045	8.5	0.4789	0.2822	0.5824	0.741
			RG	1.2853	-7.3	1.2782	-7.8	0.1074	0.1110	0.1091	0.805
0.50	4.0	1.3863	GG	1.6563	19.5	1.5847	14.3	0.4873	0.5055	0.5971	0.945
			NA	0.6327	-54.4	0.6349	-54.2	0.0354	0.0326	0.0347	0.000
			MS	1.5407	11.1	1.4472	4.4	0.4363	0.4511	0.4404	0.925
			GH	1.5903	14.7	1.5255	10.0	0.5249	0.3138	0.6286	0.706
			RG	1.2423	-10.4	1.2252	-11.6	0.1481	0.1546	0.1546	0.789
			GG	1.6563	19.5	1.5847	14.3	0.4873	0.5055	0.5971	0.945
			NA	0.3156	-77.2	0.3164	-77.2	0.0238	0.0229	0.0236	0.000

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 4

Simulation results for the multiple-covariate rare disease case with independent covariates and dependent measurement error. β^* is the true value of β . Bias(%) is the relative bias, i.e. $Bias(\%) = 100 \times (\hat{\beta} - \beta^*)/\beta^*$. IQR is 0.74 times the interquartile range of the $\hat{\beta}$ values. SE is the mean of the estimated standard error of $\hat{\beta}$. SD is the empirical standard deviation of the $\hat{\beta}$ values. CR is the empirical coverage rate of the asymptotic 95% confidence interval. Methods considered: MS = modified score, CH = Chen, RC = regression calibration, CC = complete case, NA = naive.

Corr(X, W)	exp(β^*)	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.90	1.5	0.4055	MS	0.4098	1.1	0.4057	0.1	0.0558	0.0564	0.0546	0.948
			GH	0.4157	2.5	0.4113	1.4	0.2128	0.1334	0.2410	0.760
			RG	0.4045	-0.2	0.4017	-0.9	0.0482	0.0511	0.0504	0.949
			GG	0.4717	16.3	0.4533	11.8	0.3626	0.3782	0.4538	0.961
			NA	0.3339	-17.6	0.3367	-16.9	0.0417	0.0408	0.0406	0.571
0.70	1.5	0.4055	MS	0.4141	2.1	0.4062	0.2	0.0735	0.0833	0.0800	0.957
			GH	0.4250	4.8	0.4472	10.3	0.3061	0.2026	0.4028	0.733
			RG	0.4036	-0.5	0.4043	-0.3	0.0710	0.0690	0.0678	0.953
			GG	0.4717	16.3	0.4533	11.8	0.3626	0.3782	0.4538	0.961
			NA	0.1984	-51.1	0.1994	-50.8	0.0331	0.0315	0.0320	0.000
0.50	1.5	0.4055	MS	0.4312	6.4	0.4123	1.7	0.1164	0.1460	0.1453	0.957
			GH	0.4299	6.0	0.4346	7.2	0.3686	0.2346	0.4725	0.753
			RG	0.4084	0.7	0.4070	0.4	0.1048	0.1017	0.1004	0.953
			GG	0.4717	16.3	0.4533	11.8	0.3626	0.3782	0.4538	0.961
			NA	0.1019	-74.9	0.1031	-74.6	0.0260	0.0227	0.0233	0.000
0.90	2.5	0.9163	MS	0.9414	2.7	0.9183	0.2	0.0839	0.1050	0.1148	0.956
			GH	0.9741	6.3	0.9437	3.0	0.2253	0.1617	0.2700	0.782
			RG	0.8828	-3.7	0.8821	-3.7	0.0492	0.0562	0.0548	0.891
			GG	1.0757	17.4	1.0305	12.5	0.3716	0.4151	0.4676	0.945
			NA	0.7284	-20.5	0.7276	-20.6	0.0411	0.0399	0.0400	0.004
0.70	2.5	0.9163	MS	0.9504	3.7	0.9223	0.7	0.1275	0.1366	0.1429	0.935
			GH	1.0096	10.2	0.9775	6.7	0.3164	0.2295	0.4197	0.754
			RG	0.8686	-5.2	0.8685	-5.2	0.0795	0.0833	0.0803	0.883
			GG	1.0757	17.4	1.0305	12.5	0.3716	0.4151	0.4676	0.945
			NA	0.4271	-53.4	0.4276	-5.3	0.0298	0.0302	0.0318	0.000
0.50	2.5	0.9163	MS	1.0023	9.4	0.9252	1.0	0.2244	0.2382	0.2352	0.928
			GH	1.0173	11.0	0.9848	7.5	0.3487	0.2594	0.4739	0.750
			RG	0.8800	-4.0	0.8685	-5.2	0.1307	0.1311	0.1323	0.906
			GG	1.0757	17.4	1.0305	12.5	0.3716	0.4151	0.4676	0.945
			NA	0.2219	-75.8	0.2228	-75.7	0.0235	0.0221	0.0238	0.000
0.90	4.0	1.3863	MS	1.4155	2.1	1.3822	-0.3	0.1549	0.1720	0.1881	0.942
			GH	1.4825	6.9	1.4324	3.3	0.3299	0.2126	0.4456	0.764

