

Summer 2019

## **Extension of Risk-Based Measure of Time-Varying Prognostic Discrimination for Survival Models**

Shujie Chen

Follow this and additional works at: <https://scholarcommons.sc.edu/etd>



Part of the [Biostatistics Commons](#)

---

### **Recommended Citation**

Chen, S.(2019). *Extension of Risk-Based Measure of Time-Varying Prognostic Discrimination for Survival Models*. (Master's thesis). Retrieved from <https://scholarcommons.sc.edu/etd/5431>

This Open Access Thesis is brought to you by Scholar Commons. It has been accepted for inclusion in Theses and Dissertations by an authorized administrator of Scholar Commons. For more information, please contact [dillarda@mailbox.sc.edu](mailto:dillarda@mailbox.sc.edu).

EXTENSION OF RISK-BASED MEASURE OF TIME-VARYING PROGNOSTIC  
DISCRIMINATION FOR SURVIVAL MODELS

by

Shujie Chen

Bachelor of Science  
Nanjing Medical University, 2013

---

Submitted in Partial Fulfillment of the Requirements  
for the Degree of Master of Science in Public Health in  
Biostatistics

The Norman J. Arnold School of Public Health  
University of South Carolina

2019

Accepted by:

Jiajia Zhang, Director of Thesis

Alexander McLain, Reader

Yuan Wang, Reader

Cheryl L. Addy, Vice Provost and Dean of the Graduate School

© Copyright by Shujie Chen, 2019  
All Rights Reserved.

## ACKNOWLEDGMENTS

Many people have helped me a lot during my MSPH program. First, I would like to thank my advisor, Dr. Jiajia Zhang, for her constant encouragement, patient guidance, and instructive advice. Her guidance on my thesis has allowed me to investigate the performance of hazard discrimination in different formats including the Cox proportional hazards model and time dependent proportional hazards model under different censoring. Without her consistent and illuminating guidance, this thesis could not have reached its current form.

Second, I would like to extend my heartfelt gratitude to Dr. Alexander McLain and Dr. Yuan Wang for their insightful comments and invaluable suggestions. I am greatly indebted to both of you for the considerable amount of time, experience, and knowledge you have provided during my research.

I would like to express my sincere gratitude to the faculty and staff of the Department of Epidemiology and Biostatistics for direct and indirect help to me. Thank you for sharing knowledge, manner and passion in Biostatistics.

Finally, my thanks would go to my parents and my sister for their generous love and support. I could not have done it without you.

## ABSTRACT

The Cox proportional hazards (PH) model and time dependent PH model are the most popular survival models in survival analysis. The hazard discrimination summary  $HDS(t)$  proposed by Liang and Heagerty [2017] is used to evaluate the mean hazard difference between cases and controls at time  $t$ . Liang and Heagerty [2017] evaluated the discrimination performance under the PH model and time dependent PH model with right censoring.

In this thesis, first, we further investigate their method via comprehensive simulations including 1) We extend the simulation in Liang and Heagerty [2017] under the PH model by adding more scenarios such as different distributions, censoring proportions under the PH model; and 2) similarly, more situations were added to time dependent PH model such as different time dependent functions. Second, we develop an estimation method of  $HDS(t)$  for the PH model with interval censored data. Third, we apply the proposed method to HIV data from Health Sciences South Carolina (HSSC).

# TABLE OF CONTENTS

ACKNOWLEDGMENTS . . . . .	iii
ABSTRACT . . . . .	iv
LIST OF TABLES . . . . .	vii
LIST OF FIGURES . . . . .	ix
CHAPTER 1 INTRODUCTION AND NOTATION . . . . .	1
1.1 Background . . . . .	1
1.2 Motivation and Outline of Thesis . . . . .	4
CHAPTER 2 DEFINITION OF $\widehat{HDS}(t)$ . . . . .	6
2.1 Risk-Based Measure for Binary Outcomes . . . . .	6
2.2 Cumulative-Risk-Based Measure for Survival Outcomes . . . . .	7
2.3 Incident-Risk-Based Measure for Survival Outcomes . . . . .	7
CHAPTER 3 MODELS AND ESTIMATION . . . . .	11
3.1 $HDS(t)$ in Terms of Marginal Expectations . . . . .	11
3.2 $HDS(t)$ under the Cox PH Model . . . . .	12
3.3 $HDS(t)$ under Time Dependent PH Model . . . . .	15
CHAPTER 4 SIMULATIONS TO EVALUATE $\widehat{HDS}(t)$ AND $\widehat{HDS}^{LC}(t)$ . . . . .	18

4.1	Cox PH Model with Right Censoring . . . . .	19
4.2	Time Dependent PH Model with Right Censoring . . . . .	20
4.3	Cox PH Model with Interval Censoring . . . . .	26
CHAPTER 5 REAL DATA ANALYSIS . . . . .		38
5.1	Right Censored HSSC Data . . . . .	38
5.2	Interval Censored HSSC Data . . . . .	42
CHAPTER 6 CONCLUSIONS AND FUTURE STUDY . . . . .		46
BIBLIOGRAPHY . . . . .		48

## LIST OF TABLES

Table 4.1	Parameters for different survival distributions under the Cox PH model . . . . .	20
Table 4.2	$HDS(t)$ under the Cox PH model for right censored data with censoring rate 25%; sample size = 500: . . . . .	21
Table 4.3	$HDS(t)$ under the Cox PH model for right censored data with censoring rate 25%; exponential survival distribution: . . . . .	22
Table 4.4	$HDS(t)$ under the Cox PH model for right censored data with sample size = 500; exponential survival distribution: . . . . .	23
Table 4.5	Parameters for different survival distributions under time dependent PH model . . . . .	24
Table 4.6	$HDS^{LC}(t)$ under time dependent PH Model with coefficient proportional to $t$ ; with censoring rate 25%; sample size = 500: . . . . .	25
Table 4.7	$HDS^{LC}(t)$ under time dependent PH model with coefficient proportional to $t$ ; with censoring rate 25%; exponential survival distribution . . . . .	26
Table 4.8	$HDS^{LC}(t)$ under time dependent PH model with coefficient proportional to $t$ ; with sample size = 500; exponential survival distribution: . . . . .	27
Table 4.9	$HDS^{LC}(t)$ under time dependent PH model with coefficient proportional to $\log(t)$ ; with censoring rate 25%; sample size = 500: . . . . .	28
Table 4.10	$HDS^{LC}(t)$ under time dependent PH model with coefficient proportional to $\log(t)$ ; with censoring rate 25%; exponential survival distribution: . . . . .	29
Table 4.11	$HDS^{LC}(t)$ under time dependent PH model with coefficient proportional to $\log(t)$ ; with sample size = 500; exponential survival distribution: . . . . .	30



Table 4.12	$HDS^{LC}(t)$ under time dependent PH model with coefficient proportional to $\log(t)$ ; censoring rate 25%; $n = 500$ ; exponential survival distribution: . . . . .	31
Table 4.13	$HDS^{LC}(t)$ under time dependent PH model with coefficient proportional to $t^2$ ; with censoring rate 25%; sample size = 500: . . . . .	32
Table 4.14	$HDS^{LC}(t)$ under time dependent PH model with coefficient proportional to $t^2$ ; with censoring rate 25%; exponential survival distribution: . . . . .	33
Table 4.15	$HDS^{LC}(t)$ under time dependent PH model with coefficient proportional to $t^2$ ; with sample size = 500; exponential survival distribution: . . . . .	34
Table 4.16	$HDS(t)$ the Cox PH model for interval censored data with right censoring rate 25%; sample size = 500: . . . . .	35
Table 4.17	$HDS(t)$ under the Cox PH model for interval censored data with right censoring rate 25%; exponential survival distribution: . . . . .	36
Table 4.18	$HDS(t)$ under the Cox PH model for interval censored data with sample size = 500; exponential survival distribution: . . . . .	37
Table 5.1	Characteristics of suppression data . . . . .	39
Table 5.2	Supremum Test for Proportionals Hazards Assumption . . . . .	40
Table 5.3	Choices of Optimal Bandwidth . . . . .	41
Table 5.4	Characteristics of ART data . . . . .	43
Table 5.5	$\log(\text{HR})$ and 95% C.I. for estimates using ART data . . . . .	44

## LIST OF FIGURES

Figure 5.1	Kaplan-Meier curve for suppression data . . . . .	39
Figure 5.2	Time-varying Log(HR) Estimates for suppression data . . . . .	41
Figure 5.3	$HDS^{LC}(t)$ for suppression data . . . . .	42
Figure 5.4	$HDS^{LC}(t)$ Ratio for suppression data . . . . .	43
Figure 5.5	$HDS(t)$ for ART data . . . . .	44
Figure 5.6	$HDS(t)$ Ratio for ART data . . . . .	45

# CHAPTER 1

## INTRODUCTION AND NOTATION

### 1.1 BACKGROUND

Survival data, for which the outcome variable of interest is time to event, are commonly encountered in many fields, such as public health, engineering, and so on. The event could be a negative individual experience such as occurrence of disease, death, failure, etc, or a positive event like "time to recovery". Time could be measured in days, weeks, months, etc. Censoring, a unique term in survival data, occurs when the exact event time is unobserved. There are three types of censoring: left censoring, interval censoring, and right censoring. Left (right) censoring occurs when the event happens before (after) the first (last) observation times. Interval censoring refers to the case when the event times occur between two adjacent observation times. We list the detail notations and real data example in the following sections.

#### 1.1.1 RIGHT CENSORED DATA

Right censored data is most commonly seen in practice, which may be caused by the end of the study, and loss to follow up.

#### NOTATION

Let  $T_1, \dots, T_n$  be i.i.d observed times and  $\delta_1, \dots, \delta_n$  be i.i.d censoring indicators. For subjects who have failures,  $\delta_i = 1$ , and  $T_i$  is the exact time to event. For subjects who are censored,  $\delta_i = 0$  and  $T_i$  is defined as the censoring time. Let  $M_i$  be the vector

of predictors for the  $i^{th}$  subject. Usually, we assume that conditional on  $M$ ,  $T$  and  $\delta$  are independent, which is referred to as noninformative censoring.

#### REAL EXAMPLE

Let us consider the time to AIDS-related cancer (ARC) onset among HIV patients. The start time is defined as the HIV diagnose date, which is denoted as  $T_{0i}$ ,  $i = 1, \dots, n$ . The observed date  $T_{1i}$  is define as the minimum of ARC diagnose date or last observation date. For patients who have ARC eventually,  $T_{1i}$  is the ARC diagnose date and  $\delta_i = 1$ . For patients whose are cancer free at the last observation time,  $T_{1i}$  is the last observation time and  $\delta_i = 0$ . Thus,  $\{T_{1i} - T_{0i}, \delta_i\}$  is the observed follow up for ARC onset for the  $i^{th}$  patient.

#### 1.1.2 INTERVAL CENSORED DATA

By interval censoring, time to event of interest is known only to lie within an interval instead of being observed exactly. Common examples occur in medical or health studies with periodic follow-up. Current status data is a special case of interval censoring, where each subject is observed only once for the status of the occurrence of the event of interest. Current status data is referred to as case I interval-censored data and the general case as case II interval-censored data.

#### NOTATION

Let  $L_i, R_i$  be the observed time interval and  $\delta_{i1}, \delta_{i2}, \delta_{i3}$  be the interval censoring indicators for the  $i^{th}$  subject. For subjects who are left censored,  $\delta_{i1} = 1, \delta_{i2} = 0, \delta_{i3} = 0, L_i = NA$ , and  $R_i$  is the time to first observed event. For subjects who are right censored,  $\delta_{i1} = 0, \delta_{i2} = 0, \delta_{i3} = 1, R_i = NA$  and  $L_i$  is defined as the time to last observation.  $\delta_{i1} = 0, \delta_{i2} = 1, \delta_{i3} = 0$  are censoring indicators for those who are

censored between  $L_i$  and  $R_i$ . Let  $M_i$  be the vector of predictors for the  $i^{th}$  subject. It is also commonly assumed that conditional on  $M$ ,  $L$ ,  $R$  and  $\delta_1, \delta_2, \delta_3$  are independent.

#### REAL EXAMPLE

We are interested in time to an undetectable viral load (UVL, viral load level  $< 50$  copies/mL) among HIV patients who received antiretroviral therapy (ART). The date of initiating ART is defined as the start-time point  $T_{0i}, i = 1, \dots, n$ . Regular lab visit at  $\{T_{1i}, T_{2i}, \dots, \}$  are recorded to monitor the patients' HIV status. For patients who finally have UVL but not at the first observation, we define the date of first UVL as  $T_{2i}$ , the date of last observation before  $T_{2i}$  as  $T_{1i}$ . Then, we define the time interval as  $\{L_i = T_{1i} - T_{0i}, R_i = T_{2i} - T_{0i}\}$ , and the censoring indicators as  $\delta_{1i} = 0, \delta_{2i} = 1, \delta_{3i} = 0$ . For patients who have UVL at the first observation, we define the first observation date as  $T_{2i}$ . Then, time interval is  $\{L_i = 0, R_i = T_{2i} - T_{0i}\}$ , and the censoring indicators are  $\delta_{1i} = 1, \delta_{2i} = 0, \delta_{3i} = 0$ . We also define the last observation date as  $T_{1i}$  for patients who still haven't had UVL at the end of the study. Then, time interval is  $\{L_i = T_{2i} - T_{0i}, R_i = NA\}$ , and the censoring indicators are  $\delta_{1i} = 0, \delta_{2i} = 0, \delta_{3i} = 1$ . Thus,  $\{L_i, R_i, \delta_{1i}, \delta_{2i}, \delta_{3i}\}$  is the observed outcome for the  $i^{th}$  patient.

