# Survival Analysis with Covariate Measurement Error

By

Ming Zhu

Doctor of Philosophy

Biostatistics

---

Yijian Huang, Ph.D.

Advisor

---

Nancy G. Kutner, Ph.D.

Committee Member

---

Brent A. Johnson, Ph.D.

Committee Member

---

Robert H. Lyles, Ph.D.

Committee Member

Accepted:

---

Lisa A. Tedesco, Ph.D.

Dean of the James T. Laney School of Graduate Studies

---

Date

**Survival Analysis with Covariate Measurement Error**

By

Ming Zhu

B.S., University of Science and Technology of China, 2002

M.S., National University of Singapore, 2004

Advisor: Yijian Huang, Ph.D.

An abstract of

A dissertation submitted to the Faculty of the

James T. Laney School of Graduate Studies of Emory University

in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

in Biostatistics

2011

# Abstract

In many medical research studies, survival time is typically the primary outcome of interest. The Cox proportional hazards model and the accelerated failure time model are two popular methods to investigate the relationship between covariates and possibly right-censored survival time. However, in many clinical trials, the true covariates may not always be accurately measured due to natural biological fluctuation or instrument error. For regression analysis in general, naively using mismeasured covariates in conventional inference procedures may incur substantial estimation bias. In this dissertation research, we aim to resolve several issues in survival data analysis with covariate measurement error, with particular emphasis on the aforementioned two models.

The first topic focuses on the analysis of recurrent events data, in which the event of interest may occur more than once for each subject during the follow up period. We proposes an estimation procedure for recurrent events data under the accelerated failure time model in which some covariates are not accurately measured. With replicated mismeasured covariates available, the proposed estimation procedure requires no distributional assumptions on either the true covariates or the error except for the boundedness of the latter. The resulting regression coefficient estimators are shown to be consistent and asymptotically normal. The performance of the proposed procedure is investigated by extensive numerical studies with practical sample size. In addition, an application to data from a clinical trial is provided to illustrate the proposed method.

The second topic considers the Cox proportional hazards model for univariate survival data. In the presence of covariate measurement error, several functional modeling methods have been proposed under the situation where the distribution of the measurement error is known. Among them are parametric corrected score (Nakamura, 1992) and conditional score (Tsiatis & Davidian, 2001). Although both methods are consistent, each suffers from severe problem of multiple roots or absence of appropriate root when the measurement error is substantial. The problem persists even when the sample size is practically large. We conduct a detailed investigation on the pathological behaviors of parametric corrected score and propose an approach of incorporating additional estimating functions to remedy these pathological behaviors. The estimation and inference are then accomplished by means of quadratic inference function. Extensive simulation studies are conducted to evaluate the performance of proposed method.

In the third topic, we consider the Cox proportional hazards model with covariate measurement error where the error distribution is completely unspecified, but replicated mismeasured covariates are available instead. A consistent nonparametric corrected score (Huang & Wang, 2000) has been proposed for Cox proportional hazards model with replicated mismeasured covariates. But it also suffers from finite-

sample pathological behaviors similar to that of the parametric corrected score when the measurement error is substantial. To address this issue, we develop a similar technique as in the second topic for the nonparametric corrected score and evaluate its performance by extensive simulation study.

# Survival Analysis with Covariate Measurement Error

By

Ming Zhu

B.S., University of Science and Technology of China, 2002

M.S., National University of Singapore, 2004

Advisor: Yijian Huang, Ph.D.

A dissertation submitted to the Faculty of the

James T. Laney School of Graduate Studies of Emory University

in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

in Biostatistics

2011

# Acknowledgement

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

## 1.1   Background

In many clinical research studies, survival time is typically the primary outcome of interest. The proportional hazards model (Cox, 1972) is one of the most popular methods to investigate the relationship between covariates and possibly right-censored survival times. An important alternative to Cox model is the accelerated failure time (AFT) model. In conventional statistical inference procedures, it is common to assume that all covariates are observed and measured accurately without any error. However, in medical research fields, the concerns of measurement errors often arise because of natural biological fluctuation or instrument error. For regression analysis in general, it is well known that naively using mismeasured covariates in conventional inference procedures may incur substantial estimation bias and several statistical methods have been suggested to address covariate measurement error; see Section 1.3

for a brief review and the monograph of Carroll et al. (2006) for a more comprehensive summary.

The goal of this dissertation research is to develop statistical methods to address several issues on covariate measurement error in survival data analysis. We focus on two types of data: recurrent events data and univariate survival data.

The first topic studies the recurrent events data under the accelerated failure time model. The accelerated failure time model is an important alternative to the Cox model and is appealing in that the model has a direct physical interpretation. But to our knowledge, there is little development on the accelerated failure time model for recurrent events data with covariate measurement error. In this topic, we develop an estimation procedure for recurrent events data under the accelerated failure time model, with the availability of replicated mismeasured covariates.

The second and third topics concern the Cox proportional hazards model for univariate survival data. In the presence of covariate measurement error, several consistent approaches have been proposed for Cox model, namely conditional score (Tsiatis & Davidian, 2001), parametric corrected score (Nakamura, 1990), and non-parametric corrected score (Huang & Wang, 2000). But each of these methods suffers from finite-sample pathological behaviors when the measurement error is substantial. Nonetheless, neither the nature nor the severity of these pathological behaviors are well understood in the literature. In these two topics, we investigate the pathological behaviors of these methods and propose new estimation procedure to tackle the pathological behaviors. In the second topic, we consider the situation that the distribution of measurement error is known. In the third topic, the error distribution is

no longer known but replicated mismeasured covariates are available.

In this chapter, we first present two real data sets that motivate our dissertation research, the Nutritional Prevention of Cancer (NPC) trial and the AIDS Clinical Trial Group (ACTG) 175 study. After that, we will give a brief introduction of general measurement error problems and existing methods in dealing with them. Reviews of relevant existing literature specific to survival data will also be provided. Finally, we will present an outline of this dissertation.

## 1.2 Motivating Examples

The methods proposed in this dissertation are motivated by two problems. The first one, which motivates our research on the accelerated failure time model for recurrent events data, is the Nutritional Prevention of Cancer (NPC) trial (Clark et al., 1996). Our research on the Cox proportional hazards model is motivated by the AIDS Clinical Trial Group (ACTG) 175 study.

### 1.2.1 Nutritional Prevention of Cancer (NPC) Trial

The Nutritional Prevention of Cancer (NPC) trial is a randomized double-blind, placebo-controlled clinical trial to evaluate the long term safety and efficacy of daily 200-$\mu$g supplement of selenium (Se) in preventing two types of skin cancers: basal cell carcinoma (BCC) and squamous cell carcinoma (SCC). From 1983 through 1991, a total of 1,312 patients with previous histories of BCC or SCC were recruited and

randomized into the treatment group in which patients were supposed to take daily 200-$\mu$g supplement of selenium (Se), or the control group in which placebo pills were given. Patients were then followed for a period of over twelve years. Follow up visits to clinics were scheduled for every 6 months or more frequent if necessary to exam for new dermatological problems or selenium toxicity. The primary end points of this trial were new occurrences of BCC or SCC. Possible censoring events include death and end of routine dermatological exams.

Clark et al. (1996) reported an adverse but statistically non-significant effect of selenium supplement in prevention of new occurrence of BCC or SCC. Their findings are contrary to prior expectation since highly significant positive benefits were found on selenium supplement in preventing a number of other types of cancers. The findings of Clark et al. (1996) might be contributed to two reasons. First, they did not consider recurrent BCC or SCC. During the study, each patient might experience multiple occurrences of BCC or SCC. But in the analysis, they considered only time to the first occurrence of BCC or SCC. Second, the baseline plasma Se level as an important prognostic risk factor for BCC and SCC is subject to substantial measurement error. But the authors did not consider the possible effect of measurement error on statistical analysis.

One natural question arises is how to incorporate these additional information of time to multiple occurrences of BCC or SCC and consistently estimate the covariate effect for recurrent events data under the circumstance of covariate measurement error. We will develop methods in Chapter 2 to address this question.

### 1.2.2 AIDS Clinical Trial Group (ACTG) 175 Study

AIDS Clinical Trial Group (ACTG) 175 study is a randomized clinical trial to evaluate four treatments, zidovudine alone (ZDV), zidovudine plus didanosine (ZDV + ddI), zidovudine plus zalcitabine (ZDV + ddC), and zalcitabine alone (ddC), in HIV-infected subjects with a initial screening CD4 count between 200 and 500. As a note, a healthy individual without HIV infection usually has a CD4 count of between 800 and 1200. A total of 2,467 HIV-infected volunteers participated in the ACTG 175 study and, among them, 1,067 of whom had never taken any antiretroviral therapy at study entry; see Hammer et al. (1996) for a more detailed description of this study.

We are interested in assessing the effect of baseline CD4 count on time to AIDS or death in patient. The Cox proportional hazards model is the most popular regression model for right-censored survival data and it has been widely applied in a large number of HIV/AIDS studies. However, the issue of covariate measurement error arises often in practice. For example, the CD4 count has no gold standard measurement and is subject to substantial measurement error. In ACTG 175 study, 1,036 antiretroviral naive patients had two duplicated baseline CD4 count measurements prior to the start of treatment and within 3 weeks of randomization. The duplicated baseline CD4 count measurements were taken from different blood samples. Figure 1.1 illustrates the issue of measurement error on ACTG 175 study.

Although there exists a rich literature on Cox regression with covariate measurement error, all available consistent methods suffer from severe finite-sample pathological behaviors. For example, the estimating functions may have multiple zero-crossings, no zero-crossing, or a single zero-crossing that is far away from the true

Figure 1.1: Replicated baseline CD4 count measurements from 1,036 antiretroviral-naive patients in the ACTG 175 study.

parameter. Though the pathological behaviors have been noticed in the literature, they were never well understood or systematically studies.

The above observation prompts the second research question of how to improve the finite-sample behaviors of an estimating function for Cox regression with substantial covariate measurement error.

## 1.3   Covariate Measurement Error

In literature, there are two broad classes of measurement error models, namely, *classical measurement error model* and *Berkson measurement error model* (Berkson, 1950).

Suppose covariates $\mathbf{Z} = (\mathbf{Z}_a^T, \mathbf{Z}_e^T)^T$ where $\mathbf{Z}_a$ are those covariates that can be accurately measured and $\mathbf{Z}_e$ are covariates prone to measurement error and cannot be accurately measured. Though $\mathbf{Z}_e$ cannot be measured directly, we can observe them through their surrogates, $\mathbf{W}_e$. Under the classical additive measurement error model,

$$\mathbf{W}_e = \mathbf{Z}_e + \boldsymbol{\epsilon}_e$$

where $\boldsymbol{\epsilon}_e$ is a mean-zero random noise. In most applications, it is assumed that $\boldsymbol{\epsilon}_e$ is independent of both $\mathbf{Z}_a$ and $\mathbf{Z}_e$, and $\boldsymbol{\epsilon}_e \sim N(0, \sigma_e^2)$.

On the other hand, Berkson measurement error model is defined as

$$\mathbf{Z}_e = \mathbf{W}_e + \boldsymbol{\epsilon}_e$$

The true covariate equals mismeasured covariates plus measurement error. Same as in classical error model, $\epsilon_e$ are usually assumed to be independent of $\mathbf{Z}_a$ and $\mathbf{Z}_e$.

In this dissertation research, we will consider the classical measurement error model. Let $\mathbf{W} = (\mathbf{Z}_a, \mathbf{W}_e)$, it is well known that naively replacing $\mathbf{Z}$ by $\mathbf{W}$ in regression analysis could lead to substantial bias in estimation of some or all regression coefficients (Carroll et al., 2006; Fuller, 1987).

Here is a simple example to illustrate the effect of measurement error in naive analysis. Consider a simple linear regression with classical additive error, $Y = \beta_0 + \beta_z Z + \epsilon$, where $\text{Var}(Z) = \sigma_z^2$ and $\text{Var}(\epsilon) = \sigma_\epsilon^2$. Covariate $Z$ is unobservable; instead, a surrogate $W = Z + U$ is observed. $E(U|Z) = 0$, and $\text{Var}(U|Z) = \sigma_u^2$. If we fit a naive regression model of $Y$ on $W$, then the estimator is a consistent estimator of $\beta_{z*} = \lambda \beta_z$, where

$$\lambda = \frac{\sigma_z^2}{\sigma_z^2 + \sigma_u^2} < 1$$

and

$$\text{Var}(Y|W) = \sigma_\epsilon^2 + \frac{\beta_z^2 \sigma_u^2 \sigma_z^2}{\sigma_z^2 + \sigma_u^2} > \sigma_\epsilon^2$$

The naive estimation procedure gives an estimator that is attenuated to zero. As the example shows, measurement error not only produces a biased parameter estimate but also increase the residual variance. In a more complex model, the effect of measurement error in statistical inference will be much more complicated.

The simple measurement error models introduced above can be extended to accommodate more complex applications. For example, Kipnis et al. (1999) introduced a complex measurement error model for some nutritional studies which allows for

both bias and variance components.

To perform measurement error analysis, some additional information is needed. This additional information could be known error distribution, or in the case when error distribution is unknown, some forms of extra data. Depending on the source, these extra data can be categorized as *internal* data which are obtained for a subset of primary study and *external* data which are obtained from independent studies. In each category, there are three types of data: *validation data*, in which we can observe $(\mathbf{Z}, \mathbf{W})$ directly for a subset of data; *replication data*, in which replicates of $\mathbf{W}$ are available for each subject; *instrumental data*, in which an additional instrumental variable $\mathbf{T}$ is observed in additional to $\mathbf{W}$ (correlated with $\mathbf{Z}$, but may or may not be unbiased for $\mathbf{Z}$). In medical research fields, replication data are more commonly available and there exists some literature on measurement error problems with replication data.

## 1.4   Functional Modeling and Structural Modeling

A large number of statistical models have been proposed to tackle various measurement error problems. Traditionally, based on the property of the unobserved true covariates $\mathbf{Z}_e$, a distinction was made between *classical functional models*, in which the $\mathbf{Z}_e$ are regarded as unknown fixed constants or parameters, and *classical structural models*, in which the $\mathbf{Z}_e$ are regarded as random variables. Carroll et al. (2006) instead suggested to make a distinction between *functional modeling*, in which $\mathbf{Z}_e$ may be either fixed or random, but in the latter case, no or only minimal assumptions

are made about the distribution of $\mathbf{Z}_e$, and *structural modeling* in which parametric distributional models are imposed on $\mathbf{Z}_e$.

The structural modeling approach plays an important role in application (Carroll et al., 2006, chapter 8), but the robustness of inference to parametric model assumptions is of concern. The functional modeling approach is more appealing in the sense that the estimation procedures does not need specify a true covariate distribution and thus not subject to misspecification.

Available functional modeling methods might be categorized into three classes. The first includes the conditional score method (Stefanski & Carroll, 1987) and locally efficient score method (Tsiatis & Ma, 2004). Conditional score requires the measurement error to be normally distributed and the idea is to condition away the nuisance parameters based on certain sufficient statistics. Stefanski & Carroll (1987) obtained an unbiased score function for generalized linear measurement error models. Later, Tsiatis & Ma (2004) proposed a class of semiparametric estimators that require a pilot distribution of unobserved error-prone covariates to be specified. It is shown that such estimators are consistent no matter what the pilot distribution is and are efficient if computed under the truth. It was shown later by Ma & Tsiatis (2006) that when implemented in generalized liner model and with normal measurement error, the locally efficient score method (Tsiatis & Ma, 2004) is equivalent to conditional score. Therefore, their method could be viewed as an extension of conditional score.

The second class is the corrected score method (Nakamura, 1990; Stefanski, 1989; Huang & Wang, 2001); it is also referred to as parametric correction in the literature. The corrected score method requires the existence of an estimating function that

produces consistent estimators in the absence of measurement error. It is called the original or reference estimating function. A corrected score is an estimating function based on observed error-contaminated data that has the same limit as the reference estimating function. If the reference estimating function admits consistent estimates only, the corrected estimating function inherits this property in a compact parameter space containing the true value.

The third class is nonparametric correction (Huang & Wang, 1999, 2000, 2001, 2006). Similar to the parametric correction approach, one starts with a reference estimating function based on the underlying true covariates and then construct an estimating function with error-contaminated covariates which share the same limit as reference. The most important difference, however, is that the parametric correction method requires parametric distributional assumptions on the measurement error, whereas the nonparametric correction spares them with the availability of replicated mismeasured covariates.

To avoid any confusion, we will refer to the second class as *parametric corrected score* and the third class as *nonparametric corrected score* later in this dissertation.

# 1.5 Measurement Error Techniques in Survival Analysis

## 1.5.1 Recurrent Events Data

For recurrent events data without covariate measurement error, a number of regression models have been proposed in the literature as an extension of the Cox proportional hazards model. Andersen & Gill (1982) introduced the counting process model with a Cox-type intensity function for recurrent events and established an elegant large sample asymptotic theory through martingale theory. Wei et al. (1989) proposed a marginal approach for the analysis of recurrent failures. Later, Pepe & Cai (1993) and Lawless et al. (1997) proposed multiplicative models for the rate and mean functions of arbitrary counting processes. More recently, Lin et al. (1998) proposed to extend the univariate accelerated failure time model for the analysis of recurrent events. They developed a class of consistent and asymptotic normal rank estimators. Comparing to the Cox model, the accelerated failure time model has an appealing feature of a direct physical interpretation.

In the presence of covariate measurement error, Jiang et al. (1999) developed a method using replicated measurements of error-prone covariates and applied their method to the NPC trial data. Their method is based on a discrete-time proportional means model (see Lawless & Nadeau, 1995). Parametric assumptions are imposed on both the true covariates and the errors. Later, Hu & Lin (2004) considered an extended proportional hazards model. Though no parametric assumption is imposed on the covariates, their approach requires the errors to be symmetrically distributed.

To our knowledge, there is little development under the accelerated failure time model framework for recurrent events data with covariate measurement error.

## 1.5.2 Univariate Survival Data

There exist a large collection of literature dealing with the problem of measurement error in the analysis of Cox-type models for univariate survival data. Prentice (1982) introduced the induced hazard function and proposed an induced partial likelihood function for the Cox regression model with rare events and utilized the regression calibration approach to estimate the parameters. Clayton (1991) proposed a modification of regression calibration that does not require events to be rare. Another regression calibration method was later developed by Wang et al. (1997).

Though regression calibration is used frequently to yield a reasonable approximation, it is well known that the regression calibration method is inconsistent and may have large bias under certain circumstances. Zhou & Pepe (1995) developed a consistent estimator for regression coefficients. But the requirement of discrete covariates and the availability of a validation set greatly limit the applicability of their method in practice. Another approach is to correct the partial score or related estimating functions. Nakamura (1992) developed a parametric corrected score approach for normal measurement errors and it was shown by Kong & Gu (1999) that, the basic estimator of Nakamura (1992) is consistent and asymptotically normal. Tsiatis & Davidian (2001) proposed a conditional score for normal measurement error, which is asymptotically equivalent to Nakamura's parametric corrected score with normal error. For data with at least two replicates for mismeasured covariates, Huang &

Wang (2000) proposed a consistent nonparametric estimator based on a correction of the partial score function.

More recently, Song & Huang (2005) investigated the finite-sample performance of parametric corrected score and conditional score and found that they all might suffer from pathological behaviors as the error magnitude increases. They suggested a number of refinements, but the improvement is fairly modest.

## 1.6    Outline

Chapter 2 deals with parameter estimation in the accelerated failure time model for recurrent events when covariates are measured with error. We consider the classical additive measurement error model and propose a consistent estimation procedure for parameter of interest. With replicated mismeasured covariates available, the proposed estimation procedure requires no distributional assumptions on either the true covariates or the error except for the boundedness of the latter. The resulting estimators are proven to be consistent and asymptotically normal. Simulation studies indicate that the proposed model works well for practical sample sizes and moderate measurement error. An illustration with application to the NPC trial is also provided.

Chapter 3 considers the Cox proportional hazards model. In the presence of covariate measurement error, several functional modeling methods have been proposed for the proportional hazards model under the situation where the distribution of the measurement error is known. In this chapter, we study the finite-sample patholog-ical behaviors of parametric corrected score and conditional score. To address the

issue of pathological behaviors, we propose an approach to incorporate additional estimating functions. Extensive simulation studies are conducted to investigate the performance of the proposed method. Furthermore, we illustrate the practical utility of the proposed methods via an application to the ACTG 175 study.

In Chapter 4, we propose an extension of the method in Chapter 3 by allowing the error distribution to be completely unspecified. Huang & Wang (2000) developed a nonparametric corrected score for proportional hazards model when replicated measurements for mismeasured covariates are available. But the nonparametric corrected score suffers from similar finite-sample pathological behaviors as the parametric corrected score when the measurement error is substantial. We extend the method described in Chapter 3 to nonparametric corrected score and the simulation study result shows that the proposed estimation procedure is promising in resolving those finite-sample pathological behaviors.