Corr(X, W)	$\exp(\beta^*)$	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.70	4.0	1.3863	RG	1.2701	-8.4	1.2703	-8.4	0.0645	0.0654	0.0639	0.571
			GG	1.6563	19.5	1.5847	14.3	0.4873	0.5055	0.5971	0.945
			NA	1.0474	-24.4	1.0466	-24.5	0.0405	0.0405	0.0424	0.000
			MS	1.4339	3.4	1.3776	-0.6	0.3308	0.3563	0.3666	0.930
			GH	1.5297	10.3	1.4711	6.1	0.4686	0.2837	0.5652	0.732
			RG	1.2156	-12.3	1.2144	-12.4	0.1054	0.0997	0.0980	0.567
0.50	4.0	1.3863	GG	1.6563	19.5	1.5847	14.3	0.4873	0.5055	0.5971	0.945
			NA	0.5969	-56.9	0.6004	-56.7	0.0320	0.0297	0.0320	0.000
			MS	1.4985	8.1	1.3904	0.3	0.4507	0.4740	0.4678	0.926
			GH	1.5673	13.1	1.5050	8.6	0.5324	0.3130	0.6322	0.710
			RG	1.2240	-11.7	1.2001	-13.4	0.1560	0.1549	0.1561	0.758
			GG	1.6563	19.5	1.5847	14.3	0.4873	0.5055	0.5971	0.945
			NA	0.3112	-77.6	0.3120	-77.5	0.0249	0.0220	0.0234	0.000

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 5

Simulation results for the multiple-covariate rare disease case with dependent covariates and independent measurement error. β^* is the true value of β . Bias(%) is the relative bias, i.e. $Bias(\%) = 100 \times (\hat{\beta} - \beta^*)/\beta^*$. IQR is 0.74 times the interquartile range of the $\hat{\beta}$ values. SE is the mean of the estimated standard error of $\hat{\beta}$. SD is the empirical standard deviation of the $\hat{\beta}$ values. CR is the empirical coverage rate of the asymptotic 95% confidence interval. Methods considered: MS = modified score, CH = Chen, RC = regression calibration, CC = complete case, NA = naive.

Corr(X, W)	exp(β^*)	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.90	1.5	0.4055	MS	0.4095	1.0	0.4045	-0.2	0.0526	0.0570	0.0589	0.957
			GH	0.4199	3.6	0.4149	2.3	0.1793	0.1364	0.2120	0.825
			RG	0.3954	-2.5	0.3916	-3.4	0.0503	0.0482	0.0510	0.941
			GG	0.4373	7.8	0.4208	3.8	0.3744	0.3529	0.3891	0.961
			NA	0.3221	-20.6	0.3215	-20.7	0.0393	0.0378	0.0398	0.375
0.70	1.5	0.4055	MS	0.4139	2.1	0.4063	0.2	0.0787	0.0902	0.0883	0.949
			GH	0.4355	7.4	0.4363	7.6	0.2904	0.2044	0.3414	0.774
			RG	0.3768	-7.1	0.3721	-8.2	0.0596	0.0638	0.0657	0.910
			GG	0.4373	7.8	0.4208	3.8	0.3744	0.3529	0.3891	0.961
			NA	0.1861	-54.1	0.1840	-54.6	0.0304	0.0288	0.0305	0.000
0.50	1.5	0.4055	MS	0.4320	6.5	0.3959	-2.4	0.1214	0.1533	0.1429	0.941
			GH	0.4391	8.3	0.4414	8.9	0.3502	0.2381	0.3977	0.770
			RG	0.3603	-11.1	0.3601	-11.2	0.0824	0.0885	0.0870	0.906
			GG	0.4373	7.8	0.4208	3.8	0.3744	0.3529	0.3891	0.961
			NA	0.0916	-77.4	0.0905	-77.7	0.0204	0.0203	0.0212	0.000
0.90	2.5	0.9163	MS	0.9375	2.3	0.9201	0.4	0.0837	0.1190	0.1109	0.937
			GH	0.9507	3.8	0.9219	0.6	0.1921	0.1569	0.2324	0.840
			RG	0.8767	-4.3	0.8719	-4.8	0.0541	0.0543	0.0581	0.875
			GG	0.9886	7.9	0.9803	7.0	0.3273	0.3707	0.4115	0.961
			NA	0.7140	-22.1	0.7141	-22.1	0.0453	0.0373	0.0395	0.000
0.70	2.5	0.9163	MS	0.9738	6.3	0.9403	2.6	0.1825	0.2151	0.2035	0.934
			GH	0.9886	7.9	0.9670	5.5	0.2972	0.2296	0.3536	0.809
			RG	0.8191	-10.6	0.8130	-11.3	0.0862	0.0801	0.0827	0.699
			GG	0.9886	7.9	0.9803	7.0	0.3273	0.3707	0.4115	0.961
			NA	0.4042	-55.9	0.4042	-55.9	0.0306	0.0280	0.0297	0.000
0.50	2.5	0.9163	MS	1.0112	10.4	0.9172	0.1	0.2953	0.3348	0.3250	0.928
			GH	0.9965	8.8	0.9604	4.8	0.3475	0.2607	0.4162	0.805
			RG	0.7736	-15.6	0.7737	-15.6	0.1150	0.1118	0.1099	0.683
			GG	0.9886	7.9	0.9803	7.0	0.3273	0.3707	0.4115	0.961
			NA	0.1975	-78.4	0.1969	-78.5	0.0216	0.0196	0.0206	0.000
0.90	4.0	1.3863	MS	1.4361	3.6	1.3980	0.8	0.1537	0.2373	0.1928	0.938
			GH	1.4638	5.6	1.3997	1.0	0.2668	0.2105	0.3759	0.824