#### 1.1.3 EVALUATION OF SURVIVAL MODELS

Survival analysis is a collection of statistical methods for analyzing survival data. Many survival models exist to estimate the effects of potential risk factors on the outcome and even to predict survival probability among a population of interest. For example, the Cox PH model [Cox, 1972] is the most popular semiparametric survival model, which is a special case in the generalized odds-rate (GOR) model [Dabrowska and Doksum, 1988].

## 1.2 MOTIVATION AND OUTLINE OF THESIS

Recently, hazard discrimination summary ( $HDS(t)$ ) has been proposed to evaluate the discrimination performance under survival models for right censored data at a certain time [Liang and Heagerty, 2017]. It was shown that the performance of  $HDS(t)$  under the Cox PH model and time dependent PH model with exponential survival distribution for right censored data is very effective [Liang and Heagerty, 2017]. Since interval censored data is commonly seen in periodic visit, we are motivated to generalize  $HDS(t)$  to interval censored data. Based on the definition and estimation of  $HDS(t)$ , the estimation of  $HDS(t)$  under the PH model for the interval censored data is feasible.

In this thesis, first, we further investigate their method via comprehensive simulations including 1) We extend the simulation in Liang and Heagerty [2017] under the PH model by adding more scenarios such as different distributions, censoring proportions under the PH model; and 2) similarly, more situations were added to time dependent PH model such as different time dependent functions. Second, we develop a estimation method of  $HDS(t)$  for the PH model with interval censored data. Third, we apply the proposed method to HIV data from Health Sciences South Carolina (HSSC).

The outline of the thesis is as follows.

- In Chapter 2, we introduce several discrimination performance measurements for binary outcome, which are the basis of  $HDS(t)$ . Then, we introduce the definition of  $HDS(t)$  based on hazard function [Liang and Heagerty, 2017].
- In Chapter 3, we introduce two popular survival models, the Cox PH model and time dependent PH model. The expression of  $HDS(t)$  could be simplified after plugging in the hazard function under the Cox PH model. The estimator and standard error for  $HDS(t)$  can be obtained by plugging in the estimators

of the Cox PH model and time dependent PH model.

- In Chapter 4, we evaluate the performance of  $HDS(t)$  under different survival distributions, different sample sizes, and different right censoring rates for the PH model and time dependent PH model. Then, we further investigate its performance under the Cox PH model for interval censored data.
- In Chapter 5, we use two cleaned data sets from Health Sciences South Carolina (HSSC) to illustrate the usage of  $HDS(t)$ . The  $HDS(t)$  and its 95% confidence interval are calculated and used to test the hypothesis that there is no discrimination among cases and control under the Cox PH model or time dependent PH model for right censored data. Similar approaches are applied to assess the discrimination performance for the interval censored data. The  $HDS(t)$  ratio, which evaluates the discrimination performance of a specific predictor, is applied. Bootstrap method is used to construct the confidence interval for  $HDS(t)$  ratio.
- In Chapter 6, we summarize and discuss the current work and outline some future work.

## CHAPTER 2

### DEFINITION OF $\widehat{HDS}(t)$

In this chapter, the motivation and definition of  $HDS(t)$  are illustrated. Section 2.1 introduces a discrimination performance measure for binary outcomes; section 2.2 illustrates the extension of this measure to survival outcomes; section 2.3 demonstrates  $HDS(t)$ , which is an incident-risk-based measure generalized from the previous cumulative-risk-based measure.

#### 2.1 RISK-BASED MEASURE FOR BINARY OUTCOMES

For binary outcomes, we define the subjects who have events as cases and the subjects who have no event as controls. In terms of the model based on binary outcomes like logistic regression, one of the most important issues is to evaluate the discrimination performance between cases and controls. With the motivation to summarize the magnitude of mean risk between cases and controls, discrimination slope ( $DS$ ) was proposed [Yates, 1982]. Let  $M$  denotes the vector of predictors and  $D$  is the binary outcome with  $D = 1$  for cases and  $D = 0$  for controls,  $DS$  is defined as

$$DS = E\{P(D = 1|M)|D = 1\} - E\{P(D = 1|M)|D = 0\}.$$

Based on the magnitude of risk,  $DS$  measures how well the cases and controls are discriminated. The value of  $DS$  closes to zero, indicating that there is no discrimination of this model. The model performs well when  $DS$  closes to one, and there is something wrong with the model when  $DS$  closes to negative one.



## 2.2 CUMULATIVE-RISK-BASED MEASURE FOR SURVIVAL OUTCOMES

With outcome of interest is time to event, survival models become important methods to predict risk. Thus, in order to generalize simple binary outcome prognostic measure to a time varying measure, Chambless et al. [2011] proposed time-specific  $DS$ . Let  $T$  denote time to event, the discrimination slope at time  $t$  is defined as

$$DS(t) = E_{M|T \leq t}\{P(T \leq t|M)|T \leq t\} - E_{M|T > t}\{P(T \leq t|M)|T > t\}.$$

The time-specific  $DS(t)$  is the difference in mean failure probabilities until time  $t$  between subjects who already have failures by time  $t$  and subjects who are still free of outcomes at time  $t$ . The interpretation of  $DS(t)$  is similar with that of  $DS$ . We hope  $DS(t)$  is close to one over time, which indicates the survival model performs well.

Moreover,  $DS(t)$  can be used to evaluate the improvement of discrimination performance when a specific predictor is added in the survival model by calculating the difference in  $DS(t)$  of the model including this predictor as well as the model not including this predictor with other predictors fixed. This difference is named as the integrated discrimination improvement ( $IDI$ ):

$$IDI(t) = DS_{new}(t) - DS_{old}(t)$$

Here, subscripts "new" and "old" are used to denote  $DS(t)$  for the survival with this predictor and the survival model without this predictor, respectively.  $IDI(t)$  scales the improvement of discrimination performance for the predictor we may be interested in. The positive value of  $IDI(t)$  indicates adding the predictor into the survival model improves the discrimination performance.

## 2.3 INCIDENT-RISK-BASED MEASURE FOR SURVIVAL OUTCOMES

Compared to binary outcomes, survival outcomes have time information so that not only cumulative risk but also incident risk can be analyzed. However,  $DS(t)$  cannot

take incident risk into consideration. Based on the hazard function which can reflect the instantaneous risk, hazard discrimination summary ( $HDS(t)$ ) was proposed by Liang and Heagerty [2017].  $HDS(t)$  is defined as a ratio of the mean hazard between cases and controls at time  $t$ . Before illustrating the definition of  $HDS(t)$ , we introduce the definition of hazard function.

### 2.3.1 SURVIVAL FUNCTION AND HAZARD FUNCTION

The hazard function is defined as

$$\lambda(t) = \lim_{\Delta t \rightarrow 0} P\{T \in [t, t + \Delta t] | T \geq t\} / \Delta t$$

The survival function is defined as

$$S(t) = \exp\{-\Lambda(t)\} = \exp\left\{-\int_0^t \lambda(u) du\right\}$$

where  $\Lambda(t)$  is the cumulative hazard, which can be calculated via the integration of  $\lambda(t)$ . The Cox PH model [Cox, 1972] is defined based on the hazard function. More details of definition and estimation for the Cox PH model are introduced in chapter 3.

### 2.3.2 DEFINITION OF $HDS(t)$

$HDS(t)$  is defined as a ratio of expected hazards among cases to expected hazards among controls at time  $t$  [Liang and Heagerty, 2017]. It is specified as the following:

$$HDS(t) = \frac{\text{mean case hazard}}{\text{mean control hazard}} = \frac{E_{M|T=t}\{\lambda(t|M)|T=t\}}{E_{M|T>t}\{\lambda(t|M)|T>t\}} \quad (2.1)$$

Here, we define the subjects who fail at time  $t$  as cases and subjects who are still free of outcome as controls. The survival model used should have higher hazard on the cases than on the controls. When  $HDS(t)$  is around a value of one, it can be concluded that there is no discriminatory performance. The closer the  $HDS(t)$  is to

infinity, the better the performance is. Conversely, the model is in appropriate when the value of  $HDS(t)$  is close to zero.

It has been shown that the denominator can be transformed into the marginal hazard [Liang and Heagerty, 2017]:

$$\bar{H}_0(t) = E_{M|T>t}\{\lambda(t|M)|T > t\} = \lambda(t) \quad (2.2)$$

Plugging (2.2) into (2.1),  $HDS(t)$  can be expressed as the following [Liang and Heagerty, 2017]:

$$HDS(t) = E_{M|T=t}\left\{\frac{\lambda(t|M)}{\lambda(t)}|T = t\right\} \quad (2.3)$$

This expression gives  $HDS(t)$  a second interpretation that the value of  $HDS(t)$  is the increase in the average risk assigned to incident cases at time  $t$  associated with knowing marker,  $M$ , as compared to the marginal risk associated without knowing the marker [Liang and Heagerty, 2017].

### 2.3.3 $HDS(t)$ RATIO

Similar to the function of  $IDI(t)$ , the ratio of  $HDS(t)$  also evaluates the improvement of the discrimination performance of a specific predictor [Liang and Heagerty, 2017]. For the purpose of simple illustration, we use the following notation for the numerator part in (2.1):

$$\bar{H}_1(t) = E_{M|T=t}\{\lambda(t|M)|T = t\} \quad (2.4)$$

The  $HDS(t)$  ratio can be simplified as follows:

$$\begin{aligned} HDS(t; Mod1)/HDS(t; Mod2) &= \frac{\bar{H}_1(t; Mod1)/\bar{H}_0(t; Mod1)}{\bar{H}_1(t; Mod2)/\bar{H}_0(t; Mod2)} \\ &= \frac{\bar{H}_1(t; Mod1)/\lambda(t)}{\bar{H}_1(t; Mod2)/\lambda(t)} \\ &= \frac{\bar{H}_1(t; Mod1)}{\bar{H}_1(t; Mod2)} \\ &= \frac{E_{Mod1|T=t}\{\lambda(t|Mod1)|T = t\}}{E_{Mod2|T=t}\{\lambda(t|Mod2)|T = t\}} \end{aligned} \quad (2.5)$$

Here, when the model  $Mod2$  with the same markers as  $Mod1$  except  $M^*$ , the  $HDS(t)$  ratio could be interpreted as the ratio of the average risk predicted by  $Mod1$  to the average risk predicted by  $Mod2$  for time-specific incident cases [Liang and Heagerty, 2017].  $HDS(t)$  ratio is greater than one, indicating that the predictor  $M^*$  improves the discrimination performance of the survival model.

## CHAPTER 3

### MODELS AND ESTIMATION

In this chapter, we illustrate the estimation methods of  $HDS(t)$  under the Cox PH model for right censored and interval censored data as well as time dependent PH model for right censored data, respectively. Specifically, in section 3.1, we re-express  $HDS(t)$  in terms of marginal expectations of functions of conditional hazard  $\lambda(t|M)$  and cumulative hazard  $\Lambda(t|M)$ . The estimator for  $HDS(t)$  can be obtained by inserting the estimators for  $\lambda(t|M)$  and  $\Lambda(t|M)$ . Section 3.2 introduces estimation methods for the Cox PH model with right censored data and interval censored data. By plugging in those estimators, the estimator and standard error for  $HDS(t)$  can be obtained. Section 3.3 illustrates an estimation method for time dependent PH model with right censored data, and the estimation and standard error for localized  $HDS(t)$ .

#### 3.1 $HDS(t)$ IN TERMS OF MARGINAL EXPECTATIONS

It has been shown that  $\bar{H}_1(t)$  and  $\bar{H}_0(t)$  in (2.2) and (2.4) can be rewritten as expressions with expected functions of conditional hazard and cumulative hazard to facilitate estimation as follows [Liang and Heagerty, 2017]:

$$\bar{H}_1(t) = E_{M|T=t}\{\lambda(t|M)|T = t\} = \frac{E_M[\lambda^2(t|M) \exp\{-\Lambda(t|M)\}]}{E_M[\lambda(t|M) \exp\{-\Lambda(t|M)\}]} \quad (3.1)$$

$$\bar{H}_0(t) = E_{M|T>t}\{\lambda(t|M)|T > t\} = \frac{E_M[\lambda(t|M) \exp\{-\Lambda(t|M)\}]}{E_M[\exp\{-\Lambda(t|M)\}]} \quad (3.2)$$

Given estimators for  $\lambda(t|M)$  and  $\lambda(t|M)$ , the corresponding sample  $\bar{H}_1(t)$  and  $\bar{H}_0(t)$  are as follows:

$$\hat{H}_1(t) = \frac{\sum_{i=1}^n [\hat{\lambda}^2(t|m_i) \exp\{-\hat{\Lambda}(t|m_i)\}]}{\sum_{i=1}^n [\hat{\lambda}(t|m_i) \exp\{-\hat{\Lambda}(t|m_i)\}]} \quad (3.3)$$

$$\hat{H}_0(t) = \frac{\sum_{i=1}^n [\hat{\lambda}(t|m_i) \exp\{-\hat{\Lambda}(t|m_i)\}]}{\sum_{i=1}^n [\exp\{-\hat{\Lambda}(t|m_i)\}]} \quad (3.4)$$

### 3.2 HDS(t) UNDER THE COX PH MODEL

Let  $\lambda_0(t)$  denote baseline hazard function,  $\Lambda_0(t)$  denote the cumulative baseline hazard function, and  $\beta$  denote a vector of coefficients for predictor vector  $M$ . Then, hazard at time  $t$  is defined as  $\lambda(t) = \lambda_0(t) \exp(\beta' M)$ . Here, the  $\exp(\beta)$  is the hazard ratio (HR) when one unit increase in the corresponding  $m_i$  with other variables fixed at the same level.

#### 3.2.1 ESTIMATION GIVEN RIGHT CENSORED DATA

For the right censored data, the partial likelihood estimator of  $\beta$  [Cox, 1972] and nonparametric estimator of  $\lambda_0(t)$  [Breslow, 1972] were proposed. The reason for naming this method as partial likelihood is that we consider probabilities only for subjects who fail, and we do not consider probabilities for subjects who are censored. That is, the Cox likelihood is a product of probabilities for failures only at failure time points. If  $t_{(1)}, \dots, t_{(k)}$  are the ordered failure times with the corresponding covariates  $M_{(1)}, \dots, M_{(k)}$ , the partial likelihood of the event occurring with the  $i^{th}$  subject at time  $t_{(i)}$  can be expressed as

$$L_i(\beta) = \frac{\lambda_0(t_{(i)}) \beta' M_{(i)}}{\sum_{l \in R(t_{(i)})} \lambda_0(t_{(i)}) e^{\beta' M_l}} = \frac{\beta' M_{(i)}}{\sum_{l \in R(t_{(i)})} e^{\beta' M_l}}$$

where  $R(t_{(i)})$  is the set of the individual at risk at time  $t_{(i)}$ . The denominator for each term corresponding to time  $t_j$  ( $j = 1, 2, \dots$ ) is the sum of the hazards for those subjects still at risk at time  $t_j$ , and the numerator is the hazard for the subject who

experiences the event at  $t_j$ . Because the baseline hazard can be cancelled, the baseline hazard does not need to be specified in a Cox PH model. Assuming noninformative censoring, the log of partial likelihood is

$$\ell(\boldsymbol{\beta}) = \sum_{i=1}^k \{ \boldsymbol{\beta}' \mathbf{M}_{(i)} - \log( \sum_{l \in R(t_{(i)})} e^{\boldsymbol{\beta}' \mathbf{M}_l} ) \}$$

The coefficients can be obtained via maximizing the above formula. The Hessian matrix of the partial log likelihood is

$$\ell''(\boldsymbol{\beta}) = - \sum_{i=1}^k \left( \frac{\sum_{l \in R(t_{(i)})} e^{\boldsymbol{\beta}' \mathbf{M}_l} \mathbf{M}_l' \mathbf{M}_l}{\sum_{l \in R(t_{(i)})} e^{\boldsymbol{\beta}' \mathbf{M}_l}} - \frac{[\sum_{l \in R(t_{(i)})} e^{\boldsymbol{\beta}' \mathbf{M}_l} \mathbf{M}_l'] [\sum_{l \in R(t_{(i)})} e^{\boldsymbol{\beta}' \mathbf{M}_l} \mathbf{M}_l]}{[\sum_{l \in R(t_{(i)})} e^{\boldsymbol{\beta}' \mathbf{M}_l}]^2} \right)$$

The inverse of the Hessian matrix can be used as an approximate variance-covariance matrix for the estimates.

For the estimation of baseline cumulative hazard, Breslow [1972] proposed a non-parametric maximum likelihood estimation (NPMLE) method with  $\Lambda_0(t)$  being estimated by

$$\hat{\Lambda}_0(t) = \sum_{t_i \leq t} \frac{1}{\sum_{j \in R(t_i)} \exp(\hat{\boldsymbol{\beta}}' \mathbf{M}_j)}.$$

The estimators are available in several statistical packages such as *coxph* in R software and *phreg* in SAS software.

### 3.2.2 ESTIMATION GIVEN INTERVAL CENSORED DATA

Given the PH model with interval censored data, the likelihood can be written as

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n F(R_i | M_i)^{\delta_{i1}} \{ F(R_i | M_i) - F(L_i | M_i) \}^{\delta_{i2}} \{ 1 - F(L_i | M_i) \}^{\delta_{i3}}$$

where  $F(t | M_i) = 1 - \exp\{-\Lambda_0(t) \exp(\boldsymbol{\beta}' \mathbf{M}_i)\}$  is a conditional cumulative density function.

Assuming  $\Lambda_0(\cdot) = \sum_{l=1}^k \gamma_l b_l(\cdot)$ , Wang et al. [2016] proposed an EM based augmentation method, which is computational efficient. Here,  $b_l(\cdot)$ 's are integrated spline basis functions, which are estimated using I-spline. The degree and placement of

knots determine the  $k$  spline basis functions, where  $k$  is equal to the degree plus the number of interior knots. It has been shown that given initial values for  $\beta$  and  $\gamma$ , the estimators for  $\beta$  and  $\gamma$  are obtained through iterative process until convergence [Wang et al., 2016]. The estimated variance-covariance matrix for  $(\beta, \gamma)$  is derived by the inverse of the Hessian matrix, using Louis' method [Louis, 1982]. Interval censored data can be analyzed under the Cox PH model by using *ICsurv* package in R [Wang et al., 2016]. An alternative package is *ICGOR*, which can also give approximate estimators for the Cox PH model with interval censoring [Zhou et al., 2017].

### 3.2.3 ESTIMATION FOR $HDS(t)$ UNDER THE COX PH MODEL

As mentioned in section 3.1,  $HDS(t)$  can be expressed in terms of marginal expected functions of conditional hazard and cumulative hazard. Under the Cox PH model,  $\lambda_0(t)$  is cancelled, and  $HDS(t)$  is expressed as follows

$$HDS(t) = \frac{\overline{H}_1(t)}{\overline{H}_0(t)} = \frac{E_M[\exp\{2\beta'M - e^{\beta'M}\Lambda_0(t)\}]E_M[\exp\{-e^{\beta'M}\Lambda_0(t)\}]}{E_M[\exp\{\beta'M - e^{\beta'M}\Lambda_0(t)\}]^2} \quad (3.5)$$

The estimated  $HDS(t)$  under the Cox PH model can be obtained by plugging in the estimators  $\hat{\beta}$ ,  $\hat{\Lambda}_0(t)$ , and sample predictors  $m_i, i = 1, 2, \dots, n$ , which is

$$\widehat{HDS}(t) = \frac{\sum_{i=1}^n \exp\{2\hat{\beta}'m_i - e^{\hat{\beta}'m_i}\hat{\Lambda}_0(t)\} \sum_{i=1}^n \exp\{-e^{\hat{\beta}'m_i}\hat{\Lambda}_0(t)\}}{[\sum_{i=1}^n \exp\{\hat{\beta}'m_i - e^{\hat{\beta}'m_i}\hat{\Lambda}_0(t)\}]^2} \quad (3.6)$$

For right censored data,  $\hat{\beta}$  is the partial likelihood estimator from the Cox PH model [Cox, 1972] and  $\hat{\Lambda}_0(t)$  is the Breslow estimator for the cumulative baseline hazard [Breslow, 1972]. For interval censored data, the estimated  $\widehat{HDS}(t)$  can be obtained by plugging in the  $\hat{\beta}$  and  $\hat{\Lambda}_0(t)$  estimated from the R package *ICsurv* [Wang et al., 2016].



### 3.2.4 STANDARD ERRORS FOR $\widehat{HDS}(t)$

Let  $f_\theta = E_M[\exp\{\theta \cdot \beta' M - e^{\beta' M} \Lambda_0(t)\}]$  and  $\hat{f}_\theta = \sum_{i=1}^n \exp\{\theta \cdot \hat{\beta}' m_i - e^{\hat{\beta}' m_i} \hat{\Lambda}_0(t)\}$  for  $\theta = 0, 1, 2$ , we have  $HDS(t) = \frac{f_2 f_0}{f_1^2}$  and  $\widehat{HDS}(t) = \frac{\hat{f}_2 \hat{f}_0}{\hat{f}_1^2}$

$$\widehat{HDS}(t) = \frac{\sum_{i=1}^n \exp\{2\hat{\beta}' m_i - e^{\hat{\beta}' m_i} \hat{\Lambda}_0(t)\} \sum_{i=1}^n \exp\{-e^{\hat{\beta}' m_i} \hat{\Lambda}_0(t)\}}{[\sum_{i=1}^n \exp\{\hat{\beta}' m_i - e^{\hat{\beta}' m_i} \hat{\Lambda}_0(t)\}]^2} \quad (3.7)$$

It was proved that the estimated  $HDS(t)$  has an asymptotic normal distribution [?], which is

$$\sqrt{n}(\widehat{HDS}(t) - HDS(t)) \sim N(0, A(\sum_1 + B\sum_0 B')A')$$

where  $A$  is the Jacobian of the map  $(f_0, f_1, f_2) \rightarrow \frac{f_2 f_0}{(f_1)^2}$ ;  $\sum_1$  is the variance-covariance matrix of  $(f_0, f_1, f_2)$ ;  $\sum_0$  is the asymptotic variance matrix for  $(\beta, \Lambda_0(t))$  which was derived by Tsiatis et al. [1981] and van der Vaart et al. [2007];  $B$  is a matrix of the derivatives of  $(f_0, f_1, f_2)'$  with respect to  $\beta$  and  $\Lambda_0(t)$ , respectively. After analyzing right censored or interval censored data under the Cox PH model, plugging in the estimators of  $\beta$ ,  $\Lambda_0(t)$ , and estimated variance matrix of  $\beta$ , the standard error of estimated  $HDS(t)$  can be obtained.

### 3.3 $HDS(t)$ UNDER TIME DEPENDENT PH MODEL

The validity of the Cox PH model is based on the satisfaction of PH assumption. When the PH assumption is violated, the time dependent PH model is an alternative. The time dependent PH model is defined as

$$\lambda(t) = \lambda_0(t) \exp(\beta'(t)M).$$

Then, the HR is a function of time instead of a constant. By relaxing the coefficients, time dependent PH model is more Flexible than the Cox PH model.

### 3.3.1 ESTIMATION GIVEN RIGHT CENSORED DATA

The estimation of  $\beta_h(s)$  is based on a weighted local log partial likelihood function, as proposed by Cai and Sun [2003] and Tian et al. [2005].

$$\ell(\beta_h(s)) = (nh_n)^{-1} \sum_{i=1}^k K\left(\frac{s - t_{(i)}}{h_n}\right) \{\beta'_h(s) \mathbf{M}_{(i)} - \log(\sum_{l \in R(t_{(i)})} e^{\beta'_h(s) M_l})\} \quad (3.8)$$

where  $K(\cdot)$  is a symmetric kernel function with support  $[-1, 1]$ , mean 0, and bounded first derivative, for example, the Epanechnikov kernel  $K(u) = \frac{3}{4}(1 - u^2)$  for  $-1 \leq u \leq 1$ , otherwise,  $K(u) = 0$ ; the bandwidth  $h_n = O(n^{-v})$  with  $v > 0$ . It has been shown that uniformly consistent estimators could be obtained when  $1/4 < v < 1/2$  [Tian et al., 2005]. The variance-covariance matrix for estimators is approximately  $\mathbf{I}^{-1}\{\hat{\beta}(t), t\} \int_{-1}^1 K^2(u) du$  with second derivative of weighted log partial likelihood function in (3.8)  $\mathbf{I}(\beta(t), t)$  [Tian et al., 2005]. Similar to estimation of  $\Lambda_0(t)$  under Cox PH model, the generalized Breslow estimator for  $\Lambda_h(t)$  under time dependent model is as follows  $\hat{\Lambda}_h(t) = \sum_{t_i \leq t} \frac{1}{\sum_{j \in R(t_i)} \exp(\hat{\beta}'_h M_j)}$  [Cai and Sun, 2003].

### 3.3.2 ESTIMATOR FOR $HDS(t)$ UNDER TIME DEPENDENT PH MODEL

It has been shown that estimator for  $HDS(t)$  under time dependent PH model ( $HDS^{LC}(t)$ ) is available by replacing  $\hat{\beta}$  with  $\hat{\beta}_h(t)$  and replacing  $\hat{\Lambda}_0(t)$  with  $\hat{\Lambda}_h(t)$  [Liang and Heagerty, 2017]:

$$\widehat{HDS}^{LC}(t) = \frac{\sum_{i=1}^n \exp\{2\hat{\beta}'_h(t) m_i - e^{\hat{\beta}'_h(t) m_i} \hat{\Lambda}_h(t)\} \sum_{i=1}^n \exp\{-e^{\hat{\beta}'_h(t) m_i} \hat{\Lambda}_h(t)\}}{[\sum_{i=1}^n \exp\{\hat{\beta}'_h(t) m_i - e^{\hat{\beta}'_h(t) m_i} \hat{\Lambda}_h(t)\}]^2} \quad (3.9)$$

where  $\hat{\beta}_h(t)$  is the smoothed estimate of  $\beta(t)$  as proposed by Cai and Sun [2003], and  $\hat{\Lambda}_h(t)$  is the corresponding estimate of  $\Lambda_0(t)$  [Tian et al., 2005].

### 3.3.3 STANDARD ERRORS FOR $\widehat{HDS}^{LC}(t)$

It is also shown that the estimator  $\widehat{HDS}^{LC}(t)$  has an asymptotically normal distribution [Liang and Heagerty, 2017].

$$\sqrt{nh}(\widehat{HDS}^{LC}(t) - HDS^{LC}(t)) \sim N(0, AC\sum_2 C' A') \quad (3.10)$$

where  $\sum_2$  is a variance-covariance matrix for  $\beta$ ;  $C$  is a matrix of the derivatives of  $(f_0, f_1, f_2)'$  with respect to  $\beta$ . The estimated standard error for  $\widehat{HDS}^{LC}(t)$  is shown to be a function of estimated standard error of  $\hat{\beta}$  through delta method. By plugging the smoothed estimator of  $\beta_h(t)$ ,  $\Lambda_h(t)$ , and variance-covariance matrix of  $\hat{\beta}_h(t)$  into (3.10), the estimated standard error of  $HDS(t)$  can be obtained as proposed by Cai and Sun [2003] and Tian et al. [2005].

## CHAPTER 4

### SIMULATIONS TO EVALUATE $\widehat{HDS}(t)$ AND $\widehat{HDS}^{LC}(t)$

The performance of  $HDS(t)$  has been tested under the Cox PH model and time dependent PH model [Liang and Heagerty, 2017]. However, only the right censoring type, exponentially distributed survival function, and time varying coefficients proportional to time were considered. In practice, the survival function could follow other distributions rather than exponential distribution such as Weibull distribution. Moreover, the effects of predictors could increase with time rapidly (proportional to squared time) or slowly (proportional to  $\log(t)$ ). Apart from right censoring, interval censored data is also ubiquitous. Therefore, it is motivated to extend the estimation of  $HDS(t)$  under more situations.

To evaluate the performance of  $\widehat{HDS}(t)$  and  $\widehat{HDS}^{LC}(t)$ , we conduct comprehensive simulation in this chapter. In section 4.1, we evaluate the performance of  $\widehat{HDS}(t)$  for right censored data under the Cox PH model. In section 4.2, the performance of  $\widehat{HDS}^{LC}(t)$  for right censored data under time dependent PH model will be evaluated. Finally, in section 4.3, we extend the simulation to the interval censored data under the Cox PH model.

For each simulation setting, we assumed a Cox PH model or a time dependent PH model with different censoring rates, different sample sizes, and different distribution types. For the purpose of simple illustration, we summarize the common settings as follows

- Case I: investigate different distributions:

Assuming censoring rate as 25% and sample size  $n=500$ , we consider exponen-

tial, Weibull, Log-Log, and Log Normal distribution.

- Case II: investigate different sample sizes:

Assuming the moderate censoring rate 25% and exponential distribution, we consider different sample sizes including  $n=200, 500, 800$ .

- Case III: investigate different censoring rates:

Assuming sample size  $n=500$  and exponential distribution, we consider censoring rate at 15%, 25%, and 30%.

In the following sections, we detail the simulations under different models. The results of 1000 simulations under different cases are presented in the following sections. In each table, we report the time points to be evaluated, true value of  $HDS(t)$ , the corresponding estimated  $HDS(t)$ , standard error (SE), and coverage probability (CP).

#### 4.1 COX PH MODEL WITH RIGHT CENSORING

We use the same data generating mechanism as in Liang and Heagerty [2017]. The Cox PH model is assumed as

$$\lambda(t|M) = \lambda_0(t) \cdot \exp(0.5 \cdot M_1 + 1.5 \cdot M_2)$$

We consider exponential, Weibull, Log-Log, and Log Normal distribution, and their corresponding baseline hazard function and survival distributions are shown in Table 4.1. Here,  $\Phi$  is the cumulative density function of standard normal distribution.

We generated a random number from a uniform distribution with the support  $[0, 1]$  as the true survival probability individually. For each subject, both of the predictors follow a uniform distribution with the support  $[0, 2]$ . The censoring time is generated from a uniform distribution with the support  $[0, c]$ , while  $c$  is a constant, adjusting

Table 4.1 Parameters for different survival distributions under the Cox PH model

Distribution	Baseline Hazard $\lambda_0(t)$	Parameters
Exponential	$k$	$k=0.5$
Weibull	$k^p p t^{p-1}$	$k=0.5, p=2$
Log-Logistic	$\frac{k p (k t)^{p-1}}{1 + (k t)^p}$	$k=0.5, p=2$
Log Normal	$\frac{\partial}{\partial t} \left\{ -\log \left\{ 1 - \Phi \left( \frac{\log(t) - \mu}{\sigma} \right) \right\} \right\}$	$\mu=0.5, \sigma=2$

which censoring rate can be controlled. For whose censoring time is smaller than event time, the censoring indicator is equal to zero, otherwise, one.

The results of three cases under the Cox PH model for right censored data are reported in Table 4.2, 4.3, 4.4. From Table 4.2 we can see that the estimated  $\widehat{HDS}(t)$  shows very little bias and the coverage probability is close to 95% for each evaluated time under different distributions. From the results of different sample size (see Table 4.3), the coverage probability also seems close to 95% but under small sample size the bias is a little greater than the bias under larger sample size. Compared to the results of small censoring rate, the coverage probability under large censoring rate like 60% is relatively less than 95% after 0.9 due to the lack of sample data. (see Table 4.4).

## 4.2 TIME DEPENDENT PH MODEL WITH RIGHT CENSORING

Three time dependent PH models are assumed as follows:

Model I: HR proportional to time

$$\lambda(t|M) = \lambda_0(t) \cdot \exp(t \cdot M_1 + 0.5 \cdot M_2)$$

Model II: effect changes slowly

$$\lambda(t|M) = \lambda_0(t) \cdot \exp(\log(t) \cdot M_1 + 0.5 \cdot M_2)$$

Table 4.2  $HDS(t)$  under the Cox PH model for right censored data with censoring rate 25%; sample size = 500:

Time	Exponential				Weibull			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.83	0.13	0.97	1.79	1.81	0.12	0.90
0.050	1.84	1.84	0.13	0.97	1.80	1.81	0.12	0.90
0.100	1.83	1.84	0.12	0.95	1.81	1.82	0.13	0.91
0.200	1.76	1.76	0.09	0.89	1.83	1.85	0.13	0.91
0.300	1.66	1.66	0.08	0.90	1.84	1.85	0.12	0.92
0.400	1.56	1.57	0.07	0.92	1.79	1.81	0.11	0.92
0.500	1.49	1.49	0.06	0.95	1.71	1.72	0.09	0.92
0.600	1.43	1.43	0.06	0.96	1.60	1.61	0.07	0.94
0.700	1.38	1.39	0.06	0.96	1.50	1.50	0.06	0.95
0.800	1.35	1.35	0.05	0.96	1.41	1.42	0.06	0.97
0.900	1.32	1.32	0.05	0.96	1.34	1.35	0.05	0.98
1.000	1.29	1.30	0.05	0.96	1.29	1.30	0.05	0.98
Time	Log-Log				Log Normal			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.79	1.80	0.12	0.95	1.83	1.84	0.13	0.95
0.050	1.80	1.80	0.12	0.95	1.84	1.85	0.12	0.95
0.100	1.81	1.81	0.12	0.96	1.79	1.80	0.10	0.93
0.200	1.83	1.84	0.13	0.95	1.64	1.65	0.07	0.92
0.300	1.84	1.85	0.12	0.95	1.53	1.54	0.06	0.94
0.400	1.80	1.81	0.11	0.93	1.46	1.47	0.06	0.96
0.500	1.72	1.73	0.09	0.93	1.41	1.41	0.06	0.97
0.600	1.63	1.63	0.07	0.93	1.37	1.38	0.05	0.97
0.700	1.53	1.54	0.06	0.94	1.34	1.35	0.05	0.96
0.800	1.45	1.46	0.06	0.95	1.32	1.32	0.05	0.97
0.900	1.39	1.39	0.06	0.96	1.30	1.31	0.05	0.97
1.000	1.34	1.34	0.05	0.98	1.28	1.30	0.05	0.97

Model III: effect changes rapidly

$$\lambda(t|M) = \lambda_0(t) \cdot \exp(t^2 \cdot M_1 + 0.5 \cdot M_2)$$

The baseline hazard under different distributions are specified in Table 4.5. We assumed the predictor  $M_1$  follows a uniform distribution with the support  $[1, 3]$ , and the predictor  $M_2$  follows a standard normal distribution. The process of censoring time is similar to that in section 4.1.

Table 4.3  $HDS(t)$  under the Cox PH model for right censored data with censoring rate 25%; exponential survival distribution:

Time	n=200				n=500			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.85	0.20	0.96	1.82	1.83	0.13	0.97
0.050	1.84	1.87	0.20	0.96	1.84	1.84	0.13	0.97
0.100	1.83	1.86	0.19	0.95	1.83	1.84	0.12	0.95
0.200	1.76	1.78	0.15	0.94	1.76	1.76	0.09	0.89
0.300	1.66	1.67	0.12	0.93	1.66	1.66	0.08	0.90
0.400	1.56	1.58	0.10	0.94	1.56	1.57	0.07	0.92
0.500	1.49	1.50	0.10	0.95	1.49	1.49	0.06	0.95
0.600	1.43	1.44	0.09	0.97	1.43	1.43	0.06	0.96
0.700	1.38	1.39	0.08	0.96	1.38	1.39	0.06	0.96
0.800	1.35	1.36	0.08	0.97	1.35	1.35	0.05	0.96
0.900	1.32	1.33	0.08	0.97	1.32	1.32	0.05	0.96
1.000	1.29	1.30	0.08	0.97	1.29	1.30	0.05	0.96
Time	n=800				n=1000			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.83	0.10	0.95	1.82	1.83	0.09	0.95
0.050	1.84	1.84	0.10	0.95	1.84	1.84	0.09	0.95
0.100	1.83	1.84	0.09	0.94	1.83	1.84	0.08	0.94
0.200	1.76	1.76	0.07	0.92	1.76	1.76	0.06	0.93
0.300	1.66	1.66	0.06	0.93	1.66	1.66	0.05	0.93
0.400	1.56	1.57	0.05	0.95	1.56	1.57	0.05	0.94
0.500	1.49	1.49	0.05	0.95	1.49	1.49	0.04	0.95
0.600	1.43	1.43	0.04	0.96	1.43	1.43	0.04	0.96
0.700	1.38	1.39	0.04	0.97	1.38	1.39	0.04	0.96
0.800	1.35	1.35	0.04	0.98	1.35	1.35	0.04	0.98
0.900	1.32	1.32	0.04	0.98	1.32	1.32	0.03	0.98
1.000	1.29	1.29	0.04	0.98	1.29	1.29	0.03	0.98

An Epanechnikov kernel  $K(u) = \frac{3}{4}(1 - u^2)I_{|u| \leq 1}$  scaled by a bandwidth of  $h$ , was used for all  $\widehat{HDS}^{LC}(t)$  calculations. Since it has been shown that consistent estimators can be obtained by choosing bandwidth  $h_n = O(n^{-v})$  with  $1/4 < v < 1/2$  [Tian et al., 2005], we chose  $h = 0.26, 0.20, 0.19, 0.18$  for data with  $n = 200, 500, 800, 1000$ , respectively.

Tables 4.6, 4.7, 4.8 show the results of three cases under time dependent PH model with coefficients proportional to time. As is presented in Table 4.6, the  $\widehat{HDS}(t)$  shows



Table 4.4  $HDS(t)$  under the Cox PH model for right censored data with sample size = 500; exponential survival distribution:

Time	CR=15%				CR=25%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.83	0.12	0.96	1.82	1.83	0.13	0.97
0.050	1.84	1.85	0.12	0.95	1.84	1.84	0.13	0.97
0.100	1.83	1.85	0.12	0.93	1.83	1.84	0.12	0.95
0.200	1.76	1.77	0.09	0.92	1.76	1.76	0.09	0.89
0.300	1.66	1.66	0.07	0.92	1.66	1.66	0.08	0.90
0.400	1.56	1.57	0.06	0.93	1.56	1.57	0.07	0.92
0.500	1.49	1.49	0.06	0.94	1.49	1.49	0.06	0.95
0.600	1.43	1.44	0.06	0.95	1.43	1.43	0.06	0.96
0.700	1.38	1.39	0.05	0.96	1.38	1.39	0.06	0.96
0.800	1.35	1.35	0.05	0.96	1.35	1.35	0.05	0.96
0.900	1.32	1.32	0.05	0.96	1.32	1.32	0.05	0.96
1.000	1.29	1.29	0.05	0.95	1.29	1.30	0.05	0.96
Time	CR=30%				CR=60%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.83	0.13	0.96	1.82	1.84	0.16	0.94
0.050	1.84	1.85	0.13	0.96	1.84	1.86	0.16	0.94
0.100	1.83	1.85	0.12	0.95	1.83	1.85	0.15	0.93
0.200	1.76	1.77	0.10	0.92	1.76	1.77	0.12	0.91
0.300	1.66	1.66	0.08	0.93	1.66	1.67	0.10	0.93
0.400	1.56	1.57	0.07	0.94	1.56	1.57	0.09	0.94
0.500	1.49	1.49	0.06	0.97	1.49	1.50	0.09	0.94
0.600	1.43	1.43	0.06	0.96	1.43	1.45	0.10	0.94
0.700	1.38	1.39	0.06	0.97	1.38	1.43	0.10	0.94
0.800	1.35	1.35	0.06	0.97	1.35	1.41	0.10	0.91
0.900	1.32	1.32	0.06	0.97	1.32	1.41	0.10	0.88
1.000	1.29	1.30	0.06	0.97	1.29	1.41	0.10	0.81

little bias under exponential and Log Normal distribution and greater bias under Weibull and Log-Log distribution, but coverage probability varies between 0.86 and 0.98 across different survival distributions for the model with the coefficient proportional to time. From Table 4.7, the estimated  $HDS(t)$  has greater bias under small sample size, and the bias will be narrower as the sample size increases. Moreover, under smaller censoring rate, the bias will be less and the coverage probability will increase (see Table 4.8). The estimated  $HDS(t)$  at 1.000 is unavailable since there is

Table 4.5 Parameters for different survival distributions under time dependent PH model

Distribution	Baseline Hazard $\lambda_0(t)$	Parameters
Exponential	$k$	$k=1$
Weibull	$k^p pt^{p-1}$	$k=1, p=2$
Log-Logistic	$\frac{kp(kt)^{p-1}}{1+(kt)^p}$	$k=1, p=2$
Log Normal	$\frac{\partial}{\partial t} \left\{ -\log \left\{ 1 - \Phi \left( \frac{\log(t) - \mu}{\sigma} \right) \right\} \right\}$	$\mu=1, \sigma=2$

no observed time later than 1.000 under large censoring rate (60%).

The results of three cases under time dependent PH model with coefficients proportional to  $\log(t)$  are presented in Table 4.9, 4.10, 4.11. Since the model has a coefficient proportional to  $\log(t)$ , fewer failures will be observed between 0 and 0.5. The time points to be evaluated are moved to the period between 0.5 and 1.5. It is apparent that the bias grows as the evaluated time points away from 0.9 under each setting. The coverage probability is closer to 0.95 after 0.9 than before 0.9 (Table 4.9). The bias decreases as sample size increases or the right censoring rate decreases (Table 4.10 and Table 4.11). There seems no obvious differences in coverage probability under different sample size. However, it will not work well when right censoring rate is quite large (Table 4.11). Since the time range under the scenario with coefficients proportional to  $\log(t)$  is larger than the time range under the case with coefficients proportional to  $t$ , we conduct two additional simulations under exponential survival distribution, with a 500 sample size and a 25% censoring rate with a larger bandwidth  $h = 0.25$  and  $h = 0.30$ . Though the coverage probability before 0.80 is still far from 0.95, the bias and standard errors decrease when we use a larger bandwidth (see Table 4.12). Actually, the optimal bandwidth could be chosen by  $K$ -fold cross-validation (see real data example in section 5.2).

From Table 4.13, 4.14, 4.15, we can see that the results of the model with a coefficient proportional to  $t^2$  seem very similar to the results of the model with a

Table 4.6  $HDS^{LC}(t)$  under time dependent PH Model with coefficient proportional to  $t$ ; with censoring rate 25%; sample size = 500:

Time	Exponential				Weibull			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.32	0.11	0.89	1.28	1.44	0.21	0.85
0.050	1.28	1.31	0.10	0.89	1.28	1.41	0.19	0.86
0.100	1.28	1.30	0.09	0.89	1.29	1.37	0.15	0.88
0.200	1.28	1.29	0.08	0.93	1.29	1.34	0.11	0.90
0.300	1.28	1.29	0.09	0.94	1.31	1.33	0.10	0.92
0.400	1.29	1.30	0.09	0.94	1.32	1.33	0.09	0.93
0.500	1.30	1.31	0.10	0.94	1.33	1.33	0.09	0.94
0.600	1.32	1.33	0.12	0.93	1.33	1.34	0.09	0.96
0.700	1.34	1.35	0.14	0.94	1.34	1.34	0.11	0.95
0.800	1.36	1.38	0.17	0.94	1.34	1.34	0.13	0.96
0.900	1.37	1.41	0.22	0.95	1.33	1.34	0.16	0.96
1.000	1.38	1.46	0.29	0.93	1.31	1.34	0.21	0.95
Time	Log-Log				Log Normal			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.48	0.23	0.90	1.28	1.34	0.15	0.88
0.050	1.28	1.45	0.20	0.91	1.28	1.33	0.14	0.89
0.100	1.29	1.39	0.16	0.86	1.28	1.32	0.12	0.89
0.200	1.29	1.34	0.12	0.89	1.29	1.31	0.11	0.89
0.300	1.31	1.32	0.10	0.88	1.30	1.32	0.11	0.92
0.400	1.32	1.33	0.10	0.95	1.32	1.35	0.12	0.92
0.500	1.33	1.34	0.10	0.97	1.34	1.37	0.13	0.92
0.600	1.35	1.35	0.10	0.98	1.37	1.40	0.14	0.92
0.700	1.36	1.35	0.11	0.98	1.40	1.44	0.15	0.93
0.800	1.37	1.36	0.13	0.94	1.44	1.47	0.17	0.94
0.900	1.38	1.39	0.16	0.95	1.48	1.52	0.19	0.94
1.000	1.39	1.43	0.22	0.95	1.53	1.55	0.21	0.93

coefficient proportional to time. The only obvious difference is that the former one performs worse when the censoring rate is large (Table 4.15).

Table 4.7  $HDS^{LC}(t)$  under time dependent PH model with coefficient proportional to  $t$ ; with censoring rate 25%; exponential survival distribution

Time	n=200				n=500			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.35	0.17	0.90	1.28	1.32	0.11	0.89
0.050	1.28	1.34	0.16	0.91	1.28	1.31	0.10	0.89
0.100	1.28	1.32	0.14	0.91	1.28	1.30	0.09	0.89
0.200	1.28	1.30	0.12	0.92	1.28	1.29	0.08	0.93
0.300	1.28	1.30	0.12	0.93	1.28	1.29	0.09	0.94
0.400	1.29	1.31	0.13	0.94	1.29	1.30	0.09	0.94
0.500	1.30	1.33	0.15	0.94	1.30	1.31	0.10	0.94
0.600	1.32	1.35	0.17	0.93	1.32	1.33	0.12	0.93
0.700	1.34	1.38	0.20	0.94	1.34	1.35	0.14	0.94
0.800	1.36	1.42	0.25	0.93	1.36	1.38	0.17	0.94
0.900	1.37	1.46	0.33	0.94	1.37	1.41	0.22	0.95
1.000	1.38	1.53	0.44	0.93	1.38	1.46	0.29	0.93
Time	n=800				n=1000			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.31	0.09	0.89	1.28	1.30	0.08	0.91
0.050	1.28	1.30	0.08	0.90	1.28	1.30	0.07	0.91
0.100	1.28	1.29	0.07	0.90	1.28	1.29	0.06	0.90
0.200	1.28	1.28	0.07	0.92	1.28	1.28	0.06	0.91
0.300	1.28	1.29	0.07	0.92	1.28	1.29	0.06	0.92
0.400	1.29	1.29	0.07	0.94	1.29	1.29	0.07	0.92
0.500	1.30	1.31	0.08	0.94	1.30	1.31	0.08	0.93
0.600	1.32	1.33	0.10	0.93	1.32	1.32	0.09	0.94
0.700	1.34	1.34	0.11	0.94	1.34	1.34	0.10	0.95
0.800	1.36	1.37	0.14	0.93	1.36	1.36	0.12	0.94
0.900	1.37	1.39	0.17	0.94	1.37	1.38	0.15	0.92
1.000	1.38	1.42	0.22	0.95	1.38	1.40	0.20	0.93

### 4.3 COX PH MODEL WITH INTERVAL CENSORING

To evaluate the performance for interval censored data, we use the same survival model as used in the simulation under right censoring (section 4.1). For interval censored data, the exact event time is calculated to catch the interval between which we observe the event. To control the right censoring rate, we generated the intervals as follows. The number of observation times for each subject is equal to one plus a

Table 4.8  $HDS^{LC}(t)$  under time dependent PH model with coefficient proportional to  $t$ ; with sample size = 500; exponential survival distribution:

Time	CR=15%				CR=25%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.32	0.12	0.93	1.28	1.32	0.11	0.89
0.050	1.28	1.31	0.11	0.93	1.28	1.31	0.10	0.89
0.100	1.28	1.30	0.09	0.93	1.28	1.30	0.09	0.89
0.200	1.28	1.29	0.08	0.94	1.28	1.29	0.08	0.93
0.300	1.28	1.29	0.09	0.93	1.28	1.29	0.09	0.94
0.400	1.29	1.30	0.09	0.94	1.29	1.30	0.09	0.94
0.500	1.30	1.32	0.10	0.94	1.30	1.31	0.10	0.94
0.600	1.32	1.33	0.12	0.94	1.32	1.33	0.12	0.93
0.700	1.34	1.35	0.13	0.95	1.34	1.35	0.14	0.94
0.800	1.36	1.37	0.16	0.95	1.36	1.38	0.17	0.94
0.900	1.37	1.40	0.20	0.95	1.37	1.41	0.22	0.95
1.000	1.38	1.43	0.25	0.95	1.38	1.46	0.29	0.93
Time	CR=30%				CR=60%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.32	0.11	0.87	1.28	1.33	0.08	0.74
0.050	1.28	1.31	0.10	0.87	1.28	1.32	0.08	0.74
0.100	1.28	1.30	0.09	0.88	1.28	1.31	0.07	0.74
0.200	1.28	1.29	0.08	0.90	1.28	1.29	0.07	0.83
0.300	1.28	1.29	0.08	0.92	1.28	1.30	0.08	0.85
0.400	1.29	1.30	0.09	0.93	1.29	1.31	0.10	0.86
0.500	1.30	1.32	0.10	0.94	1.30	1.35	0.14	0.86
0.600	1.32	1.33	0.12	0.93	1.32	1.38	0.18	0.89
0.700	1.34	1.36	0.15	0.94	1.34	1.43	0.25	0.91
0.800	1.36	1.40	0.18	0.93	1.36	1.54	0.38	0.92
0.900	1.37	1.44	0.24	0.93	1.37	1.95	0.72	0.93
1.000	1.38	1.49	0.33	0.93	1.38			

random count which follows a Poisson distribution with mean  $\theta = 5$ . This ensures each subject has at least one visit, and the number of visits varies among subjects. The gap time between adjacent observations was generated based on an exponential distribution with mean  $\phi = 0.1$ . This combination of  $(\theta, \phi)$  let the right censoring rate under exponential distribution close to 25%. We use  $(\theta = 7, \phi = 1/9)$ ,  $(\theta = 7, \phi = 1/9)$ ,  $(\theta = 5, \phi = 1/14)$  to control the right censoring rate equal to 25% under Weibull, Log-Log, Log-Normal distribution, respectively. We also use  $(\theta =$

Table 4.9  $HDS^{LC}(t)$  under time dependent PH model with coefficient proportional to  $\log(t)$ ; with censoring rate 25%; sample size = 500:

Time	Exponential				Weibull			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.500	1.47	1.33	0.12	0.65	1.47	1.34	0.12	0.68
0.600	1.37	1.29	0.10	0.75	1.37	1.29	0.11	0.76
0.700	1.32	1.27	0.10	0.83	1.32	1.27	0.10	0.85
0.800	1.28	1.26	0.09	0.89	1.28	1.26	0.09	0.90
0.900	1.26	1.27	0.10	0.94	1.26	1.26	0.09	0.93
1.000	1.25	1.28	0.11	0.95	1.25	1.27	0.10	0.95
1.100	1.25	1.28	0.13	0.95	1.24	1.24	0.09	0.97
1.200	1.25	1.30	0.16	0.95	1.23	1.25	0.11	0.96
1.300	1.26	1.34	0.21	0.94	1.23	1.26	0.14	0.95
1.400	1.27	1.42	0.30	0.94	1.23	1.29	0.20	0.95
1.500	1.27	1.57	0.48	0.93	1.23	1.35	0.30	0.96
Time	Log-Log				Log Normal			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.500	1.47	1.38	0.15	0.75	1.48	1.39	0.18	0.72
0.600	1.38	1.32	0.12	0.79	1.39	1.34	0.15	0.80
0.700	1.32	1.29	0.11	0.84	1.33	1.32	0.14	0.86
0.800	1.28	1.28	0.10	0.90	1.30	1.31	0.13	0.90
0.900	1.26	1.28	0.10	0.92	1.28	1.31	0.13	0.92
1.000	1.26	1.28	0.10	0.94	1.27	1.33	0.13	0.93
1.100	1.25	1.27	0.10	0.97	1.28	1.34	0.13	0.96
1.200	1.26	1.28	0.11	0.98	1.28	1.36	0.14	0.96
1.300	1.26	1.30	0.13	0.97	1.29	1.38	0.15	0.97
1.400	1.27	1.33	0.16	0.97	1.31	1.40	0.16	0.97
1.500	1.27	1.36	0.21	0.97	1.32	1.43	0.18	0.98

7,  $\phi = 1/9$ ),  $(\theta = 4, \phi = 1/11)$ ,  $(\theta = 3, \phi = 1/28)$  to control the right censoring rate equal to 15%, 30%, and 60% under exponential distribution, respectively. Then, the observation times are calculated by the cumulative gap times. For the  $i^{th}$  subject, the observed interval  $(L_i, R_i)$  is two cumulative times between which the event time  $T_i$  lies. When  $T_i$  is less (greater) than the smallest (largest) observation time, define  $L_i$  ( $R_i$ ) as NA and let  $R_i$  ( $L_i$ ) equal to the smallest (greatest) observation time.

From Table 4.16, we can see that there is little bias in  $HDS(t)$  under different survival distributions over time. Moreover, the coverage probability is close to 95%

Table 4.10  $HDS^{LC}(t)$  under time dependent PH model with coefficient proportional to  $\log(t)$ ; with censoring rate 25%; exponential survival distribution:

Time	n=200				n=500			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.500	1.47	1.37	0.18	0.75	1.47	1.33	0.12	0.65
0.600	1.37	1.32	0.15	0.82	1.37	1.29	0.10	0.75
0.700	1.32	1.29	0.14	0.88	1.32	1.27	0.10	0.83
0.800	1.28	1.28	0.14	0.90	1.28	1.26	0.09	0.89
0.900	1.26	1.29	0.14	0.93	1.26	1.27	0.10	0.94
1.000	1.25	1.29	0.15	0.95	1.25	1.28	0.11	0.95
1.100	1.25	1.29	0.16	0.96	1.25	1.28	0.13	0.95
1.200	1.25	1.31	0.19	0.96	1.25	1.30	0.16	0.95
1.300	1.26	1.35	0.24	0.95	1.26	1.34	0.21	0.94
1.400	1.27	1.41	0.32	0.96	1.27	1.42	0.30	0.94
1.500	1.27	1.49	0.45	0.97	1.27	1.57	0.48	0.93
Time	n=800				n=1000			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.500	1.47	1.33	0.10	0.59	1.47	1.32	0.09	0.57
0.600	1.37	1.28	0.08	0.69	1.37	1.28	0.08	0.67
0.700	1.32	1.26	0.08	0.82	1.32	1.26	0.07	0.78
0.800	1.28	1.25	0.07	0.88	1.28	1.26	0.07	0.88
0.900	1.26	1.26	0.08	0.92	1.26	1.26	0.07	0.93
1.000	1.25	1.27	0.08	0.94	1.25	1.27	0.08	0.95
1.100	1.25	1.26	0.09	0.97	1.25	1.26	0.08	0.97
1.200	1.25	1.28	0.10	0.97	1.25	1.27	0.09	0.97
1.300	1.26	1.30	0.13	0.96	1.26	1.29	0.11	0.97
1.400	1.27	1.31	0.16	0.97	1.27	1.30	0.14	0.96
1.500	1.27	1.33	0.20	0.96	1.27	1.32	0.18	0.95

most of time except at two or three time points under each distribution. The bias decreases and standard error shrinks as sample size increases, while the coverage probability seems stable (Table 4.17). It performs well when right censoring rate is not too large and works less well after 0.5 with huge right censoring rate (Table 4.18).

Table 4.11  $HDS^{LC}(t)$  under time dependent PH model with coefficient proportional to  $\log(t)$ ; with sample size = 500; exponential survival distribution:

Time	CR=15%				CR=25%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.500	1.47	1.34	0.12	0.68	1.47	1.33	0.12	0.65
0.600	1.37	1.29	0.11	0.76	1.37	1.29	0.10	0.75
0.700	1.32	1.27	0.10	0.85	1.32	1.27	0.10	0.83
0.800	1.28	1.26	0.09	0.90	1.28	1.26	0.09	0.89
0.900	1.26	1.26	0.09	0.93	1.26	1.27	0.10	0.94
1.000	1.25	1.27	0.10	0.95	1.25	1.28	0.11	0.95
1.100	1.25	1.27	0.10	0.98	1.25	1.28	0.13	0.95
1.200	1.25	1.28	0.12	0.97	1.25	1.30	0.16	0.95
1.300	1.26	1.30	0.14	0.97	1.26	1.34	0.21	0.94
1.400	1.27	1.32	0.18	0.96	1.27	1.42	0.30	0.94
1.500	1.27	1.34	0.22	0.93	1.27	1.57	0.48	0.93
Time	CR=30%				CR=60%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.500	1.47	1.34	0.12	0.67	1.47	1.35	0.11	0.62
0.600	1.37	1.29	0.10	0.73	1.37	1.30	0.10	0.72
0.700	1.32	1.27	0.10	0.85	1.32	1.28	0.10	0.79
0.800	1.28	1.27	0.09	0.89	1.28	1.27	0.10	0.87
0.900	1.26	1.27	0.10	0.92	1.26	1.28	0.10	0.88
1.000	1.25	1.27	0.11	0.95	1.25	1.30	0.13	0.89
1.100	1.25	1.27	0.11	0.96	1.25	1.32	0.16	0.89
1.200	1.25	1.29	0.13	0.97	1.25	1.37	0.21	0.89
1.300	1.26	1.32	0.17	0.96	1.26	1.51	0.35	0.92
1.400	1.27	1.34	0.21	0.96	1.27	1.69	0.61	0.93
1.500	1.27	1.38	0.28	0.95	1.27			



Table 4.12  $HDS^{LC}(t)$  under time dependent PH model with coefficient proportional to  $\log(t)$ ; censoring rate 25%;  $n = 500$ ; exponential survival distribution:

Time	$HDS(t)$	$h = 0.25$			$h = 0.30$		
		$\widehat{HDS}(t)$	$SE$	$CP$	$\widehat{HDS}(t)$	$SE$	$CP$
0.500	1.47	1.33	0.11	0.64	1.33	0.10	0.60
0.600	1.37	1.29	0.09	0.74	1.29	0.08	0.71
0.700	1.32	1.27	0.09	0.83	1.27	0.08	0.81
0.800	1.28	1.26	0.08	0.89	1.26	0.08	0.88
0.900	1.26	1.26	0.09	0.92	1.26	0.08	0.94
1.000	1.25	1.27	0.09	0.94	1.26	0.08	0.96
1.100	1.25	1.28	0.10	0.96	1.27	0.09	0.96
1.200	1.25	1.29	0.12	0.97	1.28	0.11	0.97
1.300	1.26	1.31	0.14	0.97	1.29	0.13	0.96
1.400	1.27	1.33	0.18	0.97	1.31	0.16	0.97
1.500	1.27	1.35	0.23	0.95	1.34	0.20	0.96

Table 4.13  $HDS^{LC}(t)$  under time dependent PH model with coefficient proportional to  $t^2$ ; with censoring rate 25%; sample size = 500:

Time	Exponential				Weibull			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.32	0.12	0.89	1.28	1.48	0.24	0.84
0.050	1.28	1.31	0.11	0.89	1.28	1.43	0.20	0.84
0.100	1.27	1.30	0.10	0.89	1.28	1.38	0.16	0.84
0.200	1.26	1.29	0.09	0.90	1.28	1.32	0.12	0.85
0.300	1.26	1.28	0.09	0.91	1.28	1.30	0.10	0.88
0.400	1.26	1.28	0.10	0.92	1.27	1.30	0.10	0.90
0.500	1.26	1.28	0.10	0.93	1.28	1.30	0.09	0.93
0.600	1.27	1.29	0.11	0.93	1.29	1.30	0.09	0.94
0.700	1.30	1.31	0.12	0.93	1.30	1.31	0.09	0.95
0.800	1.34	1.35	0.14	0.93	1.33	1.32	0.10	0.95
0.900	1.39	1.39	0.16	0.93	1.36	1.34	0.12	0.93
1.000	1.45	1.44	0.20	0.93	1.38	1.34	0.16	0.94
	Log-Log				Log Normal			
0.025	1.28	1.48	0.25	0.84	1.28	1.35	0.16	0.88
0.050	1.28	1.44	0.21	0.83	1.28	1.34	0.15	0.88
0.100	1.28	1.39	0.17	0.85	1.28	1.32	0.13	0.88
0.200	1.28	1.33	0.13	0.87	1.27	1.31	0.12	0.90
0.300	1.28	1.30	0.11	0.88	1.27	1.31	0.13	0.91
0.400	1.28	1.30	0.10	0.91	1.27	1.32	0.13	0.90
0.500	1.28	1.30	0.10	0.91	1.28	1.33	0.14	0.89
0.600	1.29	1.31	0.10	0.93	1.30	1.35	0.14	0.91
0.700	1.31	1.33	0.10	0.94	1.34	1.39	0.15	0.92
0.800	1.35	1.36	0.11	0.95	1.39	1.44	0.16	0.94
0.900	1.40	1.39	0.13	0.95	1.47	1.51	0.17	0.95
1.000	1.46	1.43	0.16	0.93	1.57	1.61	0.19	0.95

Table 4.14  $HDS^{LC}(t)$  under time dependent PH model with coefficient proportional to  $t^2$ ; with censoring rate 25%; exponential survival distribution:

Time	n=200				n=500			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.35	0.18	0.90	1.28	1.32	0.12	0.89
0.050	1.28	1.34	0.17	0.90	1.28	1.31	0.11	0.89
0.100	1.27	1.32	0.15	0.90	1.27	1.30	0.10	0.89
0.200	1.26	1.30	0.13	0.92	1.26	1.29	0.09	0.90
0.300	1.26	1.29	0.13	0.91	1.26	1.28	0.09	0.91
0.400	1.26	1.29	0.14	0.91	1.26	1.28	0.10	0.92
0.500	1.26	1.30	0.15	0.93	1.26	1.28	0.10	0.93
0.600	1.27	1.32	0.16	0.93	1.27	1.29	0.11	0.93
0.700	1.30	1.34	0.17	0.93	1.30	1.31	0.12	0.93
0.800	1.34	1.38	0.20	0.93	1.34	1.35	0.14	0.93
0.900	1.39	1.42	0.24	0.93	1.39	1.39	0.16	0.93
1.000	1.45	1.47	0.30	0.93	1.45	1.44	0.20	0.93
Time	n=800				n=1000			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.30	0.09	0.91	1.28	1.31	0.09	0.88
0.050	1.28	1.30	0.09	0.91	1.28	1.30	0.08	0.88
0.100	1.27	1.29	0.08	0.91	1.27	1.29	0.07	0.87
0.200	1.26	1.27	0.07	0.90	1.26	1.28	0.07	0.90
0.300	1.26	1.27	0.07	0.89	1.26	1.27	0.07	0.92
0.400	1.26	1.27	0.08	0.93	1.26	1.27	0.07	0.94
0.500	1.26	1.27	0.08	0.93	1.26	1.27	0.07	0.93
0.600	1.27	1.28	0.09	0.93	1.27	1.29	0.08	0.94
0.700	1.30	1.31	0.09	0.92	1.30	1.31	0.09	0.93
0.800	1.34	1.34	0.11	0.94	1.34	1.35	0.10	0.94
0.900	1.39	1.39	0.13	0.94	1.39	1.38	0.12	0.94
1.000	1.45	1.44	0.16	0.93	1.45	1.42	0.15	0.93

Table 4.15  $HDS^{LC}(t)$  under time dependent PH model with coefficient proportional to  $t^2$ ; with sample size = 500; exponential survival distribution:

Time	CR=15%				CR=25%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.32	0.13	0.92	1.28	1.32	0.12	0.89
0.050	1.28	1.31	0.12	0.91	1.28	1.31	0.11	0.89
0.100	1.27	1.30	0.10	0.91	1.27	1.30	0.10	0.89
0.200	1.26	1.28	0.09	0.92	1.26	1.29	0.09	0.90
0.300	1.26	1.27	0.09	0.93	1.26	1.28	0.09	0.91
0.400	1.26	1.27	0.10	0.93	1.26	1.28	0.10	0.92
0.500	1.26	1.28	0.10	0.95	1.26	1.28	0.10	0.93
0.600	1.27	1.29	0.11	0.95	1.27	1.29	0.11	0.93
0.700	1.30	1.32	0.12	0.94	1.30	1.31	0.12	0.93
0.800	1.34	1.35	0.13	0.95	1.34	1.35	0.14	0.93
0.900	1.39	1.39	0.15	0.95	1.39	1.39	0.16	0.93
1.000	1.45	1.44	0.19	0.95	1.45	1.44	0.20	0.93
Time	CR=30%				CR=60%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.28	1.33	0.12	0.89	1.28	1.32	0.10	0.81
0.050	1.28	1.32	0.11	0.88	1.28	1.32	0.09	0.82
0.100	1.27	1.30	0.09	0.88	1.27	1.30	0.08	0.82
0.200	1.26	1.28	0.09	0.90	1.26	1.29	0.08	0.86
0.300	1.26	1.27	0.09	0.91	1.26	1.29	0.08	0.85
0.400	1.26	1.28	0.09	0.91	1.26	1.29	0.09	0.85
0.500	1.26	1.29	0.10	0.92	1.26	1.31	0.10	0.83
0.600	1.27	1.30	0.11	0.92	1.27	1.33	0.12	0.83
0.700	1.30	1.33	0.12	0.95	1.30	1.37	0.15	0.79
0.800	1.34	1.36	0.14	0.94	1.34	1.48	0.22	0.78
0.900	1.39	1.40	0.17	0.94	1.39	2.02	0.52	0.81
1.000	1.45	1.45	0.22	0.94	1.45	2.02	0.52	0.78

Table 4.16  $HDS(t)$  the Cox PH model for interval censored data with right censoring rate 25%; sample size = 500:

Time	Exponential				Weibull			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.84	0.14	0.96	1.79	1.83	0.13	0.95
0.050	1.84	1.86	0.14	0.96	1.80	1.83	0.13	0.95
0.100	1.83	1.86	0.12	0.94	1.81	1.84	0.14	0.95
0.200	1.76	1.78	0.07	0.80	1.83	1.87	0.15	0.95
0.300	1.66	1.67	0.06	0.85	1.84	1.87	0.13	0.91
0.400	1.56	1.57	0.07	0.93	1.79	1.82	0.08	0.80
0.500	1.49	1.50	0.07	0.96	1.71	1.72	0.07	0.80
0.600	1.43	1.44	0.07	0.98	1.60	1.61	0.07	0.90
0.700	1.38	1.39	0.06	0.98	1.50	1.50	0.07	0.95
0.800	1.35	1.35	0.06	0.97	1.41	1.42	0.07	0.98
0.900	1.32	1.32	0.05	0.96	1.34	1.35	0.06	0.97
1.000	1.29	1.29	0.05	0.96	1.29	1.30	0.05	0.96
Time	Log-Log				Log Normal			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.79	1.82	0.14	0.96	1.83	1.85	0.14	0.96
0.050	1.80	1.82	0.14	0.96	1.84	1.87	0.13	0.95
0.100	1.81	1.83	0.14	0.96	1.79	1.81	0.08	0.83
0.200	1.83	1.86	0.15	0.96	1.64	1.66	0.06	0.85
0.300	1.84	1.86	0.13	0.93	1.53	1.54	0.07	0.95
0.400	1.80	1.82	0.09	0.82	1.46	1.47	0.07	0.97
0.500	1.72	1.74	0.07	0.80	1.41	1.42	0.07	0.97
0.600	1.63	1.63	0.07	0.89	1.37	1.38	0.06	0.97
0.700	1.53	1.54	0.08	0.95	1.34	1.35	0.06	0.97
0.800	1.45	1.46	0.07	0.97	1.32	1.32	0.06	0.95
0.900	1.39	1.40	0.07	0.98	1.30	1.31	0.06	0.93
1.000	1.34	1.35	0.06	0.98	1.28	1.30	0.06	0.92

Table 4.17  $HDS(t)$  under the Cox PH model for interval censored data with right censoring rate 25%; exponential survival distribution:

Time	n=200				n=500			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.88	0.23	0.97	1.82	1.84	0.14	0.96
0.050	1.84	1.90	0.24	0.97	1.84	1.86	0.14	0.96
0.100	1.83	1.90	0.20	0.92	1.83	1.86	0.12	0.94
0.200	1.76	1.80	0.12	0.81	1.76	1.78	0.07	0.80
0.300	1.66	1.69	0.11	0.86	1.66	1.67	0.06	0.85
0.400	1.56	1.58	0.12	0.91	1.56	1.57	0.07	0.93
0.500	1.49	1.51	0.12	0.94	1.49	1.50	0.07	0.96
0.600	1.43	1.45	0.11	0.94	1.43	1.44	0.07	0.98
0.700	1.38	1.40	0.10	0.95	1.38	1.39	0.06	0.98
0.800	1.35	1.36	0.10	0.95	1.35	1.35	0.06	0.97
0.900	1.32	1.33	0.09	0.94	1.32	1.32	0.05	0.96
1.000	1.29	1.30	0.08	0.93	1.29	1.29	0.05	0.96
Time	n=800				n=1000			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.85	0.11	0.96	1.82	1.84	0.10	0.96
0.050	1.84	1.87	0.11	0.96	1.84	1.86	0.10	0.96
0.100	1.83	1.87	0.10	0.91	1.83	1.86	0.08	0.93
0.200	1.76	1.78	0.05	0.76	1.76	1.77	0.05	0.78
0.300	1.66	1.67	0.05	0.83	1.66	1.67	0.04	0.84
0.400	1.56	1.57	0.06	0.93	1.56	1.57	0.05	0.94
0.500	1.49	1.50	0.06	0.96	1.49	1.50	0.05	0.97
0.600	1.43	1.44	0.05	0.96	1.43	1.44	0.05	0.98
0.700	1.38	1.39	0.05	0.97	1.38	1.39	0.04	0.98
0.800	1.35	1.35	0.05	0.97	1.35	1.35	0.04	0.98
0.900	1.32	1.32	0.04	0.96	1.32	1.32	0.04	0.98
1.000	1.29	1.29	0.04	0.96	1.29	1.29	0.04	0.97

Table 4.18  $HDS(t)$  under the Cox PH model for interval censored data with sample size = 500; exponential survival distribution:

Time	CR=15%				CR=25%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.85	0.13	0.96	1.82	1.84	0.14	0.96
0.050	1.84	1.87	0.14	0.96	1.84	1.86	0.14	0.96
0.100	1.83	1.86	0.12	0.92	1.83	1.86	0.12	0.94
0.200	1.76	1.78	0.07	0.81	1.76	1.78	0.07	0.80
0.300	1.66	1.67	0.06	0.85	1.66	1.67	0.06	0.85
0.400	1.56	1.57	0.07	0.93	1.56	1.57	0.07	0.93
0.500	1.49	1.50	0.07	0.97	1.49	1.50	0.07	0.96
0.600	1.43	1.44	0.06	0.97	1.43	1.44	0.07	0.98
0.700	1.38	1.39	0.06	0.97	1.38	1.39	0.06	0.98
0.800	1.35	1.35	0.05	0.97	1.35	1.35	0.06	0.97
0.900	1.32	1.32	0.05	0.96	1.32	1.32	0.05	0.96
1.000	1.29	1.30	0.04	0.95	1.29	1.29	0.05	0.96
Time	CR=30%				CR=60%			
	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$	$HDS(t)$	$\widehat{HDS}(t)$	$SE$	$CP$
0.025	1.82	1.84	0.14	0.95	1.82	1.84	0.17	0.96
0.050	1.84	1.86	0.14	0.96	1.84	1.86	0.17	0.96
0.100	1.83	1.85	0.12	0.92	1.83	1.86	0.14	0.91
0.200	1.76	1.77	0.07	0.78	1.76	1.78	0.09	0.80
0.300	1.66	1.66	0.07	0.85	1.66	1.67	0.10	0.88
0.400	1.56	1.57	0.08	0.92	1.56	1.58	0.12	0.91
0.500	1.49	1.49	0.08	0.96	1.49	1.54	0.13	0.87
0.600	1.43	1.43	0.07	0.98	1.43	1.53	0.12	0.78
0.700	1.38	1.39	0.07	0.98	1.38	1.53	0.12	0.68
0.800	1.35	1.35	0.06	0.97	1.35	1.53	0.12	0.58
0.900	1.32	1.32	0.06	0.95	1.32	1.53	0.12	0.50
1.000	1.29	1.29	0.06	0.93	1.29	1.53	0.12	0.42

## CHAPTER 5

### REAL DATA ANALYSIS

South Carolina ranked the eighth highest rates of HIV diagnoses in the United States in 2017 [for Disease Control and Prevention]. A large number of studies investigated the association between HIV-related diseases and potential predictors, such as viral load level (VLD). We are particularly interested in time to first time of viral load suppression [Yehia et al., 2015]. Antiretroviral therapy (ART) is recommended for everyone who has HIV. It helps patients living with HIV (PLWHs) live longer, healthier lives and reduces the risk of HIV transmission. PLWHs are suggested to start ART as soon as possible. Once HIV patients are linked to care, initiating ART is another challenging task in HIV prevention [Palella et al., 2003].

We apply the  $HDS(t)$  and  $HDS(t)$  ratio to the Health Sciences South Carolina (HSSC) data to investigate the discrimination among HIV suppression and the adherence to treatment. In section 5.1, time to first time viral suppression are considered. The  $HDS(t)$  and  $HDS(t)$  ratio under the Cox PH model and time dependent PH model are calculated. In section 5.2, initiating ART data are used. In both sections, we test the hypothesis that the main predictor improves the discrimination performance by using the methods in chapter 3.

#### 5.1 RIGHT CENSORED HSSC DATA

We use the days from the date of first diagnosed HIV to the date of first suppression ( $VLD < 200$  copies/mL) as the outcome, and number of years of retention in care (careny) as the main predictor. We also include age at baseline, nadir CD4 cell



count (lcd4), and  $\log(\text{VLD})$  at baseline as the risk factors into the survival model. There are 1051 PLWHs in the suppression data set, for whom 498 achieved success in suppression (47.38%).

The characteristic of this sample is shown in Table 5.1. Survival probability

Table 5.1 Characteristics of suppression data

Variable	Mean	Std Dev	Minimum	Maximum
age	42.9	13.2	14.0	81.0
careny	2.0	1.6	0.0	5.0
lcd4	295.9	233.2	1.0	1746.0
log	8.5	2.9	3.7	15.9

of suppression is plotted using a nonparametric method, the Kaplan-Meier method (Kaplan and Meier, 1958). Obviously, from Figure 5.1 we see that survival probability dramatically decreases at two time points, once around 150 days and again around 2200 days (nearly 6 years). Before estimating  $HDS(t)$ , we check the PH assumption

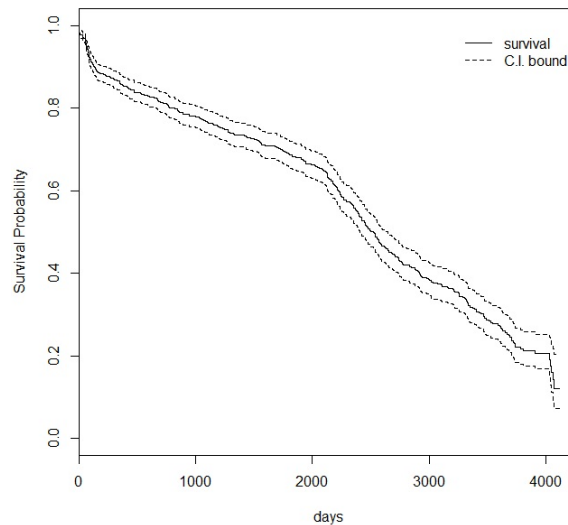


Figure 5.1 Kaplan-Meier curve for suppression data

using cumulative sums of martingale residuals [Lin et al., 1993]. From Table 5.2, we can see that only the p-value for nadir CD4 is greater than the significance level (0.05); thus, we have evidence to reject the hypothesis and conclude that the PH assumption is violated for years in care, age, and log(VLD). Since the PH assumption

Table 5.2 Supremum Test for Proportionals Hazards Assumption

Value	Maximum Absolute	Replications	Seed	Pr > MaxAbsVal
careny	3.2289	1000	1208787608	<.0001
lcd4	1.5259	1000	1208787608	0.0650
age	2.2458	1000	1208787608	0.0010
log	4.4485	1000	1208787608	<.0001

is not satisfied, we relax the proportional hazards to time dependent hazards to get an estimator for localized  $HDS^{LC}(t)$ . For the analysis of  $HDS^{LC}(t)$ , we choose the "optimal" bandwidth by using the  $K$ -fold cross-validation method, which is commonly used for nonparametric estimation [Hoover et al., 1998][Tian et al., 2005]. The data is split into  $K=13$  equal-sized parts. Given a certain bandwidth  $h$ , we estimate smoothed coefficients based on the sample data excluding the  $k^{th}$  part,  $k = 1, 2, \dots, K$ . We then calculate the "prediction error",  $PE_k(h)$ , by using the estimates to predict the  $k^{th}$  part of the data. If  $t_{(1)}, \dots, t_{(D_k)}$  are the ordered failure times in the  $k^{th}$  part of the sample,

$$PE_k(h) = - \sum_{i=1}^{D_k} \{ \hat{\beta}'_h(t_{(i)}) M_{(i)} - \log( \sum_{l \in R(t_{(i)})} \exp(\hat{\beta}'_h(t_{(i)}) M_l) ) \}$$

where  $M$  is the  $k^{th}$  part of the dataset,  $\hat{\beta}_h(t)$  is an estimate vector of the remaining  $K - 1$  parts at time  $t$ , and  $R_{(t_{(i)})}$  is the risk set at time  $t_{(i)}$ . The optimal bandwidth is such that total prediction error,  $PE(h) = \sum_{k=1}^{13} PE_k(h)$ , is minimized. We randomize the original dataset 50 times, each time repeating this cross-validation process at bandwidth  $h = 500, 510, \dots, 750$ . Most of the optimal  $h$  falls into the interval [710, 740]

(see Table 5.3). Since for almost half of the times,  $PE$  at  $h = 740$  is the smallest, we choose  $h = 740$  as the bandwidth for localized  $HDS^{LC}(t)$ .

Table 5.3 Choices of Optimal Bandwidth

h	520	700	710	730	740	750
frequency	1	1	9	14	24	1

The  $HDS(t)$  ratio is calculated by comparing the  $HDS(t)$  of the model with the number of years the patients receive retention in care and the  $HDS(t)$  of the model without the number of years retention in care. [figure]singlelinecheck=on

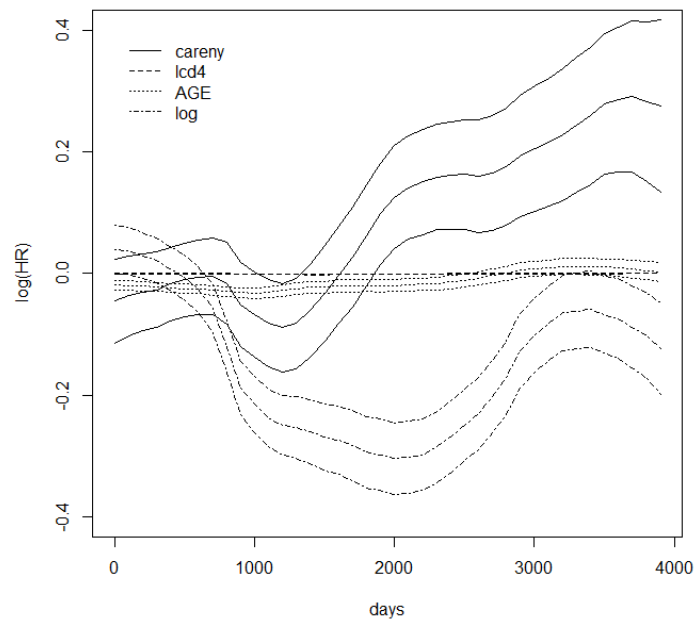


Figure 5.2 Time-varying Log(HR) Estimates for suppression data

The estimated coefficients under time dependent PH model as time changes is shown in Figure 5.2. We see that under time dependent PH model, after 1800 days the effect of number of years in care on suppression changes from negative to posi-

tive, while most of time baseline age and  $\log(\text{VLD})$  affect the hazard of suppression negatively. There is nearly no effect of nadir CD4 over time.

The estimates of  $HDS^{LC}(t)$  under time dependent PH model are presented in Figure 5.3. We see that  $HDS^{LC}(t)$  is significantly greater than one over time, which indicates that the discrimination performance of time dependent PH model including these four predictors is quite effective. We can also see that including these four predictors, the discriminatory ability of this time dependent PH model is quite strong around 2100 days.

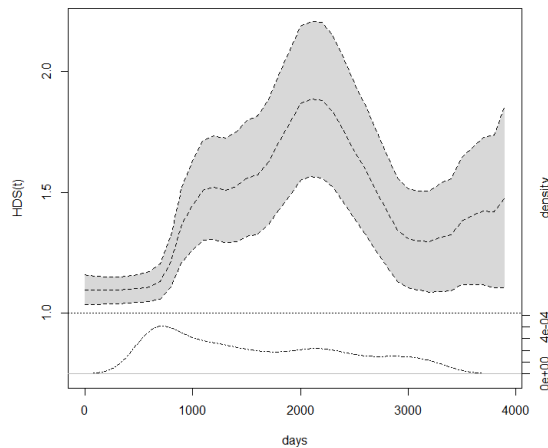


Figure 5.3  $HDS^{LC}(t)$  for suppression data

However, from Figure 5.4, the lower bound of 95% confidence interval for  $HDS^{LC}(t)$  ratio varies around one, indicating that the association of improvement in predicting survival and years of retention in care is not significant most of the time.

## 5.2 INTERVAL CENSORED HSSC DATA

The interval in the data set Data.art covers the exact time of initiating ART. We analyze ART data to evaluate the discrimination performance of the Cox PH model with initial CD4 cell count ( $\text{incd4}$ ),  $\log(\text{VLD})$ , and baseline age, and we test the

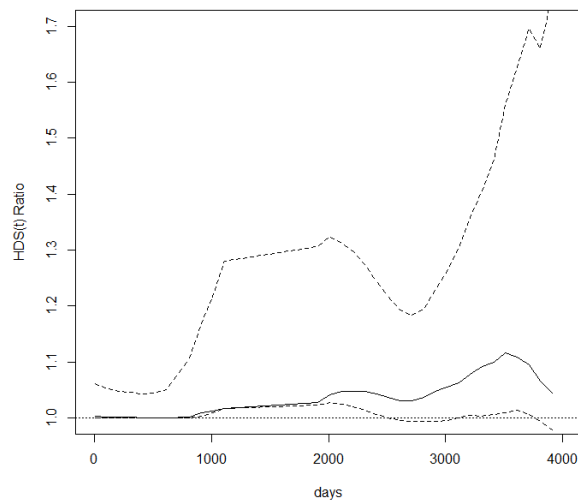


Figure 5.4  $HDS^{LC}(t)$  Ratio for suppression data

hypothesis that taking initial CD4 cell count into the Cox PH model improves the discrimination of initiating ART. 1007 PLWHs are included, among which 85, 396, 526 patients are left censored, interval censored, and right censored, respectively.

The characteristic of ART data is shown in Table 5.4. From Table 5.5, we see

Table 5.4 Characteristics of ART data

Variable	Mean	Std Dev	Minimum	Maximum
incd4	385.9	280.9	1.0	1812.0
log	8.5	3.0	3.7	15.9
age	42.9	13.1	14.0	81.0

that the 95% confidence interval (C.I.) of estimates for CD4 cell count is lower than zero. We estimate that the log(HR) of receiving ART is around -0.0009 when one cell count increase in CD4 with other predictors fixed, and we have evidence to conclude that lower initial CD4, which indicates the subject is in a worse health condition, increases the probability of initiating ART earlier. Analyzing  $HDS(t)$  under the Cox

Table 5.5  $\log(\text{HR})$  and 95% C.I. for estimates using ART data

Variable	estimate	Std Dev	LCI	UCI
incd4	-0.0009	0.0002	-0.0013	-0.0005
log	0.0222	0.0174	-0.0119	0.0563
age	0.0034	0.0037	-0.0038	0.0106

PH model by using the proposed methods, we see that the 95% confidence interval for  $HDS(t)$  covers one over time. We don't have evidence to reject the null hypothesis that there is no discriminatory performance of the Cox PH model with these three predictors (Figure 5.5). The confidence interval for  $HDS(t)$  ratio of including initial

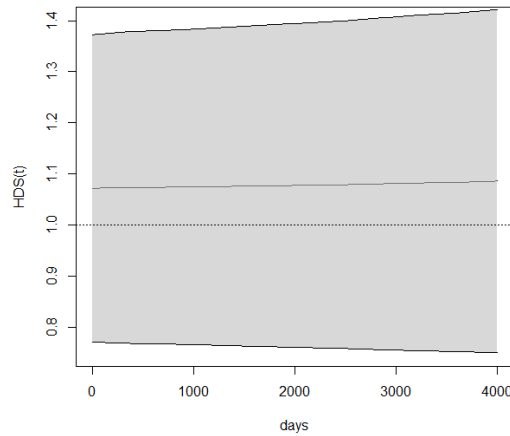


Figure 5.5  $HDS(t)$  for ART data

CD4 or not is presented in Figure 5.6. We see that the lower bound is less than one, which indicates that including initial CD4 does not improve the Cox PH model.

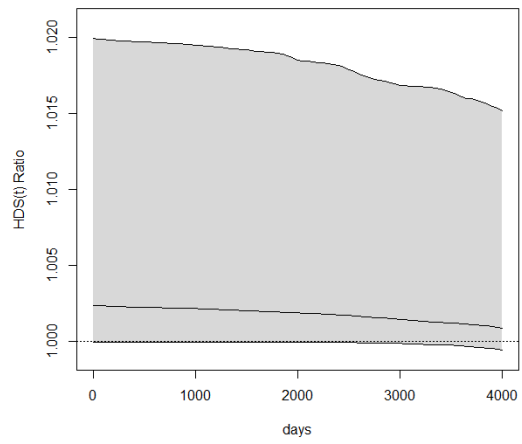


Figure 5.6  $HDS(t)$  Ratio for ART data

## CHAPTER 6

### CONCLUSIONS AND FUTURE STUDY

$HDS(t)$  is a time-varying measure which generalizes the discrimination slope to evaluate the discrimination performance for survival models proposed by Liang and Heagerty [2017]. It has been shown that the estimation for  $HDS(t)$  under the Cox PH model and time dependent  $HDS(t)$  performs effectively with exponential distribution for right censored data [Liang and Heagerty, 2017]. Based on the referenced paper, we evaluated the performance of estimation for  $HDS(t)$  under the Cox PH model and time dependent PH model with different survival distributions, different sample sizes, and different right censoring rates.

Firstly, according to the results of simulation studies,  $HDS(t)$  performs well based on data with a large sample size and a small right censoring rate. When PH assumption is violated, an alternative method is to analyze data using time dependent PH model. Using a proper bandwidth, which can be chosen by  $k$ -fold cross-validation,  $HDS^{LC}(t)$  under time dependent PH model also works well.

Secondly, we also extended the application of  $HDS(t)$  to interval censored data. It has been shown that the discrimination performance of a survival model, given interval censored data, can be evaluated by  $HDS(t)$ . Although at some time points the bias of estimated  $HDS(t)$  is not small, most of the time it works adequately with a large sample size and a small right censoring rate.

In addition, we can test whether a main predictor of interest improves the survival model through the estimation of  $HDS(t)$  ratio. Since the standard error for  $HDS(t)$  ratio is quite complicated, we can make inferences by using a bootstrap method to



construct a confidence interval of  $HDS(t)$  ratio.

Currently, there are many other survival models to predict the survival probability, such as the generalized odds rate model (GOR model), which is more generalized and flexible [Dabrowska and Doksum, 1988]. Some avenues for future study include the extension of the  $HDS(t)$  under the GOR model so that the measure of discrimination performance for survival will be more flexible. We also can relax  $HDS(t)$  to the localized  $HDS(t)$  for interval censored data so that we can evaluate the discrimination performance of survival models for interval censored data in a safe and robust way. Furthermore, when new data comes, we could better explore suppression and initiating ART for PLWHs.

## BIBLIOGRAPHY

- Norman E Breslow. Contribution to discussion of paper by dr cox. *J. Roy. Statist. Soc., Ser. B*, 34:216–217, 1972.
- Zongwu Cai and Yanqing Sun. Local linear estimation for time-dependent coefficients in cox’s regression models. *Scandinavian Journal of Statistics*, 30(1):93–111, 2003.
- Lloyd E Chambless, Christopher P Cummiskey, and Gang Cui. Several methods to assess improvement in risk prediction models: extension to survival analysis. *Statistics in medicine*, 30(1):22–38, 2011.
- David R Cox. Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2):187–202, 1972.
- Dorota M Dabrowska and Kjell A Doksum. Estimation and testing in a two-sample generalized odds-rate model. *Journal of the American Statistical Association*, 83(403):744–749, 1988.
- Centers for Disease Control and Prevention. HIV in the United States by Region, year = 2017, url = <https://www.cdc.gov/hiv/statistics/overview/geographicdistribution.html>, urldate = November 27, 2018.
- Donald R Hoover, John A Rice, Colin O Wu, and Li-Ping Yang. Nonparametric smoothing estimates of time-varying coefficient models with longitudinal data. *Biometrika*, 85(4):809–822, 1998.
- C Jason Liang and Patrick J Heagerty. A risk-based measure of time-varying prognostic discrimination for survival models. *Biometrics*, 73(3):725–734, 2017.
- Danyu Y Lin, Lee-Jen Wei, and Zhiliang Ying. Checking the cox model with cumulative sums of martingale-based residuals. *Biometrika*, 80(3):557–572, 1993.

- Thomas A Louis. Finding the observed information matrix when using the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 44(2): 226–233, 1982.
- Frank J Palella, Maria Deloria-Knoll, Joan S Chmiel, Anne C Moorman, Kathleen C Wood, Alan E Greenberg, and Scott D Holmberg. Survival benefit of initiating antiretroviral therapy in hiv-infected persons in different cd4+ cell strata. *Annals of internal medicine*, 138(8):620–626, 2003.
- Lu Tian, David Zucker, and LJ Wei. On the cox model with time-varying regression coefficients. *Journal of the American statistical Association*, 100(469):172–183, 2005.
- Anastasios A Tsiatis et al. A large sample study of cox’s regression model. *The Annals of Statistics*, 9(1):93–108, 1981.
- Aad W van der Vaart, Jon A Wellner Wellner, et al. Empirical processes indexed by estimated functions. In *Asymptotics: particles, processes and inverse problems*, pages 234–252. Institute of Mathematical Statistics, 2007.
- Lianming Wang, Christopher S McMahan, Michael G Hudgens, and Zaina P Qureshi. A flexible, computationally efficient method for fitting the proportional hazards model to interval-censored data. *Biometrics*, 72(1):222–231, 2016.
- J Frank Yates. External correspondence: Decompositions of the mean probability score. *Organizational Behavior and Human Performance*, 30(1):132–156, 1982.
- Baligh R Yehia, Alisa J Stephens-Shields, John A Fleishman, Stephen A Berry, Allison L Agwu, Joshua P Metlay, Richard D Moore, W Christopher Mathews, Ank Nijhawan, Richard Rutstein, et al. The hiv care continuum: changes over time in retention in care and viral suppression. *PloS one*, 10(6):e0129376, 2015.
- Jie Zhou, Jiajia Zhang, and Wenbin Lu. An expectation maximization algorithm for fitting the generalized odds-rate model to interval censored data. *Statistics in medicine*, 36(7):1157–1171, 2017.