In Chapter 5, we provide a summary and discuss some future work for the dissertation.

# Chapter 2

# Accelerated Failure Time Model for Recurrent Events With Errors in Covariates

## 2.1   Introduction

In clinical research studies, the event of interest may occur more than once for each subject during the follow up period. One example is the Nutritional Prevention of Cancer (NPC) trial (Clark et al., 1996; Duffield-Lillico et al., 2003), which involved multiple recurrences of skin cancers, basal cell carcinoma (BCC) and squamous cell carcinoma (SCC). For such recurrent events, a number of statistical methods have been developed to extend the univariate Cox proportional hazards model, including Andersen & Gill (1982), Wei et al. (1989), Pepe & Cai (1993) and Lawless et al.

(1997). More recently, Lin et al. (1998) proposed to extend the univariate accelerated failure time model for the analysis of recurrent events. As an appealing and well recognized feature, the accelerated failure time model has a direct physical interpretation. All these inference procedures presume the conventional setting that covariates are observed and measured accurately. However, this assumption may not be appropriate in some studies, including the aforementioned NPC trial. The baseline plasma selenium as an important prognostic risk factor for BCC and SCC is subject to substantial measurement error, due to natural biological fluctuation and instrument error.

For recurrent events data with covariate measurement error, available approaches include Jiang et al. (1999)'s method which is based on a discrete-time proportional means model utilizing replicated measurements of error-prone covariates, and an extended proportional hazards model developed by Hu & Lin (2004). But for the former approach, parametric assumptions are imposed on both the true covariates and the errors, whereas for the latter, though no parametric assumption is imposed on the covariates, it requires the errors to be symmetrically distributed. To our knowledge, there is little development on the accelerated failure time model for recurrent events data with covariate measurement error.

In this chapter, we propose a consistent estimation procedure for the accelerated failure time model based on replicated mismeasured covariates. The proposed approach bears similarity, in minimal assumption requirement, to the nonparametric corrected score method used in the other regression models (Huang & Wang, 2000, 2001, 2006). But this model of interest is not amendable to the correction strategy

and the proposed estimating function is developed on the basis of a novel identity. In Section 2.2, we present the proposed estimation procedure. Consistency and asymptotic normality of the proposed estimator will also be established. Simulation studies with practical sample size are reported in Section 2.3, along with an illustration with the NPC trial data. Further discussion is given in Section 2.4. Technical details and proofs for the large-sample study on the proposed inference procedure are collected in Section 2.5.

## 2.2    Inference Procedure

### 2.2.1    The Model and Data

Let $\tilde{N}^*(t)$ denotes the number of recurrent events by time $t$ in the absence of censoring and $\mathbf{Z}$ the p-vector covariates of interest. The accelerated failure time model (Lin et al., 1998) postulates their relationship as

$$E\{\tilde{N}^*(t)|\mathbf{Z}\} = \mu_0(e^{-\boldsymbol{\beta}_0'\mathbf{Z}}t), \tag{2.1}$$

where $\boldsymbol{\beta}_0$ is an unknown p-vector of parameter of interest, and $\mu_0(.)$ is an unspecified baseline rate function. Let $C$ be the censoring time, and the conditionally independent censoring mechanism is adopted, i.e., $C \perp \tilde{N}^*(\cdot) \,|\, \mathbf{Z}$, where $\perp$ denotes statistical independence. In the presence of censoring, $\tilde{N}^*(\cdot)$ is not fully observed but only through $\tilde{N}(\cdot) = \tilde{N}^*(\cdot \wedge C)$.

In the case of covariate measurement error, some elements of $\mathbf{Z}$ cannot be accu-

rately measured. Split covariates $\mathbf{Z} = (\mathbf{Z}_a^T, \mathbf{Z}_e^T)^T$, where $\mathbf{Z}_a$ are accurately measured and $\mathbf{Z}_e$ are error prone. Though $\mathbf{Z}_e$ cannot be measured directly, we can observe it through their surrogates, $\mathbf{W}_e$. Under the additive measurement error model, $R \geq 2$ surrogates $\{\mathbf{W}_e^{(m)} : m = 1, \ldots, R\}$ are observed, where

$$\mathbf{W}_e^{(m)} = \mathbf{Z}_e + \boldsymbol{\epsilon}_e^{(m)}, m = 1, \ldots, R. \tag{2.2}$$

These errors $\boldsymbol{\epsilon}_e^{(m)}$ are iid replicates of $\boldsymbol{\epsilon}_e$ and are independent of all other random variables. No distributional assumption is imposed on covariates $\mathbf{Z}$, nor on $\boldsymbol{\epsilon}_e$ except for boundedness. The observed data, $\{\tilde{N}_i(\cdot); \ C_i; \ \mathbf{Z}_{ai}; \ R_i; \ \mathbf{W}_{ei}^{(m)} : m = 1, \ldots, R_i\}, i = 1, \ldots, n$, consist of n iid replicates of $\{\tilde{N}(\cdot); \ C; \ \mathbf{Z}_a; \ R; \ \mathbf{W}_e^{(m)} : m = 1, \ldots, R\}$.

## 2.2.2 Proposed Estimation Procedure

Pick arbitrarily two replicates $\mathbf{W}_e^{(1)}$ and $\mathbf{W}_e^{(2)}$ from $\{\mathbf{W}_e^{(m)} : m = 1, \ldots, R\}$ and denote $\mathbf{W}^{(1)} = (\mathbf{Z}_a^T, \mathbf{W}_e^{(1)T})^T$ and $\mathbf{W}^{(2)} = (\mathbf{Z}_a^T, \mathbf{W}_e^{(2)T})^T$; $R(R-1)$ different permutations can be formed. We assume that $\boldsymbol{\epsilon}_e$ is bounded. Therefore, there exists a function of $\boldsymbol{\beta}$, say $\xi(\boldsymbol{\beta})$, such that

$$\xi(\boldsymbol{\beta}) \geq |\boldsymbol{\beta}'(\mathbf{W}^{(1)} - \mathbf{W}^{(2)})|, a.s. \tag{2.3}$$

Define counting process $N(t; \mathbf{Z}, \boldsymbol{\beta}) = \tilde{N}\{\exp(\boldsymbol{\beta}'\mathbf{Z})t\}$, and at-risk process $Y(t; \mathbf{Z}, \boldsymbol{\beta}) = I\{C \geq \exp(\boldsymbol{\beta}'\mathbf{Z})t\}$. To motivate our estimation procedure, we obtain an important

identity:

$$E\{N(t; \mathbf{W}^{(1)}, \boldsymbol{\beta}_0) | \mathbf{Z}, Y(e^{\xi(\boldsymbol{\beta}_0)}t; \mathbf{W}^{(2)}, \boldsymbol{\beta}_0) = 1, \mathbf{W}^{(2)}\}$$

$$= E\{\tilde{N}^*(e^{\boldsymbol{\beta}_0' \mathbf{W}^{(1)}} t \wedge C) | \mathbf{Z}, C \geq e^{\boldsymbol{\beta}_0' \mathbf{W}^{(2)} + \xi(\boldsymbol{\beta}_0)} t, \mathbf{W}^{(2)}\}$$

$$= E\{\tilde{N}^*(e^{\boldsymbol{\beta}_0' \mathbf{W}^{(1)}} t) | \mathbf{Z}, C \geq e^{\boldsymbol{\beta}_0' \mathbf{W}^{(2)} + \xi(\boldsymbol{\beta}_0)} t, \mathbf{W}^{(2)}\}$$

$$= E\{\tilde{N}^*(e^{\boldsymbol{\beta}_0' \mathbf{W}^{(1)}} t) | \mathbf{Z}\}$$

$$= E\{\mu_0(e^{\boldsymbol{\beta}_0' \boldsymbol{\epsilon}_e} t)\} \equiv \Lambda_0(t), \tag{2.4}$$

where we have used the fact $\boldsymbol{\beta}_0' \mathbf{W}^{(2)} + \xi(\boldsymbol{\beta}_0) \geq \boldsymbol{\beta}_0' \mathbf{W}^{(1)}$ and $\Lambda_0(t)$ is a quantity independent of $\mathbf{Z}$.

Identity (2.4) naturally leads to the following estimating function of both $\boldsymbol{\beta}_0$ and $\Lambda_0(t)$:

$$\frac{1}{n} \sum_{i=1}^{n} \int_0^{\tau} \mathcal{A} \begin{pmatrix} 1 \\ \mathbf{W}_i^{(2)} \end{pmatrix} Y_i(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}_i^{(2)}, \boldsymbol{\beta}) \mathrm{d}\{N_i(t; \mathbf{W}_i^{(1)}, \beta) - \Lambda(t)\} = 0,$$

where $\tau$ is a prespecified constant, and $\mathcal{A}$ denotes the operator averaging over all the different permutations of $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$. Further profiling out $\Lambda(t)$, we have the following estimating function for $\boldsymbol{\beta}_0$:

$$\Psi(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^{n} \int_0^{\tau} \mathcal{A} \left[ \mathbf{W}_i^{(2)} - \frac{\sum_{j=1}^{n} \mathcal{A}\{\mathbf{W}_j^{(2)} Y_j(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}_j^{(2)}, \boldsymbol{\beta})\}}{\sum_{j=1}^{n} \mathcal{A}\{Y_j(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}_j^{(2)}, \boldsymbol{\beta})\}} \right]$$
$$\times Y_i(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}_i^{(2)}, \boldsymbol{\beta}) \mathrm{d}N_i(t; \mathbf{W}_i^{(1)}, \boldsymbol{\beta}) \tag{2.5}$$

In the absence of covariate measurement error, $\mathbf{W}^{(1)} = \mathbf{W}^{(2)} = \mathbf{Z}$ and we may set $\xi(\boldsymbol{\beta}) = 0$. Then, the above estimating function reduces to Lin et al. (1998)'s

estimating function.

It is clear that $\xi(\boldsymbol{\beta})$ results in additional censoring in the proposed estimation procedure. Therefore, this function should be kept as small as possible. Suppose that $\mathbf{Z}_e$ is $q$-dimensional. Let $\boldsymbol{\epsilon}_{[k]}^{(1)}, \boldsymbol{\epsilon}_{[k]}^{(2)}, \boldsymbol{\beta}_{[k]}, k = 1, \ldots, q$, be the $k$th element of corresponding vectors. Let $M_{[k]} = \max |\boldsymbol{\epsilon}_{[k]}^{(1)} - \boldsymbol{\epsilon}_{[k]}^{(2)}|$, the function $\xi(\boldsymbol{\beta})$ can be set as $\xi(\boldsymbol{\beta}) = \sum_{k=1}^{q} |\boldsymbol{\beta}_{[k]}| M_{[k]}$. Of course, $M_{[k]}$ is typically unknown. For practical purpose, it may be empirically determined; see Section 2.3.

## 2.2.3 Asymptotic Properties

To allow for a rigorous large-sample study, we consider the following mild regularity conditions:

C1. Counting process $\tilde{N}(\cdot)$ and covariates $\mathbf{Z}$ are bounded;

C2. $C$ has a bounded density function and $\mu_0$ has a bounded second derivative;

C3. Support of $\boldsymbol{\epsilon}_e$ is bounded;

C4. Time limit $\tau$ satisfies $\lim_{n \to \infty} n^{-1} \sum_{i=1}^{n} Y_i(\tau; \mathbf{Z}_i, \boldsymbol{\beta}_0) > 0$.

**Theorem 2.1.** *Under conditions C1-C4, almost surely, a zero-crossing of $\Psi(\boldsymbol{\beta})$, say $\hat{\boldsymbol{\beta}}$, exists and converges to $\boldsymbol{\beta}_0$.*

The estimating function $\Psi(\boldsymbol{\beta})$ is a piecewise constant function of $\boldsymbol{\beta}$. Furthermore, the estimating function is not generally monotone. Therefore, the computation is not simple. When the dimension of covariates is small, $\hat{\boldsymbol{\beta}}$ can be obtained by using an

iterative bisection method (Huang, 2002) or direct grid search. When the dimension of covariates is high, specialized numerical methods such as simulated annealing (Lin & Geyer, 1992) might be more efficient.

**Theorem 2.2.** *Under conditions C1-C4, $n^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$ is asymptotically normal with mean 0.*

In terms of the interval estimation for $\boldsymbol{\beta}_0$, there are two general approaches available. One is to use bootstrap resampling which is computational intensive. An alternative method was proposed by Huang (2002) who developed a sample-based variance estimation procedure for the estimators based on nonsmooth estimation functions in general. Compared to the bootstrap resampling, Huang (2002)'s approach takes much less computing time. To adopt Huang (2002)'s approach, a consistent estimate of $\Psi(\boldsymbol{\beta}_0)$ is given in Section 2.5.

## 2.3   Numerical Studies

We investigated the performance of the proposed estimation procedure under practical sample size via extensive simulation studies. Also the procedure will be illustrated through an application to the NPC study.

For reference and comparison, the ideal, naive, and regression calibration estimators were also studied. The ideal estimator used the procedure of Lin et al. (1998) with the true covariates and, of course, it is not a realistic estimator. The naive approach uses the average of replicated surrogates in replace of true covariates in the procedure of Lin et al. (1998), whereas the regression calibration approach uses the

best linear approximation by replication data given in Carroll et al. (2006, chap. 4). For all these estimators, we used the approach of Huang (2002) for interval estimation.

For all these estimators, one needs to specify the time limit $\tau$ for their estimating function. We chose $\tau$ large enough to include all the follow-up time, as is standard practice in survival analysis.

Typically, the bound of $\boldsymbol{\epsilon}_e$ is unknown. For practical purpose, we suggest to use an empirical version of $M_{[k]}$ given in Section 2.2.2 in order to specify $\xi(\boldsymbol{\beta})$ function. Let $\mathbf{W}^{(1)}_{[k]}, \mathbf{W}^{(2)}_{[k]}, k = 1, \ldots, q$, be the $k$th element of corresponding vectors. An empirical version of $M_{[k]}$ is $\max_{i=1,\ldots,n} |\mathbf{W}^{(1)}_{i[k]} - \mathbf{W}^{(2)}_{i[k]}|$. The operator max here also applies to all possible pairs of $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$.

## 2.3.1 Simulations

Systematic numerical experiments were conducted to investigate the finite-sample performance of the proposed estimation procedure.

We first considered the accelerated failure time model with a single and error-prone covariate and set $\beta_0$=0.5. Recurrent events were generated from Poisson processes with unit baseline rate. To evaluate the performance of the proposed approach in general applications, we considered both symmetric and asymmetric distributions of the covariate and measurement error. For the true covariate, we considered i) a standard normal distribution and ii) a chi-square distribution location-shifted to mean 0 and scale-changed to variance 1. The surrogate W of Z was generated from the additive error model $W = Z + \epsilon$, where $\epsilon$ was i) normal with mean zero and variance

$\sigma_\epsilon^2 = 0.25$, or 0.5, or ii) modified chi-square distribution with mean 0 and variance 0.25 or 0.5. For ii), the modification was made on the chi-square distribution with 1 degree of freedom by truncation at 5 and then by location shift and scale change. Two replicates of $W$ were generated for each Z. Censoring times $C$ were generated from the uniform distribution on $[0, \lambda]$, where $\lambda$ was chosen to yield an average of 10 events per subject. With a sample size of 400, 1,000 data sets were generated.

Table 2.1 summarizes the simulation results. For each scenario, the mean bias, standard deviation, averaged standard error, and coverage probability of 95% Wald-type confidence interval were calculated. For the proposed estimator, median bias, a robustified standard deviation defined as inter quartile range (IQR) divided by 1.349, and median of estimated standard error were also calculated. In all different settings, naive estimator was biased toward 0 with the bias increasing as the magnitude of measurement error increasing. Naive estimator also had poor coverage probability under all scenarios. Regression calibration estimator performed well in both bias and coverage probability when the underlying true covariate was normal. However, when the true covariate was modified chi-square distribution, regression calibration estimator was biased toward 0 and the coverage probabilities were lower than nominal level. The proposed estimator performed well with small bias when the true covariate was normally distributed. However, when true covariate followed the modified chi-square distribution, the mean bias could be substantial especially with increasing error variance although median bias was still small in these circumstances. Nevertheless, the distribution of the proposed estimator was quite skewed, which is similar to other

Table 2.1: Simulation summary statistics of the estimators in the single-covariate model

| $\sigma^2_\epsilon$ | | I | NV | RC | Prop | | I | NV | RC | Prop | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | \multicolumn — $Z, \epsilon \sim$ Normal | | | | | $Z, \epsilon \sim$ Modified $\chi^2$ | | | | |
| 0.25 | B | 5 | -561 | 0 | 5 | -5 | -11 | -1110 | -613 | 394 | 20 |
| | SD | 179 | 187 | 218 | 337 | 311 | 277 | 272 | 297 | 2044 | 814 |
| | SE | 181 | 197 | 222 | 385 | 361 | 286 | 273 | 308 | 1755 | 1060 |
| | CP | 95.4 | 19.3 | 95.3 | 94.8 | | 94.5 | 1.7 | 46.1 | 93.6 | |
| 0.5 | B | 5 | -994 | 16 | 88 | 54 | -11 | -1724 | -885 | 1497 | 146 |
| | SD | 179 | 207 | 297 | 535 | 478 | 277 | 278 | 341 | 4623 | 1546 |
| | SE | 181 | 203 | 255 | 639 | 547 | 286 | 279 | 351 | 3086 | 2103 |
| | CP | 95.4 | 0.5 | 90.6 | 95.0 | | 94.5 | 0 | 28.3 | 90.4 | |
| | | $Z \sim$ Normal, $\epsilon \sim$ Modified$\chi^2$ | | | | | $Z \sim$ Modified$\chi^2, \epsilon \sim$ Normal | | | | |
| 0.25 | B | -8 | -304 | 11 | -2 | -29 | 11 | -1198 | -706 | 84 | 29 |
| | SD | 187 | 196 | 209 | 291 | 290 | 282 | 248 | 285 | 619 | 608 |
| | SE | 180 | 188 | 201 | 306 | 293 | 287 | 250 | 283 | 864 | 703 |
| | CP | 95.0 | 63.3 | 93.4 | 95.7 | | 95.7 | 0.1 | 30.1 | 96.2 | |
| 0.5 | B | -8 | -561 | 37 | 17 | -29 | 11 | -1859 | -1044 | 446 | 29 |
| | SD | 187 | 200 | 232 | 366 | 362 | 282 | 236 | 314 | 2404 | 898 |
| | SE | 180 | 195 | 221 | 405 | 391 | 287 | 242 | 303 | 1944 | 1270 |
| | CP | 95.0 | 19.2 | 93.3 | 95.1 | | 95.7 | 0 | 8.9 | 94.3 | |

Note: I, ideal method; NV, naive estimator; RC, regression calibration; Prop, the proposed estimator. Except for the second column under Prop, B: mean bias ($\times 10^4$); SD: standard deviation($\times 10^4$); SE: averaged standard error ($\times 10^4$); CP, coverage probability(%) of the 95% Wald confidence interval. For the second column under Prop, B: meadian bias ($\times 10^4$); SD: robustified standard deviation($\times 10^4$); SE: median standard error ($\times 10^4$).

correction-type estimators (e.g. Nakamura, 1990). This skewness is also reflected in standard deviation versus robustified standard deviation and averaged standard error versus median standard error. The coverage probabilities of the proposed estimator were close to nominal level under all settings.

We also considered an accelerated failure time model with two covariates. The first covariate is subject to measurement error, with two surrogate replicates available, and the second covariate is measured precisely. The true coefficients $\boldsymbol{\beta} = (\beta_1, \beta_2)^T = (0.5, 0.5)^T$. We present the results from two scenarios. In the first scenario, the true covariates followed the standard bivariate normal distribution with correlation coefficient $\rho$. The measurement error follows a normal distribution with mean 0 and variance 0.25 or 0.5. In the second scenario, the true covariates were generated from the standard bivariate normal distribution upon a marginal increasing transformation and that the first covariate has the $\chi^2(1)$ distribution location-shifted to 0 and scale-changed to 1. The measurement errors were generated from $\chi^2(1)$ distribution truncated at 5 and were transformed to have mean 0 and variance $= 0.25$ or 0.5 by scale change and location shift. For both scenarios, we considered different correlation parameter $\rho$ between two true covariates. Recurrent events were generated from Poisson processes. Censoring times were generated from the uniform distribution on $[0, \lambda]$, where $\lambda$ was chosen to yield an average of 5 events per subject. With a sample size of 400, 1,000 data sets were generated.

Table 2.2 and 2.3 summarize the simulation results for multivariate cases. As

Table 2.2: Simulation summary statistics of the estimators in the multi-covariate model with normal error

| $\rho$ | $\sigma_\epsilon^2$ | | estimators of $\beta_1$ | | | | | estimators of $\beta_2$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | I | NV | RC | Prop | | I | NV | RC | Prop | |
| 0.5 | 0.25 | B | -23 | -741 | -24 | 36 | 0 | 10 | 369 | 11 | 5 | 10 |
| | | SD | 285 | 304 | 365 | 555 | 565 | 289 | 316 | 325 | 488 | 507 |
| | | SE | 302 | 303 | 361 | 696 | 625 | 277 | 315 | 315 | 519 | 502 |
| | | CP | 95.7 | 31.5 | 94.2 | 95.3 | | 93.2 | 78.5 | 94.1 | 95.9 | |
| | 0.5 | B | -23 | -1277 | -19 | 44 | -68 | 10 | 637 | 10 | 18 | 29 |
| | | SD | 285 | 304 | 443 | 893 | 695 | 289 | 337 | 362 | 623 | 623 |
| | | SD | 302 | 299 | 414 | 1169 | 898 | 277 | 339 | 341 | 695 | 625 |
| | | CP | 95.7 | 1.9 | 92.7 | 94.4 | | 93.2 | 53.6 | 91.8 | 96.3 | |
| 0 | 0.25 | B | 4 | -556 | 4 | 46 | 29 | 1 | 8 | 7 | 16 | 10 |
| | | SD | 260 | 265 | 307 | 521 | 478 | 269 | 293 | 295 | 441 | 434 |
| | | SE | 270 | 279 | 315 | 630 | 549 | 254 | 279 | 276 | 460 | 430 |
| | | CP | 95.7 | 48.9 | 94.6 | 94.8 | | 93.7 | 93.4 | 92.7 | 94.1 | |
| | 0.5 | B | 4 | -1001 | 10 | 45 | -29 | 1 | 11 | 8 | 9 | -10 |
| | | SD | 260 | 266 | 358 | 712 | 666 | 269 | 311 | 317 | 535 | 513 |
| | | SE | 270 | 281 | 354 | 1007 | 792 | 254 | 297 | 293 | 586 | 533 |
| | | CP | 95.7 | 5.4 | 94.6 | 94.9 | | 93.7 | 93.9 | 92.7 | 95.0 | |
| -0.5 | 0.25 | B | -3 | -718 | 6 | 56 | -10 | 2 | -360 | 3 | 49 | 10 |
| | | SD | 307 | 301 | 365 | 599 | 579 | 292 | 304 | 328 | 531 | 507 |
| | | SE | 378 | 366 | 456 | 943 | 822 | 281 | 296 | 312 | 641 | 562 |
| | | CP | 98.1 | 49.3 | 98.6 | 98.3 | | 94.3 | 74.6 | 93.2 | 94.9 | |
| | 0.5 | B | -3 | -1253 | 14 | 261 | 44 | 2 | -635 | -1 | 195 | 29 |
| | | SD | 307 | 291 | 432 | 1167 | 726 | 292 | 313 | 363 | 857 | 644 |
| | | SE | 378 | 351 | 529 | 1413 | 1043 | 281 | 305 | 339 | 1003 | 779 |
| | | CP | 98.1 | 4.6 | 98.1 | 96.8 | | 94.3 | 45.2 | 93.1 | 96.2 | |

Note: Same as that of Table 2.1.

Table 2.3: Simulation summary statistics of the estimators in the multi-covariate model with $\chi^2$ error

| $\rho$ | $\sigma_\epsilon^2$ | | | estimators of $\beta_1$ | | | | | estimators of $\beta_2$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | I | NV | RC | Prop | | I | NV | RC | Prop | |
| 0.5 | 0.25 | B | -12 | -1212 | -862 | 144 | 10 | 11 | 322 | 151 | -10 | -29 |
| | | SD | 529 | 450 | 491 | 1073 | 956 | 301 | 314 | 323 | 452 | 449 |
| | | SE | 613 | 471 | 532 | 1544 | 1191 | 255 | 287 | 279 | 509 | 470 |
| | | CP | 97.4 | 26.5 | 63.4 | 96.0 | | 90.5 | 78.7 | 57.6 | 95.7 | |
| | 0.5 | B | -12 | -1872 | -1294 | 328 | 0 | 11 | 501 | 217 | -54 | 5 |
| | | SD | 529 | 418 | 497 | 2017 | 1238 | 301 | 318 | 337 | 644 | 548 |
| | | SD | 613 | 427 | 529 | 2187 | 1794 | 255 | 299 | 293 | 714 | 636 |
| | | CP | 97.4 | 0.9 | 30.2 | 95.3 | | 90.5 | 62.1 | 85.3 | 95.6 | |
| 0 | 0.25 | B | 9 | -720 | -419 | 68 | 10 | 6 | 17 | 17 | 19 | 29 |
| | | SD | 400 | 354 | 377 | 633 | 652 | 264 | 274 | 274 | 388 | 362 |
| | | SE | 412 | 368 | 395 | 826 | 724 | 257 | 266 | 265 | 401 | 394 |
| | | CP | 95.2 | 48.5 | 80.2 | 96.3 | | 93.8 | 94.3 | 94.5 | 94.9 | |
| | 0.5 | B | 9 | -1207 | -675 | 140 | 29 | 6 | 20 | 20 | 18 | 10 |
| | | SD | 400 | 339 | 384 | 866 | 825 | 264 | 283 | 284 | 468 | 463 |
| | | SE | 412 | 353 | 404 | 1707 | 1074 | 257 | 272 | 271 | 571 | 479 |
| | | CP | 95.2 | 8.0 | 58.8 | 96.4 | | 93.8 | 94.4 | 93.8 | 95.4 | |
| -0.5 | 0.25 | B | 29 | -646 | -244 | 58 | 10 | 9 | -220 | -22 | 10 | -10 |
| | | SD | 388 | 341 | 376 | 606 | 583 | 296 | 298 | 307 | 442 | 449 |
| | | SE | 475 | 422 | 481 | 975 | 858 | 264 | 273 | 276 | 484 | 430 |
| | | CP | 97.8 | 66.2 | 96.5 | 98.1 | | 91.7 | 84.0 | 91.4 | 94.0 | |
| | 0.5 | B | 29 | -1115 | -396 | 140 | -39 | 9 | -343 | -39 | 38 | -10 |
| | | SD | 388 | 389 | 461 | 1054 | 820 | 296 | 461 | 496 | 613 | 557 |
| | | SE | 475 | 401 | 512 | 1605 | 1232 | 264 | 283 | 292 | 770 | 620 |
| | | CP | 97.8 | 17.6 | 91.8 | 98.1 | | 91.7 | 69.5 | 91.6 | 92.9 | |

Note: Same as that of Table 2.1.

shown, when two covariates are correlated, the measurement error generally has impact not only on the error-prone covariate but also on that of the accurately measured one. The relative performance is similar to what was observed in the single covariate case.

### 2.3.2 Illustration with the NPC Trial Data

The NPC trial is a randomized double-blind, placebo-controlled clinical trial to evaluate the long term safety and efficacy of a daily 200-$\mu$g supplement of selenium (Se) in preventing two types of skin cancers, BCC and SCC. Patients with previous histories of BCC or SCC were recruited and randomized into the treatment group in which patients were supposed to take daily 200-$\mu$g supplement of selenium (Se), or the control group in which placebo pills were given. One of the primary end points of this trial was newly diagnosed SCC lesions.

In this analysis, we are interested in the effect of baseline plasma Se on new SCC lesions. In the NPC trial, measurements of plasma Se were taken from routine clinic visits at approximately six months intervals. For those in the control group, repeated measurements for each patient might represent replicated measurements of baseline plasma Se under the assumption of stationarity (Jiang et al., 1999; Hu & Lin, 2004). This, however, would not be true for patients in the Se group since the plasma Se level was likely to change due to the treatment. As an illustration of the proposed estimation procedure, we will restrict our attention to the patients in the control

group.

We considered selenium measurements within 24 months of randomization and included 589 patients in the control group with at least 2 such measurements in the analysis. The average number of replicates was 4.4 with a maximum of 12 replicates. One hundred sixty seven (28.4%) patients had at least one occurrence of SCC after randomization. More specifically, 83 (14.1%), 32 (5.4%), 19 (3.2%) and 33 (5.6%) patients had 1, 2, 3 and 4 or more occurrences. The average duration of follow-up was 67.7 months and the average number of occurrences was 0.67 per patient.

Table 2.4: NPC Trial SCC Data

| Method | log(Se) | |
| --- | --- | --- |
| | Estimate | SE |
| Naive I | 1.156 | 0.535 |
| Naive II | 1.280 | 0.634 |
| RC | 1.693 | 0.976 |
| Proposed | 1.743 | 1.292 |

Note: Baseline Se measurement (Naive I), Average of all Se measurements (Naive II), Regression Calibration (RC), and Proposed Estimator (Proposed)

The covariate of interest is the logarithm of the true baseline selenium. Table 2.4 shows the analysis results based on the proposed method, regression calibration method and two naive methods, one using only the baseline selenium measurement and the other using the average of all selenium measurements within 24 months of randomization. In comparison, the naive approaches yield coefficient estimates

of log(Se) with smaller magnitude. Both the regression calibration and proposed estimates are larger in magnitude. Given the fact that the distribution of error is symmetric and close to normal (Hu & Lin, 2004), this result was expected, and is consistent with our simulations.

## 2.4   Discussion

In this chapter, we have proposed a consistent estimation procedure for the accelerated failure time model with recurrent events data in the presence of covariate measurement error. Under the accelerated failure time model and additive measurement error model, this method yields consistent regression coefficient estimation with only mild boundedness assumption on the measurement error. Simulation studies showed satisfactory performance of the proposed estimation procedure in samples of moderate size and measurement error, even when the measurement error was not bounded.

Lin et al. (1998) considered a class of weighted estimating functions, including the log-rank estimating function and Gehan estimating function as special cases. Our proposed estimation procedure should be able to be extended along the same line. With a proper chosen weight function, the efficiency of estimation procedure might be improved.

The proposed estimation procedure falls into the category of functional modeling

since it does not impose any parametric distribution on covariates. In particular, the proposed estimation procedure resembles the nonparametric correction method in that both approaches are free of distributional assumptions on underlying covariates and measurement error. However, the proposed method is distinct in its motivation, as being based on a novel identity (2.4). In contrast, the nonparametric correction method relies on an original estimating function (in the absence of measurement error) to construct a corrected estimating function with the same limit.

This method requires replicated measurements on error-prone covariates, which might not always be available. For instance, patients in Se group of the NPC trial did not have replicated baseline plasma Se measurements. This limitation warrants future research.

## 2.5    Appendix: Technical Details

### 2.5.1    Proof of the Asymptotic Theory

Proof of Theorem 2.1

We first prove the uniform consistency of the estimating function $\Psi(\boldsymbol{\beta})$ given by

(2.5). By algebra, we can write $\Psi(\boldsymbol{\beta})$ as a function of four empirical processes

$$
\begin{aligned}
\Psi(\boldsymbol{\beta}) \;=\; & \widehat{\mathcal{E}}\mathcal{A}\int_0^\tau \mathbf{W}^{(2)}Y(e^{\xi(\boldsymbol{\beta})}s;\mathbf{W}^{(2)},\boldsymbol{\beta})\mathrm{d}N(s;\mathbf{W}^{(1)},\boldsymbol{\beta}) \\
& -\int_0^\tau \frac{\widehat{\mathcal{E}}\mathcal{A}\{\mathbf{W}^{(2)}Y(e^{\xi(\boldsymbol{\beta})}t;\mathbf{W}^{(2)},\boldsymbol{\beta})\}}{\widehat{\mathcal{E}}\mathcal{A}\{Y(e^{\xi(\boldsymbol{\beta})}t;\mathbf{W}^{(2)},\boldsymbol{\beta})\}} \\
& \times\mathrm{d}\widehat{\mathcal{E}}\mathcal{A}\int_0^t Y(e^{\xi(\boldsymbol{\beta})}s;\mathbf{W}^{(2)},\boldsymbol{\beta})\mathrm{d}N(s;\mathbf{W}^{(1)},\boldsymbol{\beta}) \\
\;\equiv\; & \mathbb{A}_1^{(n)}(\boldsymbol{\beta}) - \int_0^\tau \frac{\mathbb{A}_2^{(n)}(t,\boldsymbol{\beta})}{\mathbb{A}_3^{(n)}(t,\boldsymbol{\beta})}\mathrm{d}\mathbb{A}_4^{(n)}(t,\boldsymbol{\beta})
\end{aligned}
$$

where $\widehat{\mathcal{E}}$ represents sample empirical mean, e.g., $\widehat{\mathcal{E}}\mathcal{A}\{\mathbf{W}^{(2)}Y(e^{\xi(\boldsymbol{\beta})}t;\mathbf{W}^{(2)},\boldsymbol{\beta})\} = n^{-1}\sum_{i=1}^n \mathcal{A}\{\mathbf{W}_i^{(2)}\,Y_i(e^{\xi(\boldsymbol{\beta})}t;\mathbf{W}_i^{(2)},\boldsymbol{\beta})\}$.

Using the technique of Kosorok (2008), we will show in the following that all four empirical processes $\mathbb{A}_1^{(n)}$ to $\mathbb{A}_4^{(n)}$ are P-Donsker classes and thus each of them will converge to their limits uniformly in $t$ and $\boldsymbol{\beta}$.

We first show $\mathbb{A}_3^{(n)}$ is P-Donsker class. First consider the case when $\boldsymbol{\beta}$ is a scalar, we have $\xi(\boldsymbol{\beta}) = |\boldsymbol{\beta}|M_{[1]}$ where $M_{[1]}$ is defined in Section 2.2.2.

Empirical process

$$
\begin{aligned}
& Y(e^{\xi(\boldsymbol{\beta})}t;\mathbf{W}^{(2)},\boldsymbol{\beta}) \\
=\; & I(C \geq e^{\boldsymbol{\beta}\mathbf{W}^{(2)}+\xi(\boldsymbol{\beta})}t) \\
=\; & I\{\log C \geq \boldsymbol{\beta}\mathbf{W}^{(2)} + \xi(\boldsymbol{\beta}) + \log t\} \\
=\; & I\{\log C \geq \boldsymbol{\beta}(\mathbf{W}^{(2)} + M_{[1]}) + \log t\} \times I\{\log C \geq \boldsymbol{\beta}(\mathbf{W}^{(2)} - M_{[1]}) + \log t\}
\end{aligned}
$$

$I(\log C \geq \boldsymbol{\beta}(\mathbf{W}^{(2)} + M_{[1]}) + \log t)$ and $I(\log C \geq \boldsymbol{\beta}(\mathbf{W}^{(2)} - M_{[1]}) + \log t)$ are indicator functions of half-space in $\mathbb{R}^2$ and thus belongs to Vapnik-Červonenkis-class with index 4 (Kosorok, 2008, Lemma 9.12). Following the proof of Kosorok (2008, Lemma 8.12), we can show that these functions are pointwise measurable classes. Applying Theorem 9.3, Theorem 8.19 of Kosorok (2008) and the fact that indicator functions are bounded, we can see that $I(\log C \geq \boldsymbol{\beta}(\mathbf{W}^{(2)} + M_{[1]}) + \log t)$ and $I(\log C \geq \boldsymbol{\beta}(\mathbf{W}^{(2)} - M_{[1]}) + \log t)$ are P-Donsker classes indexed by $\beta$ and $t$. We now have that $Y(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}^{(2)}, \boldsymbol{\beta})$ is Donsker since products of bounded Donsker classes are also Donsker. When $\boldsymbol{\beta}$ is a $p$-element vector, $Y(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}^{(2)}, \boldsymbol{\beta})$ can be written as the product of $2^p$ indicator functions. Each of these indicator functions can be shown as Donsker class. Therefore we can show $Y(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}^{(2)}, \boldsymbol{\beta})$ is Donsker with similar argument. Thus $\mathbb{A}_3^{(n)}(t, \boldsymbol{\beta}) = \widehat{\mathcal{E}}\mathcal{A}(Y(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}^{(2)}, \boldsymbol{\beta}))$ is Donsker since sums of bounded Donsker classes are Donsker.

From Condition C3, $\mathbf{W}^{(2)}$ is P-Donsker. $\mathbf{W}^{(2)}Y(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}^{(2)}, \boldsymbol{\beta})$ is the product of two bounded P-Donsker classes and therefore it is also P-Donsker. Donsker property of $\mathbb{A}_2^{(n)}(t, \boldsymbol{\beta})$ is shown.

Now turn to $\mathbb{A}_4^{(n)}(t, \boldsymbol{\beta}) = \widehat{\mathcal{E}}\mathcal{A} \int_0^t Y(e^{\xi(\boldsymbol{\beta})}s; \mathbf{W}^{(2)}, \boldsymbol{\beta}) \mathrm{d}N(s; \mathbf{W}^{(1)}, \boldsymbol{\beta})$. Let $T_j$, $j = 1, 2, \ldots$, be the $j$th event time for the subject, and $B$ be the upper bound of counting

process $\tilde{N}(\cdot)$ (Condition C1),

$$
\begin{aligned}
\int_0^t Y(e^{\xi(\boldsymbol{\beta})}s; \mathbf{W}^{(2)}, \boldsymbol{\beta}) \mathrm{d}N(s; \mathbf{W}^{(1)}, \boldsymbol{\beta}) &= N^*(t \wedge Ce^{-\boldsymbol{\beta}'\mathbf{W}^{(2)} - \xi(\beta)}; \mathbf{W}^{(1)}, \boldsymbol{\beta}) \\
&= \tilde{N}^*(e^{\boldsymbol{\beta}'\mathbf{W}^{(1)}}t \wedge Ce^{-\boldsymbol{\beta}'\mathbf{W}^{(2)} - \xi(\boldsymbol{\beta}) + \boldsymbol{\beta}'\mathbf{W}^{(1)}}) \\
&= \sum_{k=1}^B I(T_k \le e^{\boldsymbol{\beta}'\mathbf{W}^{(1)}}t \wedge Ce^{-\boldsymbol{\beta}'\mathbf{W}^{(2)} - \xi(\boldsymbol{\beta}) + \boldsymbol{\beta}'\mathbf{W}^{(1)}})
\end{aligned}
$$

is the summation of a finite number of P-Donsker classes. Therefore $\int_0^t Y(e^{\xi(\boldsymbol{\beta})}s; \mathbf{W}^{(2)}, \boldsymbol{\beta}) \mathrm{d}N(s; \mathbf{W}^{(1)}, \boldsymbol{\beta})$ is P-Donsker. $\mathbb{A}_4^{(n)}(t, \boldsymbol{\beta})$ and $\mathbb{A}_1^{(n)}(\boldsymbol{\beta})$ are also P-Donsker.

With the uniform convergence of $\mathbb{A}_1^{(n)}$, $\mathbb{A}_2^{(n)}$, $\mathbb{A}_3^{(n)}$ and $\mathbb{A}_4^{(n)}$ and the fact that the limit of $\mathbb{A}_3^{(n)}$ is bounded away from 0, the uniform convergence of estimating function $\Psi(\boldsymbol{\beta})$ can be established. $\Psi(\boldsymbol{\beta})$ will converge to its limit $\psi(\boldsymbol{\beta})$ uniformly in $\boldsymbol{\beta}$. Note that the limit function $\psi(\boldsymbol{\beta})$ has a zero-crossing at $\boldsymbol{\beta}_0$. Thus the uniform convergence of the estimating function $\Psi(\boldsymbol{\beta})$ implies that there exists one zero-crossing of $\Psi(\boldsymbol{\beta})$ converging to $\boldsymbol{\beta}_0$.

## Proof of Theorem 2.2

First, we show the asymptotic normality of the estimating function $\Psi(\boldsymbol{\beta}_0)$.

Estimating function $\Psi(\boldsymbol{\beta})$ is a mapping from four empirical processes which have been shown to be P-Donsker. It can be shown by Lemma 12.2, Lemma 12.3 and Lemma 6.19 (chain rule) of Kosorok (2008) that the mapping is Hadamard differentiable. Applying functional delta method, we obtain that $n^{1/2}(\Psi(\boldsymbol{\beta}_0) - \psi(\boldsymbol{\beta}_0))$

converges weakly to a zero mean Gaussian process.

We then establish the asymptotic linearity of $\Psi(\boldsymbol{\beta})$ in a neighborhood of $\boldsymbol{\beta}_0$. Let $A_1$ to $A_4$ denote the limit of $\mathbb{A}_1^{(n)}$ to $\mathbb{A}_4^{(n)}$ respectively. We will first show the following property

$$\lim_{\eta\to0}\limsup_{n\to\infty}\Pr\left\{\sup_{||\boldsymbol{\beta}-\boldsymbol{\beta}_0||\le\eta}n^{1/2}||\mathbb{A}_k^{(n)}(\boldsymbol{\beta},t)-A_k(\boldsymbol{\beta},t)-\mathbb{A}_k^{(n)}(\boldsymbol{\beta}_0,t)+A_k(\boldsymbol{\beta}_0,t)||>\epsilon\right\}=0$$

$$(2.6)$$

for $k=1,\cdots,4$.

First consider $k=3$. Let $\mathbf{M}(\boldsymbol{\beta})=(sgn(\beta_{[1]})M_{[1]},\ldots,sgn(\beta_{[p]})M_{[p]})$. $\xi(\boldsymbol{\beta})$ can be re-written as $\boldsymbol{\beta}'\mathbf{M}(\boldsymbol{\beta})$. Therefore $\mathcal{A}\{Y(e^{\xi(\beta)}t;\mathbf{W}^{(2)},\beta)\}=\mathcal{A}\{Y(e^{\boldsymbol{\beta}'\mathbf{M}(\boldsymbol{\beta})}t;\mathbf{W}^{(2)},\boldsymbol{\beta})\}$. We have

$$\mathrm{Var}\left[\mathcal{A}\{Y(e^{\boldsymbol{\beta}'\mathbf{M}(\boldsymbol{\beta})}t;\mathbf{W}^{(2)},\boldsymbol{\beta})-Y(e^{\boldsymbol{\beta}_0'\mathbf{M}(\boldsymbol{\beta}_0)}t;\mathbf{W}^{(2)},\boldsymbol{\beta}_0)\}\right]$$

$$=\ \mathcal{A}\left(\mathrm{Var}\left[I\{\log C\ge\boldsymbol{\beta}'\mathbf{U}(\boldsymbol{\beta})+\log t\}-I\{\log C\ge\boldsymbol{\beta}_0'\mathbf{U}(\boldsymbol{\beta}_0)+\log t\}\right]\right)\ (\text{let }\mathbf{U}(\boldsymbol{\beta})=\mathbf{W}^{(2)}+\mathbf{M}(\boldsymbol{\beta}))$$

$$=\ \mathcal{A}(\mathrm{Var}\left[I\{\boldsymbol{\beta}'\mathbf{U}(\boldsymbol{\beta})+\log t\le\log C\le\boldsymbol{\beta}_0'\mathbf{U}(\boldsymbol{\beta}_0)+\log t\}\right.$$

$$\left.+\ I\{\boldsymbol{\beta}_0'\mathbf{U}(\boldsymbol{\beta}_0)+\log t\le\log C\le\boldsymbol{\beta}'\mathbf{U}(\boldsymbol{\beta})+\log t\}\right])$$

$$=\ \mathcal{A}(\mathrm{Var}\left[I\{\boldsymbol{\beta}'\mathbf{U}(\boldsymbol{\beta})+\log t\le\log C\le\boldsymbol{\beta}_0'\mathbf{U}(\boldsymbol{\beta}_0)+\log t\}\right]$$

$$+\mathrm{Var}\left[I\{\boldsymbol{\beta}_0'\mathbf{U}(\boldsymbol{\beta}_0)+\log t\le\log C\le\boldsymbol{\beta}'\mathbf{U}(\boldsymbol{\beta})+\log t\}\right])$$

$$\le\ \mathcal{A}(\mathrm{E}\left[I\{\boldsymbol{\beta}'\mathbf{U}(\boldsymbol{\beta})+\log t\le\log C\le\boldsymbol{\beta}_0'\mathbf{U}(\boldsymbol{\beta}_0)+\log t\}\right]$$

$$+\mathrm{E}\left[I\{\boldsymbol{\beta}_0'\mathbf{U}(\boldsymbol{\beta}_0)+\log t\le\log C\le\boldsymbol{\beta}'\mathbf{U}(\boldsymbol{\beta})+\log t\}\right])$$

$$=\ \mathcal{A}(P_1+P_2)$$

$\boldsymbol{\beta}'\mathbf{U}(\boldsymbol{\beta})$ is continuous in $\boldsymbol{\beta}$ at $\boldsymbol{\beta}_0$. When $||\boldsymbol{\beta} - \boldsymbol{\beta}_0|| \leq \delta$ and $\delta \to 0$, we have $P_1$ and $P_2$ $\to 0$. By Condition C2, then

$$\text{Var}\left[\mathcal{A}\{Y(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}^{(2)}, \boldsymbol{\beta}) - Y(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}^{(2)}, \boldsymbol{\beta}_0)\}\right] \to 0$$

The Donsker property of $\mathcal{A}\{Y(e^{\xi(\beta)}t; \mathbf{W}^{(2)}, \boldsymbol{\beta})\}$ implies the tightness of the process. Thus $\text{Var}\left[\mathcal{A}\{Y(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}^{(2)}, \boldsymbol{\beta}) - Y(e^{\xi(\boldsymbol{\beta})}t; \mathbf{W}^{(2)}, \boldsymbol{\beta}_0)\}\right] \to 0$ implies that equation (2.6) holds (Kosorok, 2008, p.15).

Now we have shown equation (2.6) holds for $k = 3$. Similarly, we can also show that equation (2.6) holds for $k = 1, 2$ and $4$.

Then with some simple algebra, one can show that

$$\lim_{\xi \to 0} \limsup_{n \to \infty} \Pr\left\{\sup_{||\boldsymbol{\beta} - \boldsymbol{\beta}_0|| \leq \xi} n^{1/2}||\Psi(\boldsymbol{\beta}, t) - \psi(\boldsymbol{\beta}, t) - \Psi(\boldsymbol{\beta}_0, t) + \psi(\boldsymbol{\beta}_0, t)|| > \epsilon\right\} = 0$$

$$(2.7)$$

After proving (2.7), together with the asymptotic normality of $\Psi(\boldsymbol{\beta})$, the asymptotic normality of $\hat{\boldsymbol{\beta}}$ follows via fairly standard arguments.

## 2.5.2 Asymptotic Variance of Estimating Function

For abbreviation of notation, we will use $Y^{(1)}$ to denote $Y(e^{\xi(\boldsymbol{\beta})}s; \mathbf{W}^{(1)}, \boldsymbol{\beta})$ and $N^{(1)}$ to denote $N(t; \mathbf{W}^{(1)}, \boldsymbol{\beta})$. Similar abbreviations will also be used for other notations.

With functional delta method, straightforward algebra gives

$$n^{1/2}\Psi(\boldsymbol{\beta}) \;=\; n^{-1/2}\sum_{i=1}^{n}\left[B_{i1}-B_{i2}\right]+o_p(1) \tag{2.8}$$

where

$$
\begin{aligned}
B_{i1} &= \mathcal{A}\int_0^\tau \left\{\mathbf{W}_i^{(2)}-\frac{\mathcal{E}\mathbf{W}^{(2)}Y^{(2)}}{\mathcal{E}Y^{(2)}}\right\}Y_i^{(2)}\mathrm{d}N_i^{(1)} \\
B_{i2} &= \mathcal{A}\int_0^\tau \left\{\frac{\mathbf{W}_i^{(2)}Y_i^{(2)}}{\mathcal{E}Y^{(2)}}-\frac{\left(\mathcal{E}\mathbf{W}^{(2)}Y^{(2)}\right)Y_i^{(2)}}{\left(\mathcal{E}Y^{(2)}\right)^2}\right\}\mathrm{d}\mathcal{E}\int Y^{(2)}\mathrm{d}N^{(1)}
\end{aligned}
$$

Thus, $n^{1/2}\Psi(\boldsymbol{\beta})$ is asymptotically a sum of iid random variables.

For fixed $\boldsymbol{\beta}$, $n^{1/2}(\Psi(\boldsymbol{\beta})-\psi(\boldsymbol{\beta}))$ is asymptotically normal with mean 0 and a co-variance matrix $\Omega(\boldsymbol{\beta})$ that can be consistently estimated by

$$\hat{\Omega}(\boldsymbol{\beta}) = n\sum_{i=1}^{n}\left\{\omega_i(\boldsymbol{\beta})-\bar{\omega}(\boldsymbol{\beta})\right\}\left\{\omega_i(\boldsymbol{\beta})-\bar{\omega}(\boldsymbol{\beta})\right\}'$$

where $\omega_i(\boldsymbol{\beta})=n^{-1}\left(B_{i3}-B_{i4}\right)$, $\bar{\omega}(\boldsymbol{\beta})=n^{-1}\sum_{i=1}^{n}\omega_i(\boldsymbol{\beta})$, and $B_{i3}$ and $B_{i4}$ are defined as

$$
\begin{aligned}
B_{i3} &= \mathcal{A}\int_0^\tau \left\{\mathbf{W}_i^{(2)}-\frac{\widehat{\mathcal{E}}\mathcal{A}\mathbf{W}^{(2)}Y^{(2)}}{\widehat{\mathcal{E}}\mathcal{A}Y^{(2)}}\right\}Y_i^{(2)}\mathrm{d}N_i^{(1)}, \text{ and} \\
B_{i4} &= \mathcal{A}\int_0^\tau \left\{\frac{\mathbf{W}_i^{(2)}Y_i^{(2)}}{\widehat{\mathcal{E}}\mathcal{A}Y^{(2)}}-\frac{(\widehat{\mathcal{E}}\mathcal{A}\mathbf{W}^{(2)}Y^{(2)})Y_i^{(2)}}{(\widehat{\mathcal{E}}\mathcal{A}Y^{(2)})^2}\right\}\times\mathrm{d}\widehat{\mathcal{E}}\mathcal{A}\int Y^{(2)}\mathrm{d}N^{(1)}.
\end{aligned}
$$

# Chapter 3

# Augmented Parametric Corrected Score for Proportional Hazards Model with Covariate Measurement Error

## 3.1   Introduction

The proportional hazards model is one of the most popular models to investigate the relationship between time to failure and covariates. However, in many clinical trials, the true covariates may not always be accurately measured. In some studies,

the magnitude of measurement error could be substantial to the extent that it is comparable to or even larger than that of the true underlying covariate. An example of substantial measurement error is the HIV viral load in HIV/AIDS studies.

As introduced in Chapter 1, available functional modeling methods for Cox proportional hazards model include the conditional score (Tsiatis & Davidian, 2001), the parametric corrected score (Nakamura, 1992; Kong & Gu, 1999), and the nonparametric corrected score (Huang & Wang, 2000, 2006; Hu & Lin, 2004). The idea of the conditional score is to condition away the nuisance parameters based on certain sufficient statistics whereas the last two classes adopt a correction strategy by constructing a corrected estimating function with error-contaminated covariates that shares the same limit as a reference estimating function with true underlying covariates. If the reference estimating function admits consistent estimates only, the corrected estimating function shall inherit this property in a compact parameter space containing the true value. Conditional score and parametric corrected score are generally different. But in the case of the Cox proportional hazards model and normal measurement error, the conditional score and parametric corrected score estimators are asymptotically equivalent.

Although all three aforementioned methods produce consistent estimators, they all suffer from finite-sample pathological behaviors especially when the measurement error is substantial. Though noticed in some literature, the pathological behaviors of these estimation procedures were never well understood. Recently, Huang (2011)

proposed an approach to incorporate additional estimating functions which constrain the derivatives of the parametric corrected score for loglinear model. This approach effectively remedies those pathological behaviors and also considerably improves the estimation efficiency. Huang's approach provides a promising general strategy to handle similarly ill-behaved estimating functions. Motivated by this approach, we will develop an estimation procedure for the Cox proportional hazards model with covariate measurement error.

In this chapter, we first conduct a detailed investigation on pathological behaviors of parametric corrected score and conditional score. After that, we propose an augmented estimation procedure in which additional estimating functions are added to the parametric corrected score for the proportional hazards model. In Section 3.2, we briefly describe the parametric corrected score and conditional score for the proportional hazards model and present the investigation results on the pathological behaviors when covariate measurement error is substantial. The proposed approach of incorporating additional estimating functions for the parametric corrected score is presented in Section 3.3. Simulation studies with practical sample size are reported in Section 3.4 together with an application to the ACTG 175 clinical trial data. Further discussion is given in Section 3.5. Technical details is collected in Section 3.6.

# 3.2 Parametric Corrected Score and Conditional Score and Their Pathological Behaviors

The proportional hazards model postulates that the cumulative hazard function $\Lambda(\cdot)$ of survival time $T$ of an individual with a $p$-vector of covariate $\mathbf{Z}$ has the form

$$\Lambda(dt|\mathbf{Z}) = \exp(\boldsymbol{\beta}'\mathbf{Z})\Lambda_0(dt)$$

where $\boldsymbol{\beta}$ is a $p$-vector parameter of interest and $\Lambda_0(\cdot)$ is an unspecified baseline cumulative hazard function. Let $C$ denotes the censoring time and adopt the usual independent censoring mechanism: given $\mathbf{Z}, C$ is independent of $T$.

The observed data, $(X_i, \Delta_i, \mathbf{Z}_i, i = 1, \ldots, n)$, consist of $n$ iid replicates of $\{X \equiv T \wedge C, \Delta \equiv I(T \leq C), \mathbf{Z}\}$. The standard inference procedure for Cox proportional hazards model is then to maximize the partial likelihood or, equivalently, to solve estimating function

$$\boldsymbol{\xi}^*(\mathbf{b}, \tilde{\Lambda}_0(\cdot)) = n^{-1} \sum_{i=1}^{n} \int_0^\tau \begin{pmatrix} 1 \\ \mathbf{Z}_i \end{pmatrix} \left\{ dN_i(t) - Y_i(t) \exp(\mathbf{b}'\mathbf{Z}_i) d\tilde{\Lambda}_0(t) \right\}, \qquad (3.1)$$

where $N_i(t) = I(X_i \leq t, \Delta_i = 1)$ is the counting process, $Y_i(t) = I(X_i \geq t)$ is the at-risk process and $\tau$ is a positive constant such that $\Pr(T \geq \tau) > 0$. Profiling out

$\tilde{\Lambda}_0(\cdot)$, the estimating function for $\boldsymbol{\beta}$ alone is

$$\boldsymbol{\xi}(\mathbf{b}) = n^{-1} \sum_{i=1}^{n} \int_0^{\tau} \left\{ \mathbf{Z}_i - \frac{\sum_{j=1}^{n} Y_j(t) \mathbf{Z}_j \exp(\mathbf{b}'\mathbf{Z}_j)}{\sum_{j=1}^{n} Y_j(t) \exp(\mathbf{b}'\mathbf{Z}_j)} \right\} dN_i(t). \qquad (3.2)$$

Estimating function (3.2) is actually the usual partial score function.

### 3.2.1 Parametric corrected score and conditional score

Split covariates $\mathbf{Z} = (\mathbf{Z}_a^T, \mathbf{Z}_e^T)^T$ where $\mathbf{Z}_a$ are those covariates that can be accurately measured and $\mathbf{Z}_e$ are covariates prone to measurement error and cannot be accurately measured. Though $\mathbf{Z}_e$ cannot be measured directly, we can observe them through their surrogates $\mathbf{W}_e$. Under the classical additive measurement error model, $\mathbf{W}_e = \mathbf{Z}_e + \boldsymbol{\varepsilon}_e$, where $\boldsymbol{\varepsilon}_e$ is the error vector and $\boldsymbol{\varepsilon}_e$ is assumed to be independent of $(T, C, \mathbf{Z})$. In this chapter, we consider the situation that the distribution of $\boldsymbol{\varepsilon}_e$ is known; The situation where distribution of $\boldsymbol{\varepsilon}_e$ is unknown is studied in Chapter 4.

Let $\mathbf{W} = (\mathbf{Z}_a, \mathbf{W}_e)$ and $\boldsymbol{\varepsilon} = (\mathbf{0}, \boldsymbol{\varepsilon}_e)$. The observed data now consist of $(X_i, \Delta_i, \mathbf{W}_i, i = 1, \dots, n)$ in the presence of covariate measurement error. It is well known that naively replacing $\mathbf{Z}$ by $\mathbf{W}$ in estimating functions (3.1) or (3.2) could incur substantial estimation bias. Denote the cumulant-generating function of $\boldsymbol{\varepsilon}$ as $\Omega(\mathbf{b}) \equiv \log \mathcal{E}\{\exp(\mathbf{b}'\boldsymbol{\varepsilon})\}$ and its derivative $\dot{\Omega}(\mathbf{b}) \equiv \partial\Omega(\mathbf{b})/\partial\mathbf{b}$. The parametric corrected score estimating

function is given by

$$
\begin{aligned}
\boldsymbol{\eta}^*(\mathbf{b}, \tilde{\Lambda}_0(\cdot)) \;=\; & n^{-1} \sum_{i=1}^{n} \int_0^\tau \binom{1}{\mathbf{W}_i - \dot{\Omega}(0)} dN_i(t) \\
& - \int_0^\tau Y_i(t) \exp\{\mathbf{b}'\mathbf{W}_i - \Omega(\mathbf{b})\} \binom{1}{\mathbf{W}_i - \dot{\Omega}(\mathbf{b})} d\tilde{\Lambda}_0(t).
\end{aligned}
\tag{3.3}
$$

Further profiling out $\tilde{\Lambda}_0(\cdot)$, we obtain the corrected score for $\boldsymbol{\beta}$ (Nakamura, 1992):

$$
\boldsymbol{\eta}(\mathbf{b}) = n^{-1} \sum_{i=1}^{n} \int_0^\tau \left\{ \mathbf{W}_i + \dot{\Omega}(\mathbf{b}) - \dot{\Omega}(0) - \frac{\sum_{j=1}^{n} Y_j(t)\mathbf{W}_j \exp(\mathbf{b}'\mathbf{W}_j)}{\sum_{j=1}^{n} Y_j(t)\exp(\mathbf{b}'\mathbf{W}_j)} \right\} dN_i(t)
\tag{3.4}
$$

which has the same limit as reference (3.2) asymptotically for each and every finite $\mathbf{b}$. The estimation is then to find the zero crossing of the above estimating function. The consistency and asymptotic normality of corrected score estimator are later established by Kong & Gu (1999).

A root-consistent estimating function is an estimating function such that every zero-crossing is consistent and a normalized estimating function is root-consistent if its limit has a unique root at the estimand (Huang and Wang, 1999). By definition, reference (3.2) is a root-consistent estimating function and the new estimating function (3.4) shall inherit the root-consistency from (3.2). The root-consistency of (3.4) assures that in a compact parameter space containing the true parameter $\boldsymbol{\beta}$, the parametric corrected score will admit a unique root asymptotically and the root is consistent and asymptotically normal.

As for the conditional score estimating function for the Cox proportional haz-

ards model (Tsiatis & Davidian, 2001), when the measurement error is normally distributed and the variance matrix is $\boldsymbol{\Sigma}$, it can be written as:

$$
\begin{aligned}
\boldsymbol{\eta}_{con}(\mathbf{b}) \; = \; & n^{-1} \sum_{i=1}^{n} \int_{0}^{\tau} \Big\{ \mathbf{W}_i + \boldsymbol{\Sigma}\mathbf{b} \\
& - \frac{\sum_{j=1}^{n} Y_j(t)[\mathbf{W}_j + \boldsymbol{\Sigma}\mathbf{b}dN_j(t)] \exp(\mathbf{b}'[\mathbf{W}_j + \boldsymbol{\Sigma}\mathbf{b}dN_j(t)])}{\sum_{j=1}^{n} Y_j(t) \exp(\mathbf{b}'[\mathbf{W}_j + \boldsymbol{\Sigma}\mathbf{b}dN_j(t)])} \Big\} \, dN_i(t)
\end{aligned}
\tag{3.5}
$$

In fact, it can be shown that estimators from parametric corrected score and conditional score are asymptotically equivalent in the case of normal measurement error (Tsiatis & Davidian, 2001).

## 3.2.2 Pathological behaviors

When the magnitude of measurement error is small, the asymptotic results of parametric corrected score and conditional score provide a good approximation for practical purposes. However, when the measurement error increases, the pathological behaviors may start to arise (Song & Huang, 2005). These pathological behaviors include multiple zero-crossings or a single wrong zero-crossing that is inappropriate. These pathological behaviors may cause serious concerns when the measurement error is substantial and limit the applicability of parametric corrected score and conditional score in practice. But these pathological behaviors were never well understood or thoroughly investigated in the literature. In this section, we will conduct a detailed investigation of pathological behaviors for these two methods.

We first consider the parametric corrected score. Consider a single-covariate model with normal measurement error. In the absence of measurement error, the partial likelihood score function $\xi(b)$ is monotonically decreasing and has a unique root. The asymptotic result suggests that the parametric corrected score $\eta(b)$ should be monotonically decreasing as well in a compact parameter space containing the true parameter when the sample size is large. One may speculate the parametric corrected score to have an overall decreasing trend over the entire parameter space. But surprisingly, the overall trend of parametric corrected score in an unbounded parameter space is increasing. Function $\eta(b)$ takes a value of $-\infty$ when $b = -\infty$ and a value of $+\infty$ when $b = +\infty$. This observation suggests that the parametric corrected score $\eta(b)$ has an odd number of roots. In our numerical studies, only single- and triple-root patterns have been observed and two typical plots of parametric corrected score are illustrated in Figure 3.1. The same root patterns were observed in Huang's (2011) investigation of loglinear model.

An alternative way of constructing corrected score is to first construct corrected partial likelihood and then take derivative. So if we characterize a root by increase or decrease of $\eta(\cdot)$ around it, the increasing and decreasing roots correspond to local minimizers and maximizers of corresponding corrected partial likelihood function, respectively. Therefore, an increasing root is considered as an inappropriate one. In the case of single-root pattern, the only root is increasing and thus inappropriate. For the triple-root pattern, an appropriate root exists since there is only one decreasing root. The single-root pattern is considered as root-finding failure.

Figure 3.1: Observed root patterns of the parametric corrected score $\eta(b)$. The true $\beta$ is -1 and these two corrected curves correspond to the same profile score (with true covariates). Portion of a corrected curve is thickened to indicate negative derivative.

Table 3.1: Prevalence (%) of single-root pattern for the parametric corrected score with $E(X) = 0$, $\mathrm{Var}(X) = 1$, $\beta = -1$, and $\varepsilon \sim \mathrm{Normal}(0,1)$.

| Censoring Rate | Distribution of X | Size | | | |
|---|---|---|---|---|---|
| | | 100 | 200 | 400 | 800 |
| 20% | Normal | 68.0 | 52.1 | 38.0 | 20.3 |
| | Modified Chi-square | 65.8 | 57.9 | 52.3 | 40.1 |
| | Uniform | 64.5 | 51.9 | 37.6 | 19.6 |
| 40% | Normal | 62.4 | 49.0 | 36.3 | 17.8 |
| | Modified Chi-square | 64.4 | 58.0 | 50.9 | 43.1 |
| | Uniform | 60.4 | 49.5 | 37.1 | 19.3 |
| 60% | Normal | 61.0 | 46.5 | 32.1 | 16.0 |
| | Modified Chi-square | 62.7 | 57.3 | 50.9 | 42.6 |
| | Uniform | 58.5 | 47.9 | 36.8 | 21.2 |

We conduct a simulation study to exam the prevalence of single-root pattern. We consider a single covariate Z and generate it from various distributions: A) standard normal distribution, B) modified chi-square distribution, and C) uniform distribution with mean 0 and variance 1. To generate the modified chi-square distribution, the chi-square distribution with 1 degree of freedom was first truncated at 5 and then location-shifted to mean 0 and rescaled to variance 1. The measurement error follows standard normal distribution. The true coefficient was taken to be $\beta = -1$ and the baseline hazard is constant 1. Censoring was generated from a uniform distribution on $[0, \mu]$, where $\mu$ is chosen so that the censoring rate is ranged from 20% to 60%. These set-ups represent a practical scenario with substantial error contamination on the covariate. The results based on 1,000 iterations are reported in Table 3.1. The prevalence of single-root patten is similar across different scenarios with various censoring rates. When the sample size is 100, the percentage of single-root pattern under all three distributions are close to or over 60%. Even when the sample size increases to 800, the prevalence of single-root pattern is still quite high.

For the conditional score, the patterns are more complicated. When the absolute value of $\beta$ gets large, the estimating function fluctuates around zero and finally approaches zero as $\beta$ goes to infinity. Therefore the conditional score may have many zero-crossings. When the $\beta$ is not so extreme, two general patterns for the conditional scores are observed and plots from two simulation data sets are shown in Figure 3.2. As explained in previous section, conditional score and corrected score are asymptotically equivalent. Therefore an appropriate zero-crossing for conditional

score should be a decreasing one as well. In the first pattern, the conditional score has one zero-crossing close to the true parameter. This single zero-crossing is decreasing and thus an appropriate one. In the second patten, the conditional score appears to have no proper zero-crossing near the truth though it may have multiple zero-crossings at extreme values of $\beta$. For conditional score, we define the estimator as the local point having the smallest $\ell^2$-norm in the neighborhood of the naive estimation. If the estimator is not a zero-crossing, we consider it as root-finding failure. Table 3.2 summarizes the prevalence of root-finding failure for the conditional score. The same simulation set up as in the corrected score is used. With a sample size of 100, the root-finding failure rate for conditional score varies from 3% to 5% for normal covariate and from 14% to 25% for modified chi-square covariate. As the sample size increases to 800, the failure rate drops to 0.2% for normal covariate but remains at



Figure 3.2: Observed root patterns of the conditional score estimating function. The true $\beta$ is -1.

least 14% for modified chi-square covariate.

Table 3.2: Prevalence (%) of root-finding failure for the conditional score with $E(X) = 0$, $\text{Var}(X) = 1$, $\beta = -1$, and $\varepsilon \sim \text{Normal}(0,1)$.

| Censoring Rate | Distribution of X | Size | | | |
|---|---|---|---|---|---|
| | | 100 | 200 | 400 | 800 |
| 20% | Normal | 2.7 | 1.5 | 0.6 | 0.3 |
| | Modified Chi-square | 14.4 | 12.2 | 15.1 | 13.9 |
| | Uniform | 3.1 | 1.6 | 0.6 | 0.1 |
| 40% | Normal | 4.6 | 2.3 | 0.8 | 0.2 |
| | Modified Chi-square | 20.0 | 16.7 | 20.9 | 18.7 |
| | Uniform | 3.1 | 2.1 | 0.3 | 0.5 |
| 60% | Normal | 4.7 | 2.5 | 0.7 | 0.3 |
| | Modified Chi-square | 24.6 | 21.7 | 25.7 | 25.0 |
| | Uniform | 5.7 | 2.7 | 1.4 | 1.0 |

The above investigation results show that, in the presence of substantial measurement error, both parametric corrected score and conditional score suffer from severe finite-sample pathological behaviors. Therefore improvements are required for these methods to have practical applicability. We observe that an appropriate zero-crossing of an estimating function should be a decreasing one. This observation suggests that the trend of estimating function is also informative and could be taken into account in the estimate determination. By Taylor expansion, the trend of estimating function may be quantified by its derivative. Recognizing this feature, Huang (2011) proposed an approach to incorporate additional estimating functions which constrain the derivatives of the corrected score for the loglinear model. The estimation and inference are then accomplished by means of empirical likelihood. This approach effectively remedies the pathological behaviors of corrected score for loglinear model

and also considerably improves the estimation efficiency. However, in the case of the Cox proportional hazards model, we are unable to construct additional estimating functions that effectively constrain the derivative of the parametric corrected score or conditional score because of the very nature of these two estimating functions. Nevertheless, Huang's approach provides an insight into a new approach to address pathological behaviors of estimating functions. If we could identify additional estimating functions that do not share the same wrong roots as the original estimating function, then by combining the original and additional estimating functions, pathological behaviors could be reduced or eliminated. But if either the original estimating function or the additional estimating functions vanish to zero, then wrong root sharing would easily arise. We have shown in the simulation that the conditional score would vanish to zero when the absolute value of parameter becomes large. Thus, the trend pattern of the parametric corrected score is more desirable than that of the conditional score and our method will be developed for the parametric corrected score.

## 3.3   Improving Corrected Score

### 3.3.1   Augmented Estimation Method

Motivated by Huang's (2011) trend-constrained corrected score, we first establish the following result:

**Theorem 3.1.** *Under the proportional hazards model and the classical additive mea-*

*surement error model,*

$$
\mathcal{E}\left[\int_0^\tau \frac{\partial^{k_1+\cdots+k_p}\exp\{\mathbf{b}'\mathbf{W}-\Omega(\mathbf{b})\}}{\partial b_1^{k_1}\cdots b_p^{k_p}}dN_i(t)\bigg|_{\mathbf{b}=\mathbf{0}}\right. \tag{3.6}
$$
$$
\left.\int_0^\tau \frac{\partial^{k_1+\cdots+k_p}Y(t)\exp\{\mathbf{b}'\mathbf{W}-\Omega(\mathbf{b})\}}{\partial b_1^{k_1}\cdots b_p^{k_p}}d\Lambda_0(t)\bigg|_{\mathbf{b}=\boldsymbol{\beta}}\right]=0,
$$

*for $k_l \geq 0, l = 1,\ldots,p$, where $b_l, l = 1,\ldots,p$, is the lth element of $\mathbf{b}$.*

Proof. Given

$$
\mathcal{E}[\exp\{\mathbf{b}'\mathbf{W}-\Omega(\mathbf{b})\}|\mathbf{Z}] = \exp(\mathbf{b}'\mathbf{Z})
$$

under additive measurement error model, one may obtain

$$
\mathcal{E}\left[\int_0^\tau \frac{Y\partial^{k_1+\cdots+k_p}\exp\{\mathbf{b}'\mathbf{W}-\Omega(\mathbf{b})\}}{\partial b_1^{k_1}\cdots b_p^{k_p}}d\Lambda_0(t)\bigg|\mathbf{Z}\right] = \prod_{l=1}^p Z_l^{k_l}\int_0^\tau Y(t)\exp(\mathbf{b}'\mathbf{Z})d\Lambda_0(t)
$$

and

$$
\mathcal{E}\left[\int_0^\tau \frac{\partial^{k_1+\cdots+k_p}\exp\{\mathbf{b}'\mathbf{W}-\Omega(\mathbf{b})\}}{\partial b_1^{k_1}\cdots b_p^{k_p}}dN(t)\bigg|_{\mathbf{b}=\mathbf{0}}\bigg|\mathbf{Z}\right] = \prod_{l=1}^p Z_l^{k_l}N(\tau).
$$

Then given the fact that $M(t) = N(t) - \int_0^t Y(u)\exp(\boldsymbol{\beta}'\mathbf{Z})d\Lambda_0(u)$ is a mean zero

martingale, equation (3.6) is implied by the above two equations.

Equation (3.6) is useful in constructing additional estimating equations for $\boldsymbol{\beta}$.

When $\sum_{l=1}^p k_l = 0$ and 1, one may obtain the usual parametric corrected score. When

$\sum_{l=1}^p k_l = 2$, the additional estimating functions are the upper triangular elements of

the following symmetric matrix:

$$n^{-1} \sum_{i=1}^{n} \int_0^{\tau} \left[ \{\mathbf{W}_i - \dot{\Omega}(0)\}^{\otimes 2} - \ddot{\Omega}(0) \right] dN_i(t)$$
$$- \int_0^{\tau} Y_i(t) \exp\{\mathbf{b}'\mathbf{W}_i - \Omega(\mathbf{b})\} \left[ \{\mathbf{W}_i - \dot{\Omega}(\mathbf{b})\}^{\otimes 2} - \ddot{\Omega}(\mathbf{b}) \right] d\tilde{\Lambda}_0(t).$$

By profiling out $\tilde{\Lambda}_0(\cdot)$, we obtain

$$n^{-1} \sum_{i=1}^{n} \int_0^{\tau} \left[ \{\mathbf{W}_i - \dot{\Omega}(0)\}^{\otimes 2} - \ddot{\Omega}(0) \right. \tag{3.7}$$
$$\left. - \frac{\sum_{j=1}^{n} Y_j(t) \exp\{\mathbf{b}'\mathbf{W}_j - \Omega(\mathbf{b})\} \left[ \{\mathbf{W}_j - \dot{\Omega}(\mathbf{b})\}^{\otimes 2} - \ddot{\Omega}(\mathbf{b}) \right]}{\sum_{i=j}^{n} Y_j(t) \exp\{\mathbf{b}'\mathbf{W}_j - \Omega(\mathbf{b})\}} \right] dN_i(t)$$

The additional estimating functions would be helpful if both parametric corrected score and additional estimating function are close to 0 around the truth (not necessarily having roots) and the additional estimating function is not close to 0 when the parametric corrected score is close to 0 at any point far away from the truth. Figure 3.3 shows four typical patterns of parametric corrected score and corresponding additional estimating function based on equation (3.7) for a single-covariate model with true parameter $\beta = -1$. Plot (a) is the ideal scenario. The parametric corrected score have three zero-crossings. The additional estimating function shares the same decreasing zero-crossing as the parametric corrected score. Moreover, they do not share any wrong zero-crossings. In plot (b), the parametric corrected score have three zero-crossings and two of them are close to each other. In this case, discrimi-

Figure 3.3: Parametric corrected score and additional estimation function based on equation (3.7) for a single-covariate model with $\beta = -1$.

nating the two roots around the truth is not very important since they are close to each other anyway. In plots (c) and (d), the parametric corrected score has a single wrong zero-crossing but the additional estimating function is not close to 0 near this wrong zero-crossing. Both estimating functions are close to 0 around the truth.

## 3.3.2   Estimation and inference

With the additional estimating functions, we have more estimating functions than the number of parameters. Available methods to synthesize estimating functions that exceed the number of parameters include empirical likelihood (Qin & Lawless, 1994) and quadratic inference function (QIF) method (Lindsay & Qu, 2003). In this research, we shall use the quadratic inference function method to determine the estimate since the estimating functions are not sums of iid terms, thus the empirical likelihood would be computational difficult.

Let $\boldsymbol{\varphi}(\mathbf{b})$ denotes the estimating functions. $\boldsymbol{\varphi}(\mathbf{b})$ is comprised of the original parametric corrected score and additional estimating functions. The quadratic inference function takes the form

$$Q(\mathbf{b}; \hat{\mathbf{C}}) = \boldsymbol{\varphi}'(\mathbf{b})\hat{\mathbf{C}}^{-1}(\mathbf{b})\boldsymbol{\varphi}(\mathbf{b}), \tag{3.8}$$

where $\hat{\mathbf{C}}(\mathbf{b})$ is any consistent estimator for the asymptotic variance of $\boldsymbol{\varphi}(\mathbf{b})$. Then the estimator is defined as the minimizer of (3.8) and is consistent for the true value

of **b**. Furthermore, the estimator is asymptotically efficient in the class of consistent estimators based on linear combination of parametric corrected score and additional estimating functions. The construction of a quadratic inference function helps to solve both aspects of pathological behaviors of the parametric corrected score. Firstly, in the case of multiple zero-crossings, the introduction of additional estimating functions helps to pick up the right zero-crossing out from multiple ones if the additional estimating functions do not share the same wrong roots as the original parametric corrected score. Secondly, the problem of no appropriate zero-crossing could be solved by minimizing the quadratic inference function.

We name the method to incorporate additional estimating functions based on (3.6) as *the augmented parametric corrected score*. One important special case is the method incorporating the upper triangular elements of matrix (3.7) and was termed *the second order augmented parametric corrected score*. Augmented parametric corrected scores with higher order are also available, with additional estimating functions corresponding to $k_l$ such that $\sum_{l=1}^{p} k_l > 2$. In Section 3.4, we will conduct extensive simulation studies to evaluate the performance of the augmented parametric corrected score and its applicability in practice.

Figure 3.4 plots quadratic inference functions corresponding to the two datasets in Figure 3.1. The second order augmented parametric corrected score is adopted. The formula for a consistent estimator of the asymptotic variance of $\boldsymbol{\varphi}(\mathbf{b})$ is given in Section 3.6.

Figure 3.4: Quadratic inference functions for the two datasets in Figure 3.1.

The interval estimation can be easily achieved by inverting the hypothesis testing statistics. As an inference function, $Q(\mathbf{b}; \hat{\mathbf{C}})$ has similar properties as the log-likelihood function (Lindsay & Qu, 2003):

(a) $Q(\boldsymbol{\beta}_0) - Q(\hat{\boldsymbol{\beta}})$ is asymptotically chi-squared with degree of freedom $p$;

(b) the profile test statistics $Q(\boldsymbol{\psi}_0, \hat{\boldsymbol{\lambda}}_0) - Q(\hat{\boldsymbol{\psi}}, \hat{\boldsymbol{\lambda}})$, where $(\boldsymbol{\psi}, \boldsymbol{\lambda})$ is a partition of the parameter of $\boldsymbol{\beta}$ and $\hat{\boldsymbol{\lambda}}_0$ is the profiled estimate of $\boldsymbol{\lambda}$ given $\boldsymbol{\psi} = \boldsymbol{\psi}_0$, is asymptotically chi-squared as a test of $H : \boldsymbol{\psi} = \boldsymbol{\psi}_0$, with degree of freedom equal to the dimension of $\boldsymbol{\psi}$.

## 3.4 Numerical Studies

### 3.4.1 Simulations

Extensive simulations studies were conducted to evaluate the performance of augmented parametric corrected score. For reference and comparison, the ideal, naive, regression calibration, conditional score, and parametric corrected score were also studied. The ideal estimator used the ordinary partial score function with true covariates and, of course, it is not a realistic estimator. The naive approach uses the mismeasured surrogates in place of true covariates in the partial score function. For the regression calibration method, $\mathbf{Z}_e$ is replaced by $E(\mathbf{Z}_e|\mathbf{Z}_a, \mathbf{W}_e)$ in the partial score function. For the proposed approach, the optimization algorithm of Nelder & Mead (1965) will be used.

As shown in previous section, both conditional score and parametric corrected score have high prevalence of root-finding failure. To permit a fair comparison, we propose the following re-defined conditional score and re-defined parametric corrected score. If an appropriate zero-crossing could be found, the estimators will take the value of zero-crossing. But if root-finding failure occurs, the estimators will be defined as the local minimizer in the $\ell^2$-norm of the estimating functions closest to the naive estimator. Operationally, we will use the following modified Newton-Raphson algorithm. We start with the naive estimator and calculate the Newton-Raphson step size. Since the goal is to find a root or local minimizer, we need prevent overshooting.

In our simulation, we cap the step size at 0.2. During each iteration, we compare the $\ell^2$-norm evaluated at new estimator to that evaluated at current estimator. If the $\ell^2$-norm evaluated at new estimator is smaller, then the new estimator will be accepted and the algorithm continues to the next iteration. Otherwise, we will halve the step size and calculate the new estimator again. Iterations will be repeated until that i) the absolute value of estimating function is less than $10^{-6}$ or ii) the step size has been halved for more than 10 time during any single iteration. If criteria i) is satisfied, the algorithm converges to a zero-crossing and a root is identified in this case. If criteria ii) is satisfied, the algorithm converges to a local minimizer and root-finding failure occurs. Comparing to the original definition that estimators take the values of zero-crossings of estimating function, this new definition actually benefits the conditional score and corrected score. As shown in Figure 3.1 and 3.2, in the case of root-finding failure, zero-crossings of conditional score and corrected score are at extreme and far away from the true value.

In the simulation, we consider both single- and double-covariate models. In the single-covariate models, the true covariate X is of mean 0 and variance 1. The true regression coefficient was set to be -1 and the baseline hazard is constant 1. The measurement error follows the standard normal distribution. Two different distributions of X were studied: A) standard normal distribution and B) modified chi-square distribution. To generate the modified chi-square distribution, the chi-square distribution with 1 degree of freedom was first truncated at 5 and then location-shifted to mean 0 and rescaled to variance 1. Censoring time was generated from a uniform

distribution on $[0, \mu]$ and we will consider two different settings of censoring rate at 20% and 60%.

In the double-covariate models, true covariate $X$ follows bivariate normal distribution with mean $(-1, 1)$ and correlation coefficient of 0.5. The first covariate was subject to a standard normal measurement error, whereas the second covariate was accurately measured. The regression coefficients were set to $(-1, 1)$ and the baseline hazard is constant 1.. Censoring time was also generated from a uniform distribution on $[0, \mu]$, where $\mu$ is chosen so that the censoring rate is 20% or 60%.

Sample sizes 100, 200, 400, 800, and 1,600 were investigated. For each scenario, 1,000 samples were simulated. We report the results on point and interval estimation separately.

Table 3.3 and 3.4 summarize the simulation results on the estimators in the single-covariate models with censoring rates of 20% and 60% respectively. The quantile-quantile plots are shown in Figure 3.4 and 3.5. For each scenario, the mean bias, and standard deviation were calculated. For augmented parametric corrected score, three sets of additional estimating functions were considered where $k = 2, 3, 4$. As expected, the naive estimator has substantial bias under both scenarios. The regression calibration estimator shows moderate bias, with larger bias in modified chi-square covariate case than normal covariate case. The re-defined conditional score shows slight bias under both scenarios, probably due to its left skewness. The quantile-quantile plots show that the re-defined conditional score deviates from normality considerably even

Table 3.3: Simulation summary statistics for the single-covariate models with 20% censoring rate: Ideal, naive(NV), regression calibration (RC), re-defined conditional score (ConS), re-defined parametric corrected score (CS), and first k augmented parametric corrected score (ACS:k), k = 2, 3, 4.

| | Ideal | | NV | | RC | | ConS | | | CS | | | ACS: 2 | | ACS: 3 | | ACS: 4 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| size | B | SD | B | SD | B | SD | F | B | SD | F | B | SD | B | SD | B | SD | B | SD |
| | | | | | | | | | Case A | | | | | | | | | |
| 100 | -15 | 154 | 593 | 99 | 148 | 254 | 2.7 | -376 | 920 | 68.0 | -11 | 287 | 46 | 431 | -43 | 414 | -103 | 381 |
| 200 | -8 | 104 | 602 | 65 | 189 | 163 | 1.5 | -251 | 719 | 52.1 | -89 | 273 | -60 | 330 | -130 | 348 | -172 | 351 |
| 400 | -5 | 73 | 601 | 48 | 195 | 113 | 0.6 | -202 | 575 | 38.0 | -146 | 282 | -99 | 287 | -153 | 313 | -184 | 327 |
| 800 | 1 | 49 | 604 | 32 | 206 | 76 | 0.3 | -99 | 351 | 20.0 | -129 | 267 | -87 | 246 | -83 | 226 | -123 | 256 |
| 1600 | 1 | 36 | 605 | 23 | | | 0 | -43 | 200 | 5.6 | -75 | 206 | -45 | 163 | -60 | 178 | -62 | 166 |
| | | | | | | | | | Case B | | | | | | | | | |
| 100 | -25 | 204 | 639 | 91 | 242 | 247 | 14.4 | -485 | 1221 | 65.8 | 145 | 241 | 110 | 387 | 41 | 409 | 9 | 409 |
| 200 | -5 | 136 | 640 | 63 | 267 | 155 | 12.2 | -396 | 1022 | 57.9 | 47 | 241 | -10 | 351 | -81 | 381 | -128 | 381 |
| 400 | -6 | 98 | 642 | 43 | 280 | 101 | 15.1 | -329 | 895 | 52.3 | -21 | 237 | -10 | 285 | -76 | 306 | -122 | 323 |
| 800 | -3 | 70 | 641 | 32 | 281 | 72 | 13.9 | -249 | 740 | 40.1 | -70 | 245 | -34 | 241 | -81 | 237 | -95 | 241 |
| 1600 | 0 | 47 | 644 | 21 | | | 12.5 | -139 | 486 | 30.3 | -82 | 242 | -9 | 152 | -29 | 133 | -42 | 139 |

Note: F: root-finding failure (%); B: mean bias ($\times 10^3$); SD: standard deviation($\times 10^3$).

Table 3.4: Simulation summary statistics for the single-covariate models with 60% censoring rate: Ideal, naive(NV), regression calibration (RC), re-defined conditional score (ConS), re-defined parametric corrected score (CS), and first k augmented parametric corrected score (ACS:k), k = 2, 3, 4.

| size | Ideal B | SD | NV B | SD | RC B | SD | ConS F | B | SD | CS F | B | SD | ACS: 2 B | SD | ACS: 3 B | SD | ACS: 4 B | SD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Case A | | | | | | | | | | |
| 100 | -24 | 206 | 550 | 134 | 57 | 335 | 4.7 | -429 | 1014 | 61.0 | -27 | 301 | 94 | 478 | -24 | 502 | -90 | 551 |
| 200 | -15 | 142 | 562 | 88 | 109 | 206 | 2.5 | -272 | 738 | 46.5 | -101 | 303 | 10 | 298 | -99 | 347 | -188 | 395 |
| 400 | -8 | 97 | 563 | 63 | 119 | 142 | 0.7 | -197 | 578 | 32.1 | -139 | 305 | -27 | 235 | -96 | 257 | -152 | 294 |
| 800 | -2 | 65 | 567 | 42 | 132 | 96 | 0.3 | -100 | 379 | 16.0 | -113 | 271 | -33 | 189 | -72 | 199 | -103 | 223 |
| 1600 | 1 | 46 | 567 | 30 | | | 0 | -51 | 238 | 4.9 | -68 | 205 | -17 | 128 | -38 | 125 | -57 | 148 |
| | | | | | | | | Case B | | | | | | | | | | |
| 100 | -66 | 351 | 687 | 124 | 351 | 288 | 24.6 | -259 | 1154 | 62.7 | 247 | 297 | 196 | 459 | 219 | 531 | 260 | 568 |
| 200 | -24 | 239 | 691 | 84 | 370 | 190 | 21.7 | -291 | 1042 | 57.3 | 133 | 287 | 81 | 359 | 44 | 429 | 54 | 496 |
| 400 | -16 | 159 | 693 | 57 | 382 | 124 | 25.7 | -202 | 841 | 50.9 | 57 | 268 | 37 | 322 | -22 | 368 | -52 | 423 |
| 800 | -9 | 112 | 693 | 40 | 385 | 86 | 25.0 | -172 | 684 | 42.6 | -8 | 263 | 1 | 280 | -55 | 276 | -94 | 306 |
| 1600 | 0 | 78 | 697 | 28 | | | 21.5 | -77 | 412 | 34.0 | -29 | 258 | -3 | 232 | -29 | 194 | -57 | 207 |

Note: Same as in that of Table 3.3

Figure 3.5: Quantile-quantile plots for $\beta$ in the single-covariate models with 20% censoring rate, where $\beta = -1$. Red, yellow, green, blue, and black correspond to sample sizes 100, 200, 400, 800, and 1,600.

Figure 3.6: Quantile-quantile plots for $\beta$ in the single-covariate models with 60% censoring rate, where $\beta = -1$. Red, yellow, green, blue, and black correspond to sample sizes 100, 200, 400, 800, and 1,600.

when the sample size is 1,600. It also has a much larger standard deviation comparing to the re-defined parametric corrected score and the augmented parametric corrected score. The re-defined parametric corrected score is unbiased for both scenarios and the standard deviation is small. All three augmented corrected scores are consistent under both scenarios, and they become less biased as sample size increases. In comparison with higher-order augmented corrected scores, the second order augmented corrected score seems more favorable overall in terms of bias and standard error. Compared to the re-defined parametric corrected score, the second order augmented corrected score has a larger standard deviation when the sample size is small. But the standard deviation of augmented corrected score decreases rapidly as the sample size increases and is smaller than that of the re-defined corrected score when then sample size is 1,600.

Table 3.5 and 3.6 show the simulation results for the double-covariate models. As expected, when two covariates are correlated, the measurement error generally has impact not only on the mismeasured covariate but also on that of the accurately measured one. The relative performance of all estimators is similar to what was observed in the single-covariate models. For multiple-covariate models, each set of additional estimating functions contains more than one element. For example, when $k = 2$, the additional estimating functions includes three elements in the upper triangle of matrix (3.7). Various augmented corrected score could be constructed depending on the additional estimating functions chosen. Since the second order augmented corrected score shows a more favorable overall performance in the single-covariates

models, we consider only the second order augmented corrected score and augmented corrected scores using a subset of the three elements in this simulation. Simulation results show that the second order augmented corrected score performs better than two other augmented corrected scores. The bias of the second order augmented corrected score reduced quickly as sample size increased. Meanwhile, it also has the smallest standard deviation.

Table 3.5: Simulation summary statistics for the double-covariate models with 20% censoring rate: Ideal, naive(NV), regression calibration (RC), re-defined conditional score (ConS), re-defined parametric corrected score (CS), and augmented parametric corrected score (ACS:1/3, ACS:2/3, ACS:2)

| size | | Ideal B SD | NV B SD | RC B SD | ConS F | B SD | CS F | B SD | ACS: 1/3 B SD | ACS: 2/3 B SD | ACS: 2 B SD |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 100 | $\beta_1$ | -23 175 | 637 101 | 66 381 | 57.5 | 215 348 | 72.0 | 267 322 | 170 532 | 245 514 | 308 538 |
| | $\beta_2$ | 23 171 | -400 142 | -113 246 | | -67 336 | | -100 327 | -60 456 | -97 404 | -142 399 |
| 200 | $\beta_1$ | -13 114 | 648 66 | 136 224 | 46.9 | 155 315 | 62.8 | 186 294 | -18 382 | 42 375 | 92 339 |
| | $\beta_2$ | 16 117 | -412 101 | -154 161 | | -35 310 | | -47 286 | 54 349 | 18 327 | -20 276 |
| 400 | $\beta_1$ | -5 82 | 649 48 | 168 144 | 33.4 | 65 297 | 46.2 | 71 280 | -89 301 | -56 277 | -28 255 |
| | $\beta_2$ | 4 80 | -416 68 | -177 103 | | 5 257 | | 15 243 | 85 269 | 60 233 | 40 210 |
| 800 | $\beta_1$ | -2 55 | 651 33 | 181 93 | 20.0 | 22 224 | 30.5 | -3 252 | -99 282 | -80 262 | -56 220 |
| | $\beta_2$ | 1 56 | -420 52 | -185 74 | | 16 186 | | 44 203 | 88 255 | 72 230 | 52 186 |
| 1600 | $\beta_1$ | -2 40 | 652 23 | | 8.0 | -13 178 | 10.5 | -43 207 | -71 219 | -67 204 | -61 190 |
| | $\beta_2$ | 1 39 | -419 34 | | | 21 134 | | 45 155 | 59 180 | 57 163 | 51 145 |

F: root-finding failure (%); B: mean bias ($\times 10^3$); SD: standard deviation($\times 10^3$). ACS: 1/3 and ACS: 2/3 correspond to the additional estimating functions containing the (1,1) element and the first row elements of matrix (3.7), respectively. ACS: 2 is the second order augmented parametric corrected score.

Table 3.8 and 3.8 report the coverage of three types of 95% confidence intervals

Table 3.6: Simulation summary statistics for the double-covariate models with 60% censoring rate: Ideal, naive(NV), regression calibration (RC), re-defined conditional score (ConS), re-defined parametric corrected score (CS), and augmented parametric corrected score (ACS:1/3, ACS:2/3, ACS:2)

| size | | Ideal | | NV | | RC | | ConS | | | CS | | | ACS: 1/3 | | ACS: 2/3 | | ACS: 2 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | B | SD | B | SD | B | SD | F | B | SD | F | B | SD | B | SD | B | SD | B | SD |
| 100 | $\beta_1$ | -32 | 242 | 610 | 140 | -13 | 513 | 43.3 | 96 | 428 | 64.8 | 239 | 320 | 225 | 668 | 326 | 666 | 392 | 660 |
| | $\beta_2$ | 33 | 245 | -352 | 206 | -40 | 336 | | -27 | 408 | | -72 | 362 | -50 | 536 | -104 | 500 | 147 | 465 |
| 200 | $\beta_1$ | -20 | 160 | 617 | 95 | 59 | 286 | 36.1 | 29 | 393 | 57.8 | 145 | 300 | 14 | 387 | 93 | 360 | 140 | 344 |
| | $\beta_2$ | 21 | 161 | -362 | 138 | -81 | 209 | | 24 | 344 | | -8 | 301 | 48 | 369 | -7 | 310 | -41 | 264 |
| 400 | $\beta_1$ | -6 | 109 | 619 | 65 | 97 | 182 | 21.5 | -2 | 293 | 40.4 | 45 | 279 | -65 | 321 | -9 | 273 | 30 | 241 |
| | $\beta_2$ | 6 | 111 | -366 | 97 | -106 | 137 | | 31 | 225 | | 36 | 254 | 79 | 301 | 35 | 246 | 6 | 197 |
| 800 | $\beta_1$ | -1 | 72 | 621 | 45 | 112 | 120 | 10.8 | -23 | 249 | 25.7 | -16 | 257 | -73 | 267 | -32 | 213 | -9 | 192 |
| | $\beta_2$ | 2 | 76 | -369 | 69 | -115 | 94 | | 34 | 198 | | 55 | 214 | 72 | 236 | 40 | 189 | 23 | 160 |
| 1600 | $\beta_1$ | -2 | 53 | 622 | 33 | | | 7.6 | -24 | 186 | 11.3 | -46 | 212 | -55 | 193 | -40 | 171 | -29 | 162 |
| | $\beta_2$ | 0 | 53 | -371 | 46 | | | | 28 | 142 | | 47 | 161 | 46 | 162 | 36 | 140 | 26 | 124 |

Note: Same as in that of Table 3.5

in the single- and double-covariate models. All three scenarios use the second order augmented corrected score. In constructing the confidence interval, we use two different approaches: inverting the hypothese testing statistics as introduced in Section 3.2.2, and the Wald-type confidence interval. For the former, we use two critical values based on the asymptotic chi-square distribution and the bootstrap calibration (Efron & Tibshirani, 1993, chap. 12). Bootstrap size of 500 is used for the bootstrap calibration. The Wald-type confidence interval and test based one with chi-square distribution critical values have poor coverage when the sample size is small, but improve with larger sample size. The coverage probability of test based confidence interval using bootstrap-calibrated critical value is close to the nominal level of 95% for all sample sizes.

Table 3.7: Coverage of 95% confidence interval for the second order augmented parametric corrected score with 20% censoring rate. C (chi-square distribution), BC (bootstrap calibration) and W (Wald-type) indicate the type of confidence interval.

| size | 100 | | | 200 | | | 400 | | | 800 | | | 1600 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | C | BC | W | C | BC | W | C | BC | W | C | BC | W | C | BC | W |
| Single-covariate: Case A | | | | | | | | | | | | | | | |
| $\beta$ | 90.2 | 96.8 | 87.0 | 91.7 | 96.6 | 91.1 | 87.1 | 96.6 | 93.8 | 87.6 | 97.1 | 94.8 | 91.8 | 95.6 | 94.9 |
| Single-covariate: Case B | | | | | | | | | | | | | | | |
| $\beta$ | 82.0 | 92.9 | 76.8 | 88.8 | 92.2 | 84.3 | 87.6 | 94.3 | 85.6 | 90.0 | 94.6 | 89.2 | 91.6 | 94.8 | 91.8 |
| Double-covariate | | | | | | | | | | | | | | | |
| $\beta_1$ | 75.8 | 94.5 | 78.4 | 86.3 | 95.4 | 87.3 | 88.3 | 96.3 | 94.0 | 89.7 | 96.7 | 95.8 | 89.4 | 96.2 | 97.2 |
| $\beta_2$ | 79.2 | 93.0 | 83.9 | 86.8 | 93.8 | 92.3 | 88.1 | 92.3 | 95.6 | 86.7 | 93.5 | 96.2 | 89.4 | 95.5 | 97.6 |

Table 3.8: Coverage of 95% confidence interval for the second order augmented parametric corrected score with 60% censoring rate. C (chi-square distribution), BC (bootstrap calibration) and W (Wald-type) indicate the type of confidence interval.

| size | 100 | | | 200 | | | 400 | | | 800 | | | 1600 | | |
|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| | C | BC | W | C | BC | W | C | BC | W | C | BC | W | C | BC | W |
| | | | | | | Single-covariate: Case A | | | | | | | | | |
| $\beta$ | 89.5 | 97.6 | 86.8 | 94.2 | 96.7 | 92.0 | 91.5 | 96.8 | 94.7 | 93.6 | 96.2 | 96.3 | 98.5 | 95.5 | 97.5 |
| | | | | | | Single-covariate: Case B | | | | | | | | | |
| $\beta$ | 87.5 | 93.7 | 73.9 | 89.7 | 92.4 | 79.1 | 88.4 | 93.2 | 80.4 | 90.4 | 93.9 | 85.8 | 91.8 | 94.3 | 88.6 |
| | | | | | | Double-covariate | | | | | | | | | |
| $\beta_1$ | 92.0 | 92.6 | 80.7 | 93.8 | 96.6 | 89.6 | 92.6 | 94.9 | 94.1 | 93.2 | 96.9 | 95.8 | 92.6 | 96.2 | 98.0 |
| $\beta_2$ | 92.0 | 93.9 | 87.3 | 93.5 | 94.5 | 93.7 | 92.3 | 92.7 | 96.0 | 91.5 | 93.4 | 96.3 | 92.7 | 95.8 | 97.1 |

## 3.4.2   Application to ACTG 175 data

We apply the proposed approach to the AIDS Clinical Trial Group (ACTG) 175 study, a randomized clinical trial to evaluate four treatments in HIV-infected patients with an initial screening CD4 counts of between 200 and 500 per cubic millimeter. A total of 2,467 patients were enrolled and an almost equal number of patients were randomized into each of the four treatment groups: zidovudine alone (ZDV), zidovudine plus didanosine (ZDV + ddI), zidovudine plus zalcitabine (ZDV + ddC), and zalcitabine alone (ddC). We are interested in assessing the effect of baseline CD4 count on time to AIDS or death in antiretroviral-naive patients. Among all study patients, 1,067 had no prior antiretroviral therapy at enrollment, among which 1,036 patients had two CD4

measurements prior to the start of treatment and within 3 weeks of randomization. For this analysis, we will consider the subset of 1,036 patients. The median length of follow-up was 32 months, and 85 events were observed.

Table 3.9: Comparison of regression coefficient estimators in the ACTG 175 data

| | log(CD4) | | ZDV + ddl | | ZDV + ddC | | ddl | |
|---|---|---|---|---|---|---|---|---|
| | Est | Var | Est | Var | Est | Var | Est | Var |
| NV | -1.838 | .1183 | -.652 | .0881 | -.895 | .1006 | -.598 | .0802 |
| ConS | -2.172 | .1698 | -.659 | .0916 | -.892 | .1024 | -.604 | .0835 |
| CS | -2.177 | .1678 | -.659 | .0900 | -.892 | .1012 | -.604 | .0819 |
| ACS | -2.177 | .1422 | -.657 | .0991 | -.876 | .1028 | -.596 | .0827 |
| | | .1678 | | .0905 | | .0964 | | .0811 |

Note: For proposed estimator, the first row of variance estimator is obtained by inverting hypothese testing statistics with bootstrap critical value; the second row of values is from sandwich variance estimator. Est: Estimated coefficient; Var: Variance.

We consider a Cox regression mode with 4 covariates: the true baseline log(CD4) and three indicators for the four treatments with ZDV group as the reference. We define the baseline log(CD4) as the average of the two log(CD4) measurements. From the duplicated measurements, we estimated the variance for error and true underlying log(CD4) to be 0.033 and 0.076 respectively. Note that the variance of measurement error is estimated using two replicated measurements of baseline log(CD). Therefore there is an additional estimating function for the variance of measurement error. Table 3.9 shows the estimators based on the naive, conditional score, parametric corrected score, and the proposed augmented corrected score. In comparison, the naive approach gives an coefficient estimator of log(CD4) with substantially smaller

magnitude. All the other approaches have similar estimates for all coefficients.

## 3.5    Discussion

For proportional hazards model with covariate measurement error, several consistent methods have been proposed under the functional modeling framework, including the conditional score and the parametric corrected score. However, when the measurement error is substantial to the extent that the errors are comparable to the true covariates in variance, both methods might experience pathological behaviors and root finding failure. Recently, Huang (2011) developed a novel approach to incorporate additional estimating functions which constrain the derivatives of the parametric corrected score. That approach proves effective and eliminates finite sample pathological behaviors of parametric corrected score for the loglinear model. Motivated by Huang's (2011) approach, we conduct an investigation on the pathological behaviors of parametric corrected score and conditional score and propose an augmented parametric corrected score for the proportional hazards model by incorporating additional estimating functions to the original parametric corrected score. Results of simulation studies show the proposed approach is effective in eliminating pathological behaviors even with small sample size and substantial measurement error. The variance of proposed estimator appears to be larger than the parametric corrected score when sample size is smaller than 400, but it decreases rapidly and become smaller than the parametric corrected score as the sample size increases to 400.

In this chapter, we have only considered the situation where the distribution of the measurement error is known. With additional data available on the measurement error, the parametric distribution imposed on the measurement error may be spared. With the availability of replicated mismeasured covariates, Huang & Wang (2000) developed a nonparametric corrected score method for the proportional hazards model. In Chapter 4, we will extend the approach of incorporating additional estimating functions to the nonparametric corrected score.

## 3.6 Appendix: Asymptotic Variance of Estimating Function

For simplicity, we consider the single-covariate model given in Section 3.2.2. The second order augmented corrected score has the form

$$\boldsymbol{\varphi}(b) = n^{-1} \sum_{i=1}^{n} \int_0^{\tau} \left\{ \begin{pmatrix} W_i + b\sigma^2 \\ W_i^2 - \sigma^2 \end{pmatrix} - \frac{\sum_{j=1}^{n} Y_j(t) \exp(bW_j) \begin{pmatrix} W_j \\ (W_j - b\sigma^2)^2 - \sigma^2 \end{pmatrix}}{\sum_{j=1}^{n} Y_j(t) \exp(bW_j)} \right\} dN_i(t).$$

With functional delta method (Huang & Wang, 2000), straightforward algebra gives

$$n^{1/2} \boldsymbol{\varphi}(b) = n^{-1/2} \sum_{i=1}^{n} (B_{i1} - B_{i2}) + o_p(1)$$

where

$$B_{i1} = \int_0^\tau \left\{ \begin{pmatrix} W_i + b\sigma^2 \\ W_i^2 - \sigma^2 \end{pmatrix} - \frac{\mathcal{E}\left\{ Y(t)\exp(bW) \begin{pmatrix} W \\ (W - b\sigma^2)^2 - \sigma^2 \end{pmatrix} \right\}}{\mathcal{E}\left\{ Y(t)\exp(bW) \right\}} \right\} dN_i(t), \text{ and}$$

$$B_{i2} = \int_0^\tau \left\{ \frac{\exp(bW_i)\begin{pmatrix} W_i \\ (W_i - b\sigma^2)^2 - \sigma^2 \end{pmatrix}}{\mathcal{E}\left\{ Y(t)\exp(bW) \right\}} \right.$$
$$\left. - \frac{\mathcal{E}\left\{ Y(t)\exp(bW) \begin{pmatrix} W \\ (W - b\sigma^2)^2 - \sigma^2 \end{pmatrix} \right\} \exp(bW_i)}{\left[ \mathcal{E}\left\{ Y(t)\exp(bW) \right\} \right]^2} \right\} Y_i(t) d\mathcal{E}N(t).$$

Thus, $n^{1/2}\boldsymbol{\varphi}(b)$ is asymptotically a sum of iid random variables.

For fixed $b$, $n^{1/2}\boldsymbol{\varphi}(b)$ is asymptotically normal with a covariance matrix $\Sigma(b)$ that can be consistently estimated by

$$\hat{\Sigma}(b) = n \sum_{i=1}^n \left\{ \omega_i(b) - \bar{\omega}(b) \right\} \left\{ \omega_i(b) - \bar{\omega}(b) \right\}'$$

where $\omega_i(b) = n^{-1}(B_{i3} - B_{i4})$, $\bar{\omega}(b) = n^{-1}\sum_{i=1}^n \omega_i(b)$, and $B_{i3}$ and $B_{i4}$ are defined as

$$B_{i3} = \int_0^\tau \left\{ \begin{pmatrix} W_i + b\sigma^2 \\ W_i^2 - \sigma^2 \end{pmatrix} - \frac{\hat{\mathcal{E}}\left\{ Y(t)\exp(bW) \begin{pmatrix} W \\ (W - b\sigma^2)^2 - \sigma^2 \end{pmatrix} \right\}}{\hat{\mathcal{E}}\left\{ Y(t)\exp(bW) \right\}} \right\} dN_i(t), \text{ and}$$

$$B_{i4} = \int_0^\tau \left\{ \frac{\exp(bW_i)\begin{pmatrix} W_i \\ (W_i - b\sigma^2)^2 - \sigma^2 \end{pmatrix}}{\hat{\mathcal{E}}\left\{ Y(t)\exp(bW) \right\}} \right.$$
$$\left. - \frac{\hat{\mathcal{E}}\left\{ Y(t)\exp(bW) \begin{pmatrix} W \\ (W - b\sigma^2)^2 - \sigma^2 \end{pmatrix} \right\} \exp(bW_i)}{\left[ \hat{\mathcal{E}}\left\{ Y(t)\exp(bW) \right\} \right]^2} \right\} Y_i(t) d\hat{\mathcal{E}}N(t).$$

Asymptotic variance of other augmented corrected scores could be derived similarly.

# Chapter 4

# Augmented Nonparametric Corrected Score for Proportional Hazards Model with Covariate Measurement Error

## 4.1 Introduction

In Chapter 3, to address the issue of finite-sample pathological behaviors of the parametric corrected score, we propose an approach to incorporate additional estimating functions. In that topic, we consider the situation where the distribution of mea-

surement error is known. But in practice, it is common that the measurement error distribution is unknown but rather replication data are available for error-prone covariates. Based on replicated mismeasured covariates, Huang & Wang (2000) proposed a nonparametric corrected score for the Cox proportional hazards model under additive measurement error model and the resulting regression coefficient estimators are shown to be consistent and are asymptotically normal.

In this chapter, we conduct an investigation on the pathological behaviors of nonparametric corrected score and extend the approach developed in Chapter 3 to the nonparametric corrected score for the proportional hazards model. In Section 4.2, we briefly describe the nonparametric corrected score method and investigate its finite sample behaviors. The proposed approach of incorporating additional estimating functions is presented in Section 4.3. Simulation study results and an application to the ACTG 175 clinical trial data is summarized in Section 4.4. Further discussion is given in Section 4.5.

## 4.2  Nonparametric Corrected Score and Patholog-ical Behaviors

### 4.2.1  Nonparametric Corrected Score

In this section, we will briefly describe the nonparametric corrected score of Huang & Wang (2000) and study its pathological behaviors when the measurement error is substantial.

First note that the partial score function (3.2) can be rewritten as a function of four empirical processes

$$\boldsymbol{\xi}(\mathbf{b}) = \int_0^\tau \left[ d\widehat{\mathcal{E}}\{\mathbf{Z}N(t)\} - \frac{\widehat{\mathcal{E}}\{Y(t)\mathbf{Z}\exp(\mathbf{b}'\mathbf{Z})\}}{\widehat{\mathcal{E}}\{Y(t)\exp(\mathbf{b}'\mathbf{Z})\}} d\widehat{\mathcal{E}}\{N(t)\} \right] \tag{4.1}$$

where $\widehat{\mathcal{E}}$ represents sample empirical mean as defined in Chapter 2. With the functional representation of $\boldsymbol{\xi}(\mathbf{b})$, we can see clearly its limit:

$$\widetilde{\boldsymbol{\xi}}(\mathbf{b}) = \int_0^\tau \left[ d\mathcal{E}\{\mathbf{Z}N(t)\} - \frac{\mathcal{E}\{Y(t)\mathbf{Z}\exp(\mathbf{b}'\mathbf{Z})\}}{\mathcal{E}\{Y(t)\exp(\mathbf{b}'\mathbf{Z})\}} d\mathcal{E}\{N(t)\} \right] \tag{4.2}$$

Split covariates $\mathbf{Z} = (\mathbf{Z}_e^T, \mathbf{Z}_a^T)^T$, where $\mathbf{Z}_a$ are accurately measured and $\mathbf{Z}_e \equiv (Z_1, \dots, Z_L)^T$ are subject to an additive measurement error $\boldsymbol{\epsilon}_e$. $\boldsymbol{\epsilon}_e = (\epsilon_1, \dots, \epsilon_L)^T$ are mutually independent and independent of all other variables. Under the additive measurement error model, for each $Z_l$, $l = 1, \dots, L$, a finite number of $R_l \geq 2$

surrogates $\mathbf{W}_l^{[R_l]} \equiv \{W_{lm} : m = 1, \ldots, R_l\}$ are observed, where $W_{lm} = Z_l + \epsilon_{lm}$.

These error $\epsilon_{lm}$ are iid replicate of $\epsilon_l$. The observed data,

$$\{X_i; Y_i; W_{li}^{[R_{li}]} : l = 1 \ldots, L; \mathbf{Z}_{ai}\}, i = 1, \ldots, n,$$

consist of $n$ iid replicates of $\{X; Y; \mathbf{W}_l^{[R_l]} : l = 1 \ldots, L; \mathbf{Z}_a\}$.

Pick arbitrarily two replicates from each $\mathbf{W}_l^{[R_l]}$ where $l = 1, \ldots, L$ to form two

vectors $\mathbf{W}^{(r)} \equiv (W_1^{(r)T}, \ldots, W_L^{(r)T}, \mathbf{Z}_a^T)^T, r = 1, 2$. Since the selection of $\mathbf{W}^{(1)}$ and

$\mathbf{W}^{(2)}$ is arbitrary, $\prod_{l=1}^{L} R_l(R_l-1)$ different permutations can be formed. The following

nonparametric corrected estimating function is proposed by Huang & Wang (2000):

$$\widehat{\boldsymbol{\xi}}(\mathbf{b}) = \int_0^\tau \left[ d\widehat{\mathcal{E}}\{\mathcal{A}\mathbf{W}^{(1)}N(t)\} - \frac{\widehat{\mathcal{E}}\{\mathcal{A}Y(t)\mathbf{W}^{(1)}\exp(\mathbf{b}'\mathbf{W}^{(2)})\}}{\widehat{\mathcal{E}}\{\mathcal{A}Y(t)\exp(\mathbf{b}'\mathbf{W}^{(1)})\}} d\widehat{\mathcal{E}}\{N(t)\} \right] \qquad (4.3)$$

where $\mathcal{A}$ denotes the operator averaging over all the different permutations of $\mathbf{W}^{(1)}$

and $\mathbf{W}^{(2)}$. Estimating function (4.3) converges to the same limit as reference (4.1)

and the resulting regression coefficient estimators are shown to be consistent and are

asymptotic normal.

## 4.2.2    Pathological Behaviors

Consider a single covariate model. In the absence of measurement error, the partial

score function is monotonically decreasing and has a unique zero-crossing. But with

measurement error, we observe five different patterns for the nonparametric corrected score estimating functions and typical plots from simulated data sets are shown in Figure 4.1.

Pattern (a) in Figure 4.1 has no zero-crossing and represents a case of root-finding failure. Pattern (b) is monotonically decreasing with one zero-crossing and is considered as the ideal case. A correct root can be easily identified for this pattern. Though pattern (c) also has only one zero-crossing, it is an increasing root and is also considered as root finding failure. Pattern (d) has two zero-crossings that are far away from each other. With a properly chosen numerical algorithm and initial value, the correct root could be identified. Pattern (e) has two zeros-crossings that are both close to the true parameter.

We conducted a simulation study to investigate the prevalence of root-finding failure in the presence of substantial measurement error. We use the same simulation set up, definition of root-finding failure and modified Newton-Raphson algorithm as described in Chapter 3. For each covariates subject to measurement error, two replicated surrogates were generated. As shown in Table 4.1, in the single-covariate models, the failure rate is about 11-14% for normal covariate model and 27-29% for modified chi-square covariate model when the sample size is 100. The root-finding failure rate decreases as the sample size increases. But even when the sample size increases to 800, the failure rate is still 11-16% for modified chi-square covariate model.

Figure 4.1: Observed root patterns of the nonparametric corrected score estimating function. The true $\beta$ is -1.

Table 4.1: Prevalence (%) of root-finding failure for the nonparametric corrected score with $E(X) = 0$, $\mathrm{Var}(X) = 1$, $\beta = -1$, and $\varepsilon \sim \mathrm{Normal}(0,1)$.

| Censoring Rate | Distribution of X | Size | | | |
|---|---|---|---|---|---|
| | | 100 | 200 | 400 | 800 |
| 20% | Normal | 10.7 | 7.4 | 2.6 | 1.1 |
| | Modified Chi-square | 27.0 | 20.9 | 14.0 | 11.2 |
| | Uniform | 12.5 | 7.9 | 3.8 | 1.3 |
| 40% | Normal | 10.9 | 6.8 | 2.7 | 1.1 |
| | Modified Chi-square | 29.2 | 23.5 | 17.3 | 14.5 |
| | Uniform | 12.3 | 8.4 | 3.5 | 1.3 |
| 60% | Normal | 13.8 | 7.5 | 2.6 | 1.0 |
| | Modified Chi-square | 28.4 | 23.9 | 22.0 | 15.7 |
| | Uniform | 13.9 | 10.0 | 4.9 | 2.2 |

## 4.3 Improving Nonparametric Corrected Score

Let $\mathbf{Z} = (Z_1, Z_2, \ldots, Z_p)^T$, where $Z_l$ is the $l$th element of vector $\mathbf{Z}$. Consider the following class of estimating equations:

$$\int_0^\tau \left[ d\widehat{\mathcal{E}}\{Z_1^{k_1} \cdots Z_p^{k_p} N(t)\} - \frac{\widehat{\mathcal{E}}\{Y(t) Z_1^{k_1} \cdots Z_p^{k_p} \exp(\mathbf{b}'\mathbf{Z})\}}{\widehat{\mathcal{E}}\{Y(t)\exp(\mathbf{b}'\mathbf{Z})\}} d\widehat{\mathcal{E}}\{N(t)\} \right], \text{where} \sum_{l=1}^p k_l = k. \tag{4.4}$$

Equations (4.4) present a class of estimating functions for parameter $\mathbf{b}$ in the absence of measurement error. When $k = 1$, (4.4) reduces to the usual partial score function. Using (4.4) as a class of reference functions, their corresponding corrected version could serve the purpose of providing additional estimating functions to the nonparametric corrected score function.

Let $\mathbf{W}^{(r)} = (W_1^{(r)}, \ldots, W_p^r)^T, r = 1, 2$. Mimicking the argument in Huang & Wang (2000), we obtain the following additional estimating functions with mismeasured covariates sharing the same limit as reference (4.4):

$$\int_0^\tau \left[ d\widehat{\mathcal{E}}\{\mathcal{A}(W_1^{(1)})^{k_1} \cdots (W_p^{(1)})^{k_p} N(t)\} \right.$$
$$\left. - \frac{\widehat{\mathcal{E}}\{\mathcal{A}Y(t)(W_1^{(1)})^{k_1} \cdots (W_p^{(1)})^{k_p} \exp(\mathbf{b}'\mathbf{W}^{(2)})\}}{\widehat{\mathcal{E}}\{\mathcal{A}Y(t) \exp(\mathbf{b}'\mathbf{W}^{(1)})\}} d\widehat{\mathcal{E}}\{N(t)\} \right], \text{where } \sum_{l=1}^p k_l = k.$$

Quadratic inference function method will then be used to combine nonparametric corrected score and these additional estimating functions. We will name this method the *augmented nonparametric corrected score*. Specifically, when $\sum_{l=1}^p k_l = k$, this method will be called the $k$-th order augmented nonparametric corrected score. Figure 4.2 illustrates the quadratic inference functions for the five typical plots in Figure 4.1.

Figure 4.2: Quadratic inference functions for the data sets in Figure 4.1.

## 4.4   Numerical Studies

### 4.4.1   Simulations

The simulation set up is same as in that of Chapter 3. For each covariate sub-ject to measurement error, two replicated surrogates were generated. For reference and comparison, the ideal, naive, regression calibration, and nonparametric corrected score estimators are also presented. The nonparametric corrected score is similarly re-defined and solved using a modified Newton-Raphson algorithm as in Chapter 3. For naive estimator, the average of two replicated measurements will be used in place of the true covariate in the usual partial score function. For regression calibration, the best linear approximation given in Carroll et al. (2006, chap. 4) will be used.

Table 4.2 and 4.3 summarize the simulation results on the estimators in the single-covariate models. Quantile-quantile plots are shown in Figure 4.3 and 4.4. Under each scenario, the mean bias, and standard deviation were calculated. In addition, median bias and a robustified standard deviation (IQR divided by 1.349) are also reported for nonparametric corrected score and the proposed approach. As expected, naive estimator incurs substantial bias under all scenarios. The regression calibration estimator shows small bias for normal covariate. But it shows substantial bias when the covariate is modified chi-square. The nonparametric corrected score estimator is consistent and shows little bias. For the single-covariate models with 60% censoring rate, the finding failure rate is about 11% for normal covariate when the sample size

Table 4.2: Simulation summary statistics for the single-covariate models with 20% censoring rate: Ideal, naive(NV), regression calibration (RC), re-defined nonparametric corrected score (NPC), and second order augmented corrected score (ACS:2).

| size | Ideal B | SD | NV B | SD | RC B | SD | NPC F | B | | SD | | ACS: 2 B | | SD | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | | | |
| | | | | | | | Case A | | | | | | | | |
| 100 | -11 | 152 | 421 | 119 | 111 | 212 | 10.7 | -349 | -97 | 842 | 466 | -633 | -105 | 2121 | 442 |
| 200 | -5 | 106 | 429 | 79 | 134 | 139 | 7.4 | -214 | -64 | 565 | 330 | -385 | -67 | 1510 | 322 |
| 400 | -4 | 72 | 435 | 55 | 148 | 95 | 2.6 | -92 | -28 | 332 | 191 | -167 | -20 | 1025 | 199 |
| 800 | -2 | 51 | 437 | 37 | 155 | 61 | 1.1 | -33 | -9 | 151 | 128 | -35 | -7 | 171 | 126 |
| 1600 | -2 | 37 | 436 | 27 | 152 | 44 | 0.1 | -21 | -9 | 101 | 96 | -22 | -14 | 103 | 97 |
| | | | | | | | Case B | | | | | | | | |
| 100 | -28 | 202 | 492 | 108 | 212 | 198 | 27.0 | -299 | -41 | 882 | 484 | -804 | -90 | 2373 | 543 |
| 200 | -10 | 136 | 496 | 75 | 236 | 127 | 20.9 | -194 | -46 | 626 | 368 | -333 | -51 | 1009 | 350 |
| 400 | -9 | 92 | 500 | 52 | 244 | 88 | 14.0 | -175 | -24 | 583 | 304 | -222 | -17 | 1102 | 233 |
| 800 | -6 | 67 | 498 | 37 | 245 | 62 | 11.2 | -115 | -35 | 346 | 244 | -126 | -24 | 646 | 161 |
| 1600 | -2 | 46 | 500 | 26 | 250 | 43 | 5.2 | -68 | -14 | 241 | 161 | -53 | -12 | 241 | 115 |

Note: F: root-finding failure (%); B: bias ($\times 10^3$); SD: standard deviation($\times 10^3$). For NPC and ACS: 2, The first column under B is mean bias and the second column is median bias; The first column under SD is the usual standard deviation and the second column is robustified standard deviation defined as IQR/1.349.

Table 4.3: Simulation summary statistics for the single-covariate models with 60% censoring rate: Ideal, naive(NV), regression calibration (RC), re-defined nonparametric corrected score (NPC), and second order augmented corrected score (ACS:2).

| size | Ideal B | SD | NV B | SD | RC B | SD | NPC F | B | | SD | | ACS:2 B | | SD | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Case A | | | | | | | |
| 100 | -24 | 206 | 382 | 160 | 52 | 273 | 13.8 | -338 | -98 | 888 | 494 | -463 | -112 | 1799 | 422 |
| 200 | -15 | 142 | 384 | 105 | 66 | 180 | 7.5 | -205 | -81 | 534 | 347 | -271 | -78 | 1025 | 304 |
| 400 | -8 | 97 | 395 | 72 | 89 | 120 | 2.6 | -79 | -23 | 303 | 202 | -85 | -27 | 413 | 186 |
| 800 | -2 | 65 | 400 | 51 | 99 | 83 | 1.0 | -30 | -5 | 173 | 139 | -24 | -2 | 173 | 128 |
| 1600 | 1 | 46 | 398 | 36 | 95 | 57 | 0 | -21 | -13 | 109 | 94 | -14 | -9 | 106 | 88 |
| | | | | | | | | Case B | | | | | | | |
| 100 | -66 | 351 | 545 | 149 | 294 | 255 | 28.4 | -215 | 55 | 922 | 523 | -410 | -22 | 1488 | 598 |
| 200 | -24 | 239 | 557 | 102 | 329 | 161 | 23.9 | -134 | 34 | 673 | 401 | -298 | -39 | 1106 | 474 |
| 400 | -16 | 159 | 562 | 70 | 338 | 111 | 22.0 | -155 | -11 | 606 | 377 | -355 | -38 | 1169 | 405 |
| 800 | -9 | 112 | 562 | 52 | 342 | 81 | 15.7 | -108 | 9 | 426 | 309 | -181 | -16 | 720 | 254 |
| 1600 | 0 | 78 | 564 | 37 | 345 | 57 | 12.2 | -81 | -2 | 305 | 236 | -107 | -10 | 500 | 179 |

Note: Same as in that of Table 4.2.

Figure 4.3: Quantile-quantile plots for $\beta$ in the single-covariate models with 20% censoring rate, where $\beta = -1$. Red, yellow, green, blue, and black correspond to sample sizes 100, 200, 400, 800, and 1,600.

Figure 4.4: Quantile-quantile plots for $\beta$ in the single-covariate models with 60% censoring rate, where $\beta = -1$. Red, yellow, green, blue, and black correspond to sample sizes 100, 200, 400, 800, and 1,600.

is 100. The failure rate deceases to 1% as the sample size increases to 800. The failure rate is much higher for modified chi-square covariates. Even with a sample size of 800, the failure rate is around 11%. Similar root-finding failure rates are observed for models with 20% censoring rate. The proposed augmented corrected score is left skewed with substantial mean bias and large standard deviation. But the median bias is close to zero and the robustified standard deviation is comparable to or smaller than that of nonparametric corrected score. The quantile-quantile plots show that the proposed estimator suffers from substantial skewness.

Table 4.4 and 4.5 report the simulation results with the double-covariate models. The relative performance of most estimators are as expected and follow a similar pattern to that in the single-covariate models.

The coverage of the 95% confidence intervals are summarized in Table 4.6 and 4.7. The test based confidence interval using chi-square distribution critical value and Wald-type one appear to have poor coverage probability, whereas the one using bootstrap-calibrated critical value has coverage probability close to nominal level of 95%.

## 4.4.2 Application to ACTG 175 data

We apply the proposed approach to the AIDS Clinical Trial Group (ACTG) 175 study. Among 1,067 antiretroviral-naive patients at enrollment, 1,036 patients had two CD4

Table 4.4: Simulation summary statistics for the double-covariate models with 20% censoring rate: Ideal, naive(NV), regression calibration (RC), re-defined nonparametric corrected score (NPC), and second order augmented corrected score (ACS:2).

| size | | Ideal | | NV | | RC | | NPC | | | | | ACS: 2 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | B | SD | B | SD | B | SD | F | B | | SD | | B | | SD | |
| 100 | $\beta_1$ | -23 | 175 | 463 | 122 | 61 | 271 | 30.5 | -47 | 69 | 518 | 383 | -500 | -106 | 1856 | 506 |
| | $\beta_2$ | 23 | 171 | -289 | 154 | -88 | 211 | | 100 | -41 | 512 | 335 | 468 | 108 | 1536 | 427 |
| 200 | $\beta_1$ | -13 | 114 | 477 | 85 | 100 | 174 | 21.0 | -24 | 38 | 386 | 278 | -300 | -70 | 1136 | 349 |
| | $\beta_2$ | 16 | 117 | -307 | 105 | -122 | 140 | | 57 | -16 | 373 | 235 | 217 | 57 | 772 | 290 |
| 400 | $\beta_1$ | -5 | 82 | 485 | 58 | 137 | 115 | 9.4 | -9 | 10 | 217 | 196 | -155 | -27 | 751 | 243 |
| | $\beta_2$ | 4 | 80 | -313 | 74 | -137 | 94 | | 23 | -8 | 188 | 167 | 124 | 34 | 612 | 190 |
| 800 | $\beta_1$ | -2 | 55 | 482 | 39 | 135 | 78 | 4.5 | -28 | -14 | 166 | 150 | -86 | -25 | 481 | 157 |
| | $\beta_2$ | 1 | 56 | -313 | 50 | -138 | 64 | | 28 | 15 | 134 | 119 | 64 | 27 | 344 | 122 |
| 1600 | $\beta_1$ | -2 | 40 | 485 | 28 | 140 | 54 | 0.6 | -25 | -11 | 119 | 106 | -16 | -16 | 201 | 111 |
| | $\beta_2$ | 1 | 39 | -313 | 37 | -141 | 46 | | 20 | 13 | 93 | 92 | 17 | 17 | 140 | 93 |

Note: Same as in that of Table 4.2.

Table 4.5: Simulation summary statistics for the double-covariate models with 60% censoring rate: Ideal, naive(NV), regression calibration (RC), re-defined nonparametric corrected score (NPC), and second order augmented corrected score (ACS:2).

| size | | Ideal | | NV | | RC | | NPC | | | | | ACS: 2 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | B | SD | B | SD | B | SD | F | B | | SD | | B | | SD | |
| 100 | $\beta_1$ | -32 | 242 | 429 | 172 | 2 | 347 | 27.8 | -78 | 28 | 556 | 373 | -510 | -74 | 2334 | 507 |
| | $\beta_2$ | 33 | 245 | -252 | 204 | -38 | 270 | | 113 | -9 | 551 | 343 | 446 | 79 | 2178 | 438 |
| 200 | $\beta_1$ | -20 | 160 | 437 | 119 | 40 | 230 | 17.3 | -73 | 13 | 452 | 291 | -326 | -76 | 1372 | 374 |
| | $\beta_2$ | 21 | 161 | -260 | 144 | -62 | 186 | | 91 | 9 | 432 | 265 | 233 | 58 | 1028 | 305 |
| 400 | $\beta_1$ | -6 | 109 | 451 | 81 | 80 | 152 | 7.7 | -35 | 9 | 304 | 230 | -138 | -36 | 847 | 237 |
| | $\beta_2$ | 6 | 111 | -268 | 103 | -80 | 127 | | 43 | 1 | 279 | 188 | 102 | 34 | 505 | 208 |
| 800 | $\beta_1$ | -1 | 72 | 452 | 54 | 81 | 101 | 2.8 | -37 | -16 | 183 | 167 | -66 | -24 | 306 | 152 |
| | $\beta_2$ | 2 | 76 | -273 | 68 | -88 | 82 | | 27 | 13 | 143 | 131 | 44 | 19 | 206 | 125 |
| 1600 | $\beta_1$ | -2 | 53 | 449 | 38 | 83 | 68 | 0.6 | -32 | -20 | 128 | 116 | -45 | -18 | 219 | 104 |
| | $\beta_2$ | 0 | 53 | -268 | 49 | -84 | 57 | | 24 | 15 | 99 | 89 | 34 | 20 | 146 | 88 |

Note: Same as in that of Table 4.2.

Table 4.6: Coverage of 95% confidence interval for the second order augmented non-parametric corrected score with 20% censoring rate. C (chi-square distribution), BC (bootstrap calibration) and W (Wald-type) indicate the type of confidence interval.

| size | 100 | | | 200 | | | 400 | | | 800 | | | 1600 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | C | BC | W | C | BC | W | C | BC | W | C | BC | W | C | BC | W |
| Single-covariate: Case A | | | | | | | | | | | | | | | |
| $\beta$ | 79.6 | 93.9 | 86.7 | 81.5 | 96.2 | 90.5 | 85.3 | 96.2 | 90.6 | 86.7 | 95.5 | 91.3 | 88.4 | 95.2 | 90.6 |
| Single-covariate: Case B | | | | | | | | | | | | | | | |
| $\beta$ | 82.7 | 96.4 | 84.6 | 83.3 | 96.9 | 87.9 | 85.3 | 96.8 | 90.7 | 85.1 | 96.1 | 91.0 | 86.2 | 95.8 | 91.6 |
| Double-covariate | | | | | | | | | | | | | | | |
| $\beta_1$ | 72.1 | 95.1 | 75.2 | 78.6 | 96.8 | 88.3 | 82.7 | 95.4 | 88.4 | 85.3 | 96.2 | 90.3 | 84.8 | 96.1 | 88.5 |
| $\beta_2$ | 73.7 | 93.0 | 77.9 | 79.4 | 95.2 | 87.5 | 84.0 | 96.9 | 88.5 | 85.3 | 95.7 | 90.5 | 86.0 | 95.4 | 89.6 |

Table 4.7: Coverage of 95% confidence interval for the second order augmented non-parametric corrected score with 60% censoring rate. C (chi-square distribution), BC (bootstrap calibration) and W (Wald-type) indicate the type of confidence interval.

| size | 100 | | | 200 | | | 400 | | | 800 | | | 1600 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | C | BC | W | C | BC | W | C | BC | W | C | BC | W | C | BC | W |
| Single-covariate: Case A | | | | | | | | | | | | | | | |
| $\beta$ | 84.0 | 94.0 | 89.7 | 84.6 | 95.8 | 93.5 | 89.6 | 95.4 | 94.1 | 88.9 | 95.0 | 93.9 | 89.3 | 95.0 | 92.0 |
| Single-covariate: Case B | | | | | | | | | | | | | | | |
| $\beta$ | 88.2 | 96.2 | 84.6 | 90.9 | 98.3 | 90.5 | 89.2 | 96.6 | 90.2 | 88.8 | 94.7 | 91.0 | 89.4 | 94.7 | 91.0 |
| Double-covariate | | | | | | | | | | | | | | | |
| $\beta_1$ | 81.2 | 95.6 | 84.6 | 81.0 | 97.2 | 89.5 | 86.8 | 93.6 | 90.3 | 88.5 | 96.3 | 91.8 | 89.0 | 96.1 | 91.7 |
| $\beta_2$ | 83.0 | 94.7 | 86.6 | 84.8 | 95.8 | 89.2 | 86.4 | 94.5 | 91.8 | 89.1 | 95.7 | 92.0 | 87.6 | 95.4 | 89.8 |

measurements prior to the start of treatment and within 3 weeks of randomization. For this analysis, we will consider the subset of 1,036 patients and use the two CD4 counts as replicated measurements of true baseline CD4 count.

We consider a Cox regression mode with 4 covariates: the true baseline log(CD4) and three indicators for the four treatments with ZDV group as the reference. Table 4.8 shows the estimators based on the naive, regression calibration, nonparametric-correction and the proposed approaches. In comparison, the naive approach gives an coefficient estimator of log(CD4) with substantially smaller magnitude. All the other approaches have similar estimates for all coefficients.

Table 4.8: Comparison of regression coefficient estimators in the ACTG 175 data

|          | log(CD4) | | ZDV + ddl | | ZDV + ddC | | ddl | |
|----------|----------|-------|-----------|-------|-----------|-------|-------|-------|
|          | Est      | Var   | Est       | Var   | Est       | Var   | Est   | Var   |
| NV       | -1.838   | .1183 | -.652     | .0881 | -.895     | .1006 | -.598 | .0802 |
| RC       | -2.235   | .1756 | -.649     | .0882 | -.890     | .1008 | -.603 | .0803 |
| NPC      | -2.200   | .1795 | -.653     | .0906 | -.870     | .1010 | -.603 | .0828 |
| Proposed | -2.186   | .2135 | -.652     | .0970 | -.877     | .1168 | -.604 | .0900 |
|          |          | .1820 |           | .0901 |           | .0999 |       | .0826 |

Note: For proposed estimator, the first row of variance estimator is obtained by inverting hypothese testing statistics with bootstrap critical value; the second row of values is from sandwich variance estimator. Est: Estimated coefficient; Var: Variance.

## 4.5 Discussion

In this chapter, we have extended the approach of incorporating additional estimating functions developed in Chapter 3 into the nonparametric corrected score for Cox proportional hazards model. With replicated measurements of error-prone covariates available, this approach does not require any assumptions in addition to the Cox proportional hazards model and the additive measurement error model. This is an appealing feature because many assumptions can be difficult to verify in practice. Simulation studies show that the proposed approach produces small bias and is promising in eliminating root finding failure. The proposed estimator suffers from left skewness, especially in the case of modified chi-square covariate. In general, the variance of proposed estimator appears to be larger than the nonparametric corrected score when sample size is smaller than 400, but it becomes smaller than the nonparametric corrected score for sample size greater than 400.

# Chapter 5

# Summary and Future Work

## 5.1 Summary

In this dissertation research, we study three measurement error problems in the analysis of survival data.

We first study the analysis of recurrent events data under the accelerated failure time model in the presence of covariate measurement error. With replicated mismeasured covariates available, we propose a estimation procedure based on a novel identity. The proposed estimation procedure requires no distributional assumptions on either the true covariates or the error except for the boundedness of the latter. The resulting regression coefficient estimators are shown to be consistent and asymptotically normal. Simulation studies show the proposed procedure performance well

with practical sample size and moderate measurement error. We apply the proposed method to NPC clinical trial to illustrate its practical utility.

We then study the Cox proportional hazards model assuming the distribution of the measurement error is known. Both existing consistent methods, parametric corrected score (Nakamura, 1992) and conditional score (Tsiatis & Davidian, 2001), suffer from severe finite-sample pathological behaviors when the measurement error is substantial. We study the finite sample pathological behaviors of these two estimators and propose an augmented parametric corrected score by incorporating additional estimating functions to the original parametric corrected score. The estimation and inference are then accomplished by means of quadratic inference function. Simulation studies show the proposed approach is effective in eliminating finite-sample pathological behaviors even with small sample size and substantial measurement error. An application to ACTG 175 study is presented.

Furthermore, we consider the Cox proportional hazards model when the error distribution is completely unspecified but replicated mismeasured covariates are available. We applied the technique in the second topic to the nonparametric corrected score and the simulation study result shows that the proposed estimation procedure is promising in resolving those pathological behaviors.

For all three measurement error problems, the naive estimators are biased, which is obvious and well expected. Prediction and hypothesis testing are, however, different stories. In general, the presence of measurement error has no effect on prediction

problem. The surrogate $\mathbf{W}$ is error-free as a measurement of itself. So if a model is built on the mismeasured covariate, predicting the response from the mismeasured covariate does not involve any bias. The effect of measurement error on hypothesis testing is much more complicated. See the monograph of Carroll et al. (2006) for a detailed discussion on the hypothesis testing problem. In general, the naive test of no effects due to $\mathbf{Z_e}$ is still valid. But the naive test of no effects due to $\mathbf{Z_a}$ is not valid except for under some restrictive assumptions.

## 5.2   Future Work

In this subsection we discuss some future research topics and possible extensions of this dissertation work.

In the first topic, the proposed method requires replicated measurements on error-prone covariates, which might not always be available in practice. For example, patients in Se group of the NPC trial did not have replicated baseline plasma Se measurements. In some other applications, due to financial constrains or various other reasons, the replication data may only be available for a subset of all subjects. This limits the applicability of the proposed method. One direction of future research is to relax this requirement and to develop an estimating function that can utilize the whole data set.

The method of incorporating additional estimating functions in the last two topics

are motivated by Huang's (2011) investigation of loglinear model. But it turns out that we have a weaker claim than in Huang's (2011) paper. In the case of loglinear model, the additional estimating functions effectively impose constrains on the derivative of corrected score. But for Cox regression model, we do not have such nice property because of the form of partial score function. Another issue is the efficiency. The proposed estimator should be asymptotically more efficient than the corrected score estimator since it involves more estimating functions. However, this research shows that having additional estimating function might not always provide better efficiency when sample size is small or moderate. How to improve the estimation efficiency is a tough problem to solve but an interesting topic for future research.

In this dissertation research, we adopt the classical measurement error model. In many applications, the measurement error models may be much more complex. For example, in nutritional studies, the observed surrogate may be a linear function of the true underlying covariate and several other covariates (Kipnis et al., 1999, 2001, 2003). One potential future research topic is to extend the proposed approaches to accommodate different measurement error models.

# Bibliography

ANDERSEN, P. K. & GILL, R. D. (1982). Cox's regression model for counting processes: a large sample study. *The Annals of Statistics* **10**, 1100–1120.

BERKSON, J. (1950). Are there two regressions? *Journal of the American Statistical Association* **45**, 164–180.

CARROLL, R. J., RUPPERT, D., STEFANSKI, L. A. & CRAINICEANU, C. (2006). *Measurement Error in Nonlinear Models: A Modern Perspective*. London: Chapman & Hall.

CLARK, L. C., COMBS, G. F., TURNBULL, B. W., SLATE, E. H., CHALKER, D. K., CHOW, J., DAVIS, L. S., GLOVER, R. A., GRAHAM, G. F., GROSS, E. G., KRONGRAD, A., LESHER, J. L., PARK, H. K., SANDERS, B. B., SMITH, C. L. & TAYLOR, J. R. (1996). Effects of selenium supplementation for cancer prevention in patients with carcinoma of the skin: A randomized controlled tria. *Journal of the American Medical Association* **276**, 1957–1963.

CLAYTON, D. G. (1991). Models for the analysis of cohort and case-control studies

with inaccurately measured exposures. In *Statistical Models for Longitudinal Studies of Health*, J. H. Dwyer, M. Feinleib & P. P. Lipsert, eds. New York: Oxford University Press.

Cox, D. R. (1972). Regression models and life-tables (with discussion). *Journal of the Royal Statistical Society, Series B* **34**, 187–220.

Duffield-Lillico, A. J., Slate, E. H., Reid, M. E., Turnbull, B. W., Wilkins, P. A., Combs, G. F., Park, H. K., Gross, E. G., Graham, G. F., Stratton, M. S., Marshall, J. R. & Clark, L. C. (2003). Selenium supplementation and secondary prevention of nonmelanoma skin cancer in a randomized trial. *Journal of the National Cancer Institute* **95**, 1477–1481.

Efron, B. & Tibshirani, R. J. (1993). *An introduction to the bootstrap.* New York: Chapman & Hall.

Fuller, W. A. (1987). *Measurement Error Models.* New York: Wiley.

Hammer, S. M., Katzenstein, D. A., Hughes, M. D., Gundacker, H., Schooley, R. T., Haubrich, R. H., Henry, W. K., Lederman, M. M., Phair, J. P., Niu, M., Hirsch, M. S. & Merigan, T. (1996). A trial comparing nucleoside monotherapy with combination therapy in hiv-infected adults with cd4 cell counts from 200 to 500 per cubic millimeter. *The Newe England Journal of Medicine* **335**, 1081–1090.

Hu, C. & Lin, D. Y. (2004). Semiparametric failure time regression with replicates

of mismeasured covariates. *Journal of the American Statistical Association* **99**, 105–118.

HUANG, Y. J. (2002). Calibration regression of censored lifetime medical cost. *Journal of the American Statistical Association* **97**, 318–327.

HUANG, Y. J. (2011). Trend-constrained corrected score for errors-in-variables. In preparation.

HUANG, Y. J. & WANG, C. Y. (1999). Nonparametric correction to errors in covariates. Tech. rep., Fred Hutchinson Cancer Research Center, Seattle.

HUANG, Y. J. & WANG, C. Y. (2000). Cox regression with accurate covariates unascertainable: A nonparametric-correction approach. *Journal of the American Statistical Association* **95**, 1209–1219.

HUANG, Y. J. & WANG, C. Y. (2001). Consistent functional methods for logistic regression with errors in covariates. *Journal of the American Statistical Association* **96**, 1469–1482.

HUANG, Y. J. & WANG, C. Y. (2006). Errors-in-covariates effect on estimating functions: Additivity in limit and non-parametric correction. *Statistica Sinica* **16**, 861–881.

JIANG, W., TURNBULL, B. W. & CLARK, L. C. (1999). Semiparametric regression models for repeated events with random effects and measurement error. *Journal of the American Statistical Association* **94**, 111–124.

KIPNIS, V., CARROLL, R. J., FREEDMAN, L. S. & LI, L. (1999). A new dietary measurement error model and its application to the estimation of relative risk: Application to four validation studies. *American Journal of Epidemiology* **150**, 642–651.

KIPNIS, V., MIDTHUNE, D., FREEDMAN, L. S., BINGHAM, S., DAY, N. E., RIBOLI, E. & CARROLL, R. J. (2003). Bias in dietary-report instruments and its implications for nutritional epidemiology. *Public Health Nutrition* **5**, 915–923.

KIPNIS, V., MIDTHUNE, D., FREEDMAN, L. S., BINGHAM, S., SCHATZKIN, A., SUBAR, A. & CARROLL, R. J. (2001). Empirical evidence of correlated biases in dietary assessment instruments and its implications. *American Journal of Epidemiology* **153**, 394–403.

KONG, F. H. & GU, M. (1999). Consistent estimation in cox proportional hazards model with covariate measurement erros. *Statistica Sinica* **9**, 953–969.

KOSOROK, M. R. (2008). *Introduction to Empirical Processes and Semiparametric Inference.* New York: Springer.

LAWLESS, J. F. & NADEAU, C. (1995). Some simple and robust methods for the analysis of recurrent events. *Technometrics* **37**, 158–168.

LAWLESS, J. F., NADEAU, C. & COOK, R. J. (1997). Analysis of mean and rate functions for recurrent events. In *Proceedings of the First Seattle Symposium in*

*Biostatistics: Survival Analysis*, D. Y. Lin & T. R. Flemming, eds. New York: Springer-Verlag.

LIN, D. Y. & GEYER, C. J. (1992). Computational methods for semiparametric linear regression with censored data. *Journal of Computational and Graphical Statistics* **1**, 77–90.

LIN, D. Y., WEI, L. J. & YING, Z. L. (1998). Accelerated failure time model for counting processes. *Biometrika* **85**, 341–350.

LINDSAY, B. G. & QU, A. (2003). Inference functions and quadratic score tests. *Statistical Science* **18**, 394–410.

MA, Y. & TSIATIS, A. A. (2006). On closed form semiparametric estimators for measurement error models. *Statistica Sinica* **16**, 183–193.

NAKAMURA, T. (1990). Corrected score function for errors in variables models: Methodology and application to generalized linear models. *Biometrika* **77**, 127–137.

NAKAMURA, T. (1992). Proportional hazards models with covariates subject to measurement error. *Biometrics* **48**, 829–838.

NELDER, J. A. & MEAD, R. (1965). A simplex algorithm for function minimization. *Computer Journal* **7**, 308–313.

PEPE, M. S. & CAI, J. (1993). Some graphical displays and marginal regression

analyses for recurrent failure times and time dependent covariates. *Journal of the American Statistical Association* **88**, 811–820.

PRENTICE, R. L. (1982). Covariate measurement errors and parameter estimation in a failure time regression model. *Biometrika* **69**, 331–342.

QIN, J. & LAWLESS, J. (1994). Empirical likelihood and general estimating equations. *The Annals of Statistics* **22**, 300–325.

SONG, X. & HUANG, Y. J. (2005). On corrected score approach for proportional hazards model with covariate measurement error. *Biometrics* **61**, 702–714.

STEFANSKI, L. A. (1989). Unbiased estimation of a nonlinear function of a noral mean with application to measurement error models. *Communications in Statistics-Theory and Methods* **18**, 4335–4358.

STEFANSKI, L. A. & CARROLL, R. J. (1987). Conditional scores and optimal scores for generalized linear measurement-error models. *Biometrika* **74**, 703–716.

TSIATIS, A. A. & DAVIDIAN, M. (2001). A semiparametric estimator for the proportional hazards model with longitudinal covariates measured with error. *Biometrika* **88**, 447–458.

TSIATIS, A. A. & MA, Y. (2004). Locally efficient semiparametric estimators for functional measurement error models. *Biometrika* **91**, 835–848.

WANG, C. Y., HSU, L., FENG, Z. D. & PRENTICE, R. L. (1997). Regression calibration in failure time regression. *Biometrics* **53**, 131–145.

WEI, L. J., LIN, D. Y. & WEISSFELD, L. (1989). Regression analysis of multivariate incomplete failure time data by modeling marginal distributions. *Journal of the American Statistical Association* **84**, 1065–1073.

ZHOU, H. & PEPE, M. S. (1995). Auxiliary covariate data in failure time regression. *Biometrika* **82**, 139–149.