Corr(X, W)	$\exp(\beta^*)$	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.70	4.0	1.3863	RG	1.2871	-7.2	1.2805	-7.6	0.0675	0.0648	0.0665	0.606
			GG	1.6191	16.8	1.5145	9.2	0.4535	0.4886	0.7240	0.969
			NA	1.0479	-24.4	1.0419	-24.8	0.0425	0.0392	0.0421	0.000
			MS	1.5063	8.7	1.4175	2.3	0.2952	0.3454	0.3020	0.920
			GH	1.5515	11.9	1.4604	5.3	0.3985	0.2887	0.5080	0.807
			RG	1.1605	-16.3	1.1516	-16.9	0.1026	0.0986	0.0994	0.383
0.50	4.0	1.3863	GG	1.6191	16.8	1.5145	9.2	0.4535	0.4886	0.7240	0.969
			NA	0.5725	-58.7	0.5714	-58.8	0.0348	0.0285	0.0310	0.000
			MS	1.4931	7.7	1.4066	1.5	0.3796	0.4436	0.4172	0.918
			GH	1.5455	11.5	1.4421	4.0	0.4568	0.3185	0.5594	0.792
			RG	1.0735	-22.6	1.0742	-22.5	0.1330	0.1334	0.1296	0.367
			GG	1.6191	16.8	1.5145	9.2	0.4535	0.4886	0.7240	0.969
			NA	0.2753	-80.1	0.2737	-80.3	0.0256	0.0197	0.0209	0.000

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 6

Simulation results for the multiple-covariate rare disease case with with dependent covariates and dependent measurement error. β^* is the true value of β . Bias(%) is the relative bias, i.e. $Bias(\%) = 100 \times (\hat{\beta} - \beta^*) / \beta^*$. IQR is 0.74 times the interquartile range of the $\hat{\beta}$ values. SE is the mean of the estimated standard error of $\hat{\beta}$. SD is the empirical standard deviation of the $\hat{\beta}$ values. CR is the empirical coverage rate of the asymptotic 95% confidence interval. Methods considered: MS = modified score, CH = Chen, RC = regression calibration, CC = complete case, NA = naive.

Corr(X, W)	exp(β^*)	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.90	1.5	0.4055	MS	0.4132	1.9	0.4034	-0.5	0.0535	0.0638	0.0717	0.949
			GH	0.4256	5.0	0.4184	3.2	0.1716	0.1510	0.2117	0.840
			RG	0.3882	-4.3	0.3845	-5.2	0.0475	0.0467	0.0494	0.930
			GG	0.4373	7.8	0.4208	3.8	0.3744	0.3529	0.3891	0.961
			NA	0.3201	-21.1	0.3190	-21.3	0.0385	0.0373	0.0389	0.363
0.70	1.5	0.4055	MS	0.4097	1.0	0.3965	-2.2	0.0741	0.0937	0.0875	0.953
			GH	0.4247	4.7	0.3979	-1.9	0.3105	0.2402	0.3989	0.832
			RG	0.3656	-9.8	0.3627	-10.5	0.0552	0.0612	0.0621	0.875
			GG	0.4373	7.8	0.4208	3.8	0.3744	0.3529	0.3891	0.961
			NA	0.1804	-55.5	0.1782	-56.0	0.0308	0.0280	0.0294	0.000
0.50	1.5	0.4055	MS	0.4249	4.8	0.3938	-2.9	0.1226	0.1663	0.1591	0.936
			GH	0.3966	-2.2	0.3846	-5.2	0.3626	0.2814	0.4833	0.816
			RG	0.3519	-13.2	0.3498	-13.7	0.0848	0.0876	0.0842	0.890
			GG	0.4373	7.8	0.4208	3.8	0.3744	0.3529	0.3891	0.961
			NA	0.0897	-77.9	0.0885	-78.2	0.0226	0.0200	0.0209	0.000
0.90	2.5	0.9163	MS	0.9330	1.8	0.9018	-1.6	0.0899	0.1134	0.1229	0.929
			GH	0.9425	2.9	0.9136	-0.3	0.2188	0.1830	0.2600	0.871
			RG	0.8400	-8.3	0.8350	-8.9	0.0495	0.0512	0.0541	0.625
			GG	0.9886	7.9	0.9803	7.0	0.3273	0.3707	0.4115	0.961
			NA	0.6922	-24.5	0.6923	-24.4	0.0386	0.0356	0.0370	0.000
0.70	2.5	0.9163	MS	0.9499	3.7	0.9150	-0.1	0.1512	0.1800	0.1843	0.921
			GH	0.9770	6.6	0.9371	2.3	0.3183	0.2731	0.4340	0.855
			RG	0.7789	-15.0	0.7735	-15.6	0.0799	0.0733	0.0742	0.484
			GG	0.9886	7.9	0.9803	7.0	0.3273	0.3707	0.4115	0.961
			NA	0.3834	-58.2	0.3838	-58.1	0.0297	0.0262	0.0273	0.000
0.50	2.5	0.9163	MS	0.9763	6.5	0.8980	-2.0	0.2604	0.3032	0.2847	0.896
			GH	0.9659	5.4	0.9285	1.3	0.3232	0.3116	0.4971	0.863
			RG	0.7554	-17.6	0.7568	-17.4	0.1147	0.1110	0.1063	0.641
			GG	0.9886	7.9	0.9803	7.0	0.3273	0.3707	0.4115	0.961
			NA	0.1929	-79.0	0.1929	-78.9	0.0214	0.0189	0.0200	0.000
0.90	4.0	1.3863	MS	1.4213	2.5	1.3515	-2.5	0.1523	0.1813	0.1713	0.870
			GH	1.4668	5.8	1.4018	1.1	0.2839	0.2350	0.4104	0.832

Corr(X, W)	$\exp(\beta^*)$	β^*	Method	Mean		Median		IQR	SE	SD	CR
				$\hat{\beta}$	Bias(%)	$\hat{\beta}$	Bias(%)				
0.70	4.0	1.3863	RG	1.2054	-13.1	1.1996	-13.5	0.0561	0.0602	0.0609	0.160
			GG	1.6191	16.8	1.5145	9.2	0.4535	0.4886	0.7240	0.969
			NA	0.9925	-28.4	0.9893	-28.6	0.0387	0.0362	0.0386	0.000
			MS	1.4304	3.2	1.3619	-1.8	0.2769	0.3131	0.3011	0.870
			GH	1.5199	9.6	1.4261	2.9	0.3573	0.3232	0.5417	0.848
			RG	1.0864	-21.6	1.0744	-22.5	0.0961	0.0883	0.0886	0.129
0.50	4.0	1.3863	GG	1.6191	16.8	1.5145	9.2	0.4535	0.4886	0.7240	0.969
			NA	0.5337	-61.5	0.5336	-61.5	0.0309	0.0259	0.0281	0.000
			MS	1.4489	4.5	1.3289	-4.1	0.3643	0.4592	0.4374	0.871
			GH	1.5204	9.7	1.4437	4.1	0.3662	0.3563	0.5615	0.848
			RG	1.0487	-24.4	1.0493	-24.3	0.1392	0.1332	0.1277	0.324
			GG	1.6191	16.8	1.5145	9.2	0.4535	0.4886	0.7240	0.969
			NA	0.2686	-80.6	0.2676	-80.7	0.0230	0.0188	0.0204	0.000

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 7

HPFS Results. SE = standard error of estimate. SE Ratio = Ratio between the standard error of the estimate and the standard error of the modified score estimate. Methods considered: MS = modified score, CH = Chen, RC = regression calibration, CC = complete case, NA = naive.

Method	Diet Score Coefficient			BMI Coefficient		
	Estimate	SE	SE Ratio	Estimate	SE	SE Ratio
Naive	0.0216	0.0027	0.1107	0.0867	0.0019	0.2346
CC	0.0788	0.0738	3.0246	0.0913	0.1335	16.4815
RC	0.0485	0.0096	0.3934	0.0867	0.0078	0.9630
CH	0.0136	0.0383	1.5697	0.0800	0.0220	2.7160
MS	0.0712	0.0244	1.0000	0.0865	0.0081	1.0000

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript