

# **Exploring Mental Health Related Emergency Department Visits: Frequency of Recurrence and Risk Factors**

by

Fei Wang

Project Submitted in Partial Fulfillment  
of the Requirements for the Degree of

Master of Science

in the

Department of Statistics and Actuarial Science  
Faculty of Science

© Fei Wang 2014

SIMON FRASER UNIVERSITY

Summer 2014

All rights reserved.

However, in accordance with the *Copyright Act of Canada*, this work may be reproduced without authorization under the conditions for "Fair Dealing." Therefore, limited reproduction of this work for the purposes of private study, research, criticism, review and news reporting is likely to be in accordance with the law, particularly if cited appropriately.

## APPROVAL

**Name:** Fei Wang  
**Degree:** Master of Science  
**Title of Project:** Exploring Mental Health Related Emergency Department Visits: Frequency of Recurrence and Risk Factors

**Examining Committee:** Dr. Tim B. Swartz  
Professor, SFU  
Chair

---

Dr. X. Joan Hu  
Professor, SFU  
Senior Supervisor

---

Dr. Q. Michelle Zhou  
Assistant Professor, SFU  
Supervisor

---

Dr. John J. Spinelli  
Adjunct Professor, SFU  
Professor, School of Population and Public Health, UBC  
Internal Examiner

**Date Approved:** August 20, 2014

## Partial Copyright Licence



The author, whose copyright is declared on the title page of this work, has granted to Simon Fraser University the non-exclusive, royalty-free right to include a digital copy of this thesis, project or extended essay[s] and associated supplemental files ("Work") (title[s] below) in Summit, the Institutional Research Repository at SFU. SFU may also make copies of the Work for purposes of a scholarly or research nature; for users of the SFU Library; or in response to a request from another library, or educational institution, on SFU's own behalf or for one of its users. Distribution may be in any form.

The author has further agreed that SFU may keep more than one copy of the Work for purposes of back-up and security; and that SFU may, without changing the content, translate, if technically possible, the Work to any medium or format for the purpose of preserving the Work and facilitating the exercise of SFU's rights under this licence.

It is understood that copying, publication, or public performance of the Work for commercial purposes shall not be allowed without the author's written permission.

While granting the above uses to SFU, the author retains copyright ownership and moral rights in the Work, and may deal with the copyright in the Work in any way consistent with the terms of this licence, including the right to change the Work for subsequent purposes, including editing and publishing the Work in whole or in part, and licensing the content to other parties as the author may desire.

The author represents and warrants that he/she has the right to grant the rights contained in this licence and that the Work does not, to the best of the author's knowledge, infringe upon anyone's copyright. The author has obtained written copyright permission, where required, for the use of any third-party copyrighted material contained in the Work. The author represents and warrants that the Work is his/her own original work and that he/she has not previously assigned or relinquished the rights conferred in this licence.

Simon Fraser University Library  
Burnaby, British Columbia, Canada

revised Fall 2013

# Abstract

This thesis project aims to provide insights into pediatric mental health care and help to improve its current practice. We explore records of mental health related emergency department visits from children and youth. The data are extracted from the provincial health administrative data systems of Alberta. We start with a descriptive data analysis, and then adopt the counting process framework to conduct statistical inference. A generalized (stratified) Cox regression model and a renewal process model are considered. We evaluate the frequency and identify important risk factors with various model specifications. We also account for the gaps of the visit process due to hospitalization. The project presents the estimates of the model parameters via likelihood and partial likelihood approaches. Robust estimates and the non-parametric bootstrap estimates for the standard errors of the parameter estimators are obtained in addition to the likelihood based standard error estimates. We summarize the analysis and outline a few problems for future investigation in the final chapter.

**Keywords** Counting process · Parametric/Semi-parametric/Non-parametric estimation procedure · Renewal process model · Stratified Cox regression model

*To my beloved parents!*

*"It is our choices, Harry, that show what we truly are, far more than our abilities."*

— *Albus Dumbledore*

HARRY POTTER AND THE CHAMBER OF SECRETS, 1999

# Acknowledgments

First of all, I would like to thank my supervisor, Dr. Joan Hu. She guided me with patience and encouragement, and helped me build up my confidence in reaching high standards, throughout the two years of my master program at SFU. I am forever grateful for having the opportunity to study and work under her supervision. Without her time devoted to it, this thesis would not have been completed.

I would also like to express my appreciation to the members of my thesis committee, Dr. Michelle Zhou and Dr. John Spinelli for their insightful comments, wonderful suggestions, and valuable time.

Special thanks to Dr. Rhonda Rosychuk for the opportunity to access to the administrative health data, and to conduct the data analysis in this project. Many thanks go to all the members in Rosychuk Biostatistics Laboratory, for their kindnesses, personal and technical helps. The data utilized in this project is provided by Alberta Health.

It brings my pleasure to record my thanks to the whole Department of Actuarial Science and Statistics, who has generously provided me with such an excellent atmosphere for studying statistics and doing research. I voice my particular appreciation to my graduate fellows Sherry Chen, Michael Grosskopf, Bobby Han, Kunasekaran Nirmalkanna, Nate Payne, Werjindra Premarathna, Pulindu Ratnasekera, Jerold Smith, Biljana Stojkova, Elena Szefer, Huijing Wang, Vicky Weng, Yi Xiong, Annie Yu, Sabrina Zhang, and who have been so kind to me.

Finally, I am deeply indebted to my family for their unconditional love, care, support and encouragement.

# Disclaimer

This study is based in part on data provided by Alberta Health. The interpretation and conclusions contained herein are those of the researchers and do not necessarily represent the views of the Government of Alberta. Neither the Government nor Alberta Health and Government express any opinion in relation to this study.



# Contents

<b>Approval</b>	<b>ii</b>
<b>Partial Copyright License</b>	<b>iii</b>
<b>Abstract</b>	<b>iv</b>
<b>Dedication</b>	<b>v</b>
<b>Quotation</b>	<b>vi</b>
<b>Acknowledgments</b>	<b>vii</b>
<b>Disclaimer</b>	<b>viii</b>
<b>Contents</b>	<b>ix</b>
<b>List of Tables</b>	<b>xii</b>
<b>List of Figures</b>	<b>xiv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background and Motivation . . . . .	1
1.2 General Framework . . . . .	3
1.3 Outline . . . . .	4

<b>2 Alberta Pediatric Mental Health Care (PMHC) Data</b>	<b>5</b>
2.1 Introduction of the PMHC Data . . . . .	5
2.1.1 Data Format . . . . .	6
2.1.2 Formulation of PMHC Data . . . . .	8
2.1.3 Potential Risk Factors . . . . .	9
2.2 Descriptive Analysis . . . . .	10
2.2.1 Summary of Descriptive Statistics . . . . .	10
2.2.2 Non-parametric Estimates of Cumulative Intensity Function . . . . .	14
2.3 Analysis under Parametric Model . . . . .	16
<b>3 Stratified Cox Regression Models</b>	<b>19</b>
3.1 Statistical Modeling . . . . .	19
3.2 Estimation Procedures . . . . .	20
3.2.1 Estimating Regression Parameters . . . . .	20
3.2.2 Estimating Baseline Intensity Functions . . . . .	22
3.3 Analysis Results . . . . .	23
3.3.1 With Time-independent Regression Coefficients . . . . .	23
3.3.2 With Time-dependent Regression Coefficients . . . . .	30
3.3.3 Main Effects together with Two Factor Interactions I . . . . .	31
3.3.4 Main Effects together with Two Factor Interactions II . . . . .	32
<b>4 Extended Renewal Process Model</b>	<b>35</b>
4.1 Statistical Modeling . . . . .	36
4.2 Estimation Procedures . . . . .	36
4.2.1 Estimation with a Semi-parametric Model . . . . .	37
4.2.2 Estimation with a Parametric Model . . . . .	38
4.3 Analysis Results . . . . .	39
4.4 Comparison with the Stratified Cox Regression Model . . . . .	43
<b>5 Final Remarks</b>	<b>44</b>
5.1 Summary . . . . .	44

5.2 Future Investigation . . . . .	45
<b>Bibliography</b>	<b>47</b>
<b>Appendix A Estimates of the Cumulative Intensity Functions</b>	<b>49</b>
<b>Appendix B Non-parametric Bootstrap Estimates of Standard Errors</b>	<b>55</b>
<b>Appendix C Estimation Formula</b>	<b>60</b>

# List of Tables

2.1	Data format of baseline information. . . . .	6
2.2	Data format of EDMH information. . . . .	7
2.3	Data format of hospitalization information. . . . .	7
2.4	Numbers of EDMH visits and distinct children and youth by fiscal year. . . . .	10
2.5	Characteristics of children and youth, according to the observations of the four risk factors at the index initial EDMH visit (Total $n = 27947$ ). . . . .	10
2.6	Summary of risk factors by the number of recurrent EDMH visits (%). . . . .	12
2.7	MLEs of the regression coefficients in Model 2. . . . .	17
3.1	MPLEs of the regression coefficients and robust estimates for the standard errors in Model 3a/3b/3c with time-independent covariate effects. . . . .	26
3.2	MPLEs of the regression coefficients and robust estimates for the standard errors in Model 3a/3b/3c with time-dependent covariate effects. . . . .	27
3.3	MPLEs of the regression coefficients and robust estimates for the standard errors in Model 3a/3b/3c with all pairs of two factor interactions and time-independent covariate effects. . . . .	28
3.4	MPLEs of the regression coefficients and robust estimates for the standard errors in Model 3a/3b/3c with all pairs of two factor interactions and time-dependent covariate effects. . . . .	29
3.5	MPLEs of the regression coefficients and the estimates of the standard errors under the fitted AG model considering two factor interactions and time-independent covariate effects. . . . .	32

3.6	MPLEs of the regression coefficients and the estimates of the standard errors under the fitted AG model considering two factor interactions and time-dependent covariate effects. . . . .	34
4.1	MLEs of the regression coefficients in Model 4a. . . . .	39
4.2	MLEs of the regression coefficients in Model 4b. . . . .	41
B.1	MPLEs of the regression coefficients and nonparametric bootstrap estimates for the standard errors in Model 3a/3b/3c with time-independent covariate effects. . . . .	56
B.2	MPLEs of the regression coefficients and nonparametric bootstrap estimates for the standard errors in Model 3a/3b/3c with time-dependent covariate effects. . . . .	57
B.3	MPLEs of the regression coefficients and nonparametric bootstrap estimates for the standard errors in Model 3a/3b/3c with all pairs of two factor interactions and time-independent covariate effects. . . . .	58
B.4	MPLEs of the regression coefficients and nonparametric bootstrap estimates for the standard errors in Model 3a/3b/3c with all pairs of two factor interactions and time-dependent covariate effects. . . . .	59

# List of Figures

2.1	ED mental health visits by age group and sex. . . . .	11
2.2	Distributions of pSES by the order of EDMH visit per subject had observed. . . . .	13
2.3	Distributions of region by the order of EDMH visit per subject had observed. . . . .	13
2.4	The log-transformed generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration. . . . .	15
2.5	The log-transformed estimates of the cumulative intensity functions against log of event times. . . . .	18
4.1	The estimates of the cumulative intensity function for an individual who had EDMH visit at time (A) $T_1 = 1000, T_2 = 2000, T_3 = 2500$ , and (B) $T_1 = 500, T_2 = 1000, T_3 = 2000$ . . . . .	40
4.2	The estimates of the cumulative rates of risk to the next EDMH visit since the index initial EDMH visit under Model 4a and Model 4b. The lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively. . . . .	42
A.1	The Generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration. . . . .	50
A.2	The Generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration, together with the Breslow estimates under Model 3c (the model considering time-independent covariate effects) in Table 3.1. The dashed lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively. . . . .	51

A.3	The Generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration, together with the Breslow estimates under Model 3c (the model considering time-dependent covariate effects) in Table 3.2. The dashed lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively. . . . .	52
A.4	The Generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration, together with the Breslow estimates under Model 3c (the two factor interaction model with time-independent covariate effects) in Table 3.5. The dashed lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively. . . . .	53
A.5	The Generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration, together with the Breslow estimates under Model 3c (the two factor interaction model with time-dependent covariate effects) in Table 3.6. The dashed lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively. . . . .	54

# Chapter 1

## Introduction

### 1.1 Background and Motivation

Emergency departments (ED) are often the first point of contact with the mental health care system for children and youth, and are considered as a safety net for the lack of inpatient and outpatient mental health services (Newton et al. 2011). The research team led by Drs. Amanda Newton and Rhonda Rosychuk has observed a significant heterogeneity in mental health presentations to ED in Alberta, and a high degree of repeated ED use. Their additional findings include "health system factors impact patient outcomes", a lack of community-based care available to children and youth, and an on-going need for pediatric mental health services; see Newton et al. (2010).

Newton et al. (2011) present a surveillance report about the analysis on the available data of the ED visits for mental health (EDMH visits) provided by Alberta Health. This study was designed to assist health care planners in recommending policies and allocating resources for children's mental health care. The report shows that, for example, more females than males presented for emergency mental health care. It is unclear, however, whether sex is a risk factor or simply reflects the overall pattern of EDMH visits, since similar pattern is observed in the study population. To better understand the need and to improve emergency mental health services for children and youth, the research team attempts to conduct comprehensive analyses related to mental health presentations to Alberta ED. One of the team's specific objectives is to evaluate the frequency of children and



youth's EDMH visits and to identify the risk factors.

When the frequency of EDMH visit over time is of primary interest, one may address the corresponding problem by conducting an analysis of recurrent events, using the well-developed methods for recurrent event analysis. For example, Andersen and Gill (1982) assume recurrent events follow a Poisson process, the well-known AG model, and propose the corresponding estimation procedures. Prentice et al. (1981) consider a stratified proportional intensity function to model event processes. Much of this later development has taken place within the general framework of counting process. Anderson et al. (1993) present theory and statistical methods for counting processes in an authoritative way. Cook and Lawless (2007) provide extensive examples and reviews.

Hospitalization records of Alberta children and youth are also available in the database maintained by Alberta Health. A patient should not be at risk for an EDMH visit during the time period when this patient is admitted to hospital. It is desirable to address the issue of observation gaps, or time periods not at risk to EDMH visits, due to hospitalization. Otherwise, those patients who have long hospital stay will likely be classified as ones with low risk to have the next EDMH visit, which may result in biased inference on the effects of associated risk factors. Recently, many notable studies analyze recurrent events which have a duration associated with them. Cook and Lawless (2007 Chap 6), for example, consider an alternating two-state process model, which defines an "active" and an "inactive" state, and models two types of recurrent events corresponding to two types of transitions between the two states at the same time. Another example is that Hu et al. (2011) extend the well-known models for recurrent events given in Prentice et al. (1981) and Anderson and Gill (1982) to a generalized Cox regression model and adjust the risk set at a time point to accommodate event duration. The approach in Hu et al. (2011) focuses on the conditional intensity function of recurrent event, and does not need to model the event duration. This project adopts the statistical method of Hu et al. (2011) to evaluate the frequency of EDMH visits, while adjusting for the gaps of the visit processes due to hospitalization.

The available data for the analysis in this project only contain information of children and youth who had EDMH visits within a fixed study period. This raises another issue: the first EDMH visit observed within this observation window may not be the first EDMH visit of the whole life of one subject. This motivates the consideration of a renewal process model, where a subject is assumed

to restore to the same physical state right after each EDMH visit. There are methods developed to conduct the gap time analysis in renewal process model. See Breslow (1972), Kalbfleisch and Prentice (2002), Cook and Lawless (2007, Chap 4). Following Hu et al. (2011), we extend the renewal process model to account for the non-negligible duration of hospitalization.

This project aims to assess the frequency of EDMH visits and to identify the associated important risk factors. As the first attempt to achieve this goal, we focus on recurrence of EDMH visits since subjects' initial EDMH visits recorded during the period from April 1 2002 to March 31 2011, based on the ED records and other information of Alberta residents aged younger than 18 years old. We adapt the counting process framework to conduct data analysis with a stratified Cox regression model and a renewal process model. Following Hu et al. (2011), models are extended to address possible gaps when the individuals are not at the risk to such visits due to hospitalization.

## 1.2 General Framework

Let  $0 < T_1 < T_2, \dots$  be the times of a subject's first, second, ... EDMH visits since his/her initial visit recorded in the database (referred to as index initial EDMH visit).

Define  $N(t) = \sum_{j=1}^{\infty} I(T_j \leq t)$ , where  $N(t)$  represents the cumulative number of the recurrent EDMH visits, up to time  $t > 0$ , that the subject in the study has since his/her index initial EDMH visit.

Note that  $N(\cdot)$  is a right-continuous counting process.

Denote further the covariate vector of the subject at time  $t$  by  $Z(t)$ . The history information of the subject at time  $t$  is denoted by  $\mathcal{H}(t) = \mathcal{N}(t) \cup \mathcal{Z}(t)$ , where  $\mathcal{N}(t) = \{N(s) : 0 \leq s < t\}$ , and  $\mathcal{Z}(t) = \{Z(s) : 0 \leq s < t\}$ .

Let  $\lambda(t|\mathcal{H}(t))$  be the conditional intensity function of the counting process  $N(\cdot)$ . Following Hu et al. (2011), we consider the generalized Cox regression model

$$\begin{aligned} \lambda(t|\mathcal{H}(t)) &= \lim_{\Delta t \rightarrow 0^+} P\{N(t) - N(t - \Delta t) = 1 | \mathcal{H}(t)\} / \Delta t \\ &= \lambda_0\{t; \mathcal{H}(t)\} \exp\{\beta(t; \mathcal{H}(t))' Z(t)\}, \quad t > 0, \end{aligned} \quad (1)$$

where  $\lambda_0\{t; \mathcal{H}(t)\}$  is an arbitrary baseline intensity function, and  $\beta(t; \mathcal{H}(t))$  is a function up to finite dimensional parameters. The dimension of  $\beta(t; \mathcal{H}(t))$  is the same as  $Z(t)$ . We refer to this model as Model 1. Specifying the components in Model 1 into different forms, one can explore various special cases of the generalized Cox regression model. For example, assume the baseline intensity function and the regression coefficients are independent of the history information, that is,  $\lambda_0\{t; \mathcal{H}(t)\} = \lambda_0(t)$  and  $\beta(t; \mathcal{H}(t)) = \beta$ . Model 1 reduces to a Poisson process model, the well-known AG model considered in Anderson and Gill (1982).

Consider a group of independent individuals with observations on the counting process starting from 0 to a censoring time. We formulate the available information associated with a subject as the right-censored realization of  $\{N(t) : t \geq 0\}$  together with  $Z(t)$  at EDMH visits. Assume non-informative censoring conditional on  $Z(\cdot)$ . Our primary goal is to estimate the baseline intensity function  $\lambda_0\{t; \mathcal{H}(t)\}$  and the regression coefficients  $\beta(t; \mathcal{H}(t))$  based on the right-censored counting process data. The estimator of  $\beta(t; \mathcal{H}(t))$  can be used to identify a risk factor by conducting a test on whether the effect of a covariate is significant. We are also interested in estimating the standard errors of the estimators, and constructing confidence intervals for the intensity function.

### 1.3 Outline

The rest of the project is organized as follows. Chapter 2 introduces Alberta Pediatric Mental Health Care (PMHC) data and provides a descriptive analysis. A parametric version of Model 1 is assumed, and the fit of it based on the PMHC data is represented in this chapter. Chapter 3 introduces one class of specifications of Model 1, a stratified Cox regression model, and reports the estimation results. Chapter 4 presents another class of specifications of Model 1, an extended renewal process model, along with inference procedures and the associated asymptotic results. Conclusions and some remarks for further investigation are given in Chapter 5.

## **Chapter 2**

# **Alberta Pediatric Mental Health Care (PMHC) Data**

According to Leitch (2007), eighty percent of mental illness begin in childhood. This paper also states that fifteen percent of Canadian children and youth live with a mental illness, but only 1 in 6 receives timely specialized services. With limited options for immediate mental health care, families go to emergency departments. In order to better understand the risk factors associated with variation in the number of EDMH visits for children and youth, the project analyzes a set of Pediatric Mental Health Care (PMHC) data extracted from the database of Alberta Health. This chapter describes the PMHC data.

### **2.1 Introduction of the PMHC Data**

The PMHC data sources are the four population-based administrative databases: Ambulatory Care Classification System (ACCS), Population Registry File (PRF), Physicians Claims File (PCF), and Hospitalizations Discharge Database (HDD). The ACCS database provides Alberta EDMH visit information. The demographic and geographic data are from the PRF database. PCF tracks physician (follow-up) visits, while HDD contains hospitalization data.

Individual level data that exist in the data sources from April 1, 2002 to March 31, 2011 were

extracted. Since some individuals may have EDMH visits before April 1, 2002, this project focuses on the time period for each individual starting from the day when the individual has his/her first EDMH visit after April 1, 2002, referred to as the index initial EDMH visit; while the time period for each individual is ended at either the day of his/her 18th birthday or March 1, 2011, whichever is earlier.

The subjects of this project are Alberta residents who had at least one EDMH visit during the study period from April 1, 2002 and March 31, 2011, and were younger than 18 years of age at the time of the EDMH visit. An Alberta resident is defined as an individual who is registered in the Alberta Health Care Insurance Plan(AHCIP). In Alberta, a child or youth is defined as someone under the age of 18 years. The resulting study group consists of 27947 subjects.

### 2.1.1 Data Format

We view one EDMH visit as an event of interest. In order to conduct the analysis, we construct three data matrices, based on the available PMHC data. The first data matrix is referred to as *Baseline Information*. It consists of all the characteristic information of the study subjects at the index initial EDMH visits. Each subject has one row in the data matrix. Table 2.1 shows the data format.

Table 2.1: Data format of baseline information.

ID	Start.Date	Age	pSES	Sex	Region
xxxxx1	yyyy-mm-dd	x	x	x	x
xxxxx2	yyyy-mm-dd	x	x	x	x
xxxxx3	yyyy-mm-dd	x	x	x	x
		.....			

The variables in Table 2.1, from the left to right, correspond to (i) the unique identification number of the study subject, (ii) the start date of the index initial EDMH visit, (iii) age at the index initial EDMH visit (in years), (iv) the proxy of the social-economic status at the index initial EDMH visit, (v) sex, and (vi) residential region at the index initial EDMH visit, respectively.

The second data matrix contains all the EDMH visit records of the study subjects. Each row represents one EDMH visit record. For those subjects who had more than one EDMH visit, there are multiple rows in the data matrix. We refer to this data matrix as *EDMH Information*. The data

format is shown in Table 2.2.

Table 2.2: Data format of EDMH information.

ID	Start.Date	Age	pSES	Sex	Region
xxxxx1	yyyy-mm-dd	x	x	x	x
xxxxx1	yyyy-mm-dd	x	x	x	x
xxxxx1	yyyy-mm-dd	x	x	x	x
xxxxx2	yyyy-mm-dd	x	x	x	x
xxxxx3	yyyy-mm-dd	x	x	x	x
xxxxx3	yyyy-mm-dd	x	x	x	x
.....					

The variables in Table 2.2, from the left to right, are (i) the unique identification number of the study subject, (ii) the start date of the EDMH visit, (iii) age at the EDMH visit (in years), (iv) the proxy of the social-economic status at the EDMH visit, (v) sex, and (vi) residential region at the EDMH visit, respectively.

The third data matrix includes the hospitalization records, referred to as *Hospitalization Information*. Some patients may have more than one hospitalization record, while some may not have any during their follow-up periods within this observation window. Therefore, not all subjects have information in this data matrix. Table 2.3 presents the format of this data matrix.

Table 2.3: Data format of hospitalization information.

ID	Start.Date	End.Date
xxxxx2	yyyy-mm-dd	yyyy-mm-dd
xxxxx4	yyyy-mm-dd	yyyy-mm-dd
xxxxx4	yyyy-mm-dd	yyyy-mm-dd
.....		

The variables in Table 2.3, from left to right, correspond to (i) the unique identification number of the study subject, (ii) the admission date of the hospitalization, and (iii) the discharge date of the hospitalization, respectively.

We set the **Start.Date** in the *Baseline Information* data set as the time origin for each subject. Then all other time-dependent variables, such as **Start.Date** and **End.Date** in the *Hospitalization Information* data set, are shifted to the time since the index initial EDMH visit correspondingly.

### 2.1.2 Formulation of PMHC Data

Consider days as the time scale. Let  $Sday$  and  $Eday$  be April 1, 2002 and March 31, 2011, the starting and ending days of the PMHC data extraction, respectively.

Denote the realization of  $\{N(\cdot), Z(\cdot)\}$  associated with subject  $i$  in the study by  $\{N_i(\cdot), Z_i(\cdot)\}$  for  $i = 1, \dots, n$ . For subject  $i$ , denote the times of EDMH visits since the index initial EDMH visit by  $T_{i1}, T_{i2}, \dots$ , where  $0 < T_{i1} < T_{i2} < \dots$ . Suppose the observation on  $N_i(\cdot)$  is subject to the right-censoring with the censoring time  $C_i$ . According to the data collection mechanism,

$$C_i = \min(B_i + 18 \times 365 - A_{i0}, Eday - A_{i0}) > 0,$$

where  $A_{i0}$  and  $B_i$  are the calendar times of the index initial EDMH visit and subject  $i$ 's birthday. Let  $Y_i^C(t) = I(C_i \geq t)$  be the censoring indicator. In addition, the PMHC database includes observations on  $Z_i(\cdot)$  at times  $0 = t_{i0} < t_{i1} < \dots < t_{iK_i}$ , where  $t_{iK_i} \leq C_i < t_{i,(K_i+1)}$ , and  $K_i$  is the number of EDMH visits of subject  $i$  observed after the index initial EDMH visit and within the observation window.

Since there is no information of the individual birthday, we utilize the age of subject at EDMH visits to estimate the birthday. Denote  $A_{ik}$  as the calendar times of the  $k$ th EDMH visit for subject  $i$ ,  $Age_{ik}$  is the corresponding age recorded. Then  $A_{ik} - (Age_{ik} + 1) \times 365 < B_i \leq A_{ik} - Age_{ik} \times 365$ . We approximate  $B_i$  by the average value of  $A_{ik} - (Age_{ik} + 1/2) \times 365$  over all  $k$  for subject  $i$ .

Let  $\mathcal{R}(t)$  be the risk set at time  $t$ . That is,  $\mathcal{R}(t)$  contains all those subjects who are at risk to EDMH visits at time  $t$ . To address the issue of non-negligible duration of hospitalization, we define the risk set at time  $t$  as a set of subjects who are not being hospitalized and have not been censored at time  $t$ . For subject  $i$ , denote the times to the admission of hospitalization by  $V_{i1}^A, V_{i2}^A, \dots$ , and the times to the discharge of hospitalization by  $V_{i1}^D, V_{i2}^D, \dots$ , where  $0 < V_{i1}^A < V_{i1}^D < V_{i2}^A < V_{i2}^D < \dots$ . Then we have

$$\mathcal{R}(t) = \{i : t \notin \cup_k (V_{ik}^A, V_{ik}^D); t \leq C_i\}.$$

Let  $Y_i^R(t)$  be the at-risk indicator of subject  $i$ , where  $Y_i^R(t) = 1$  or  $0$  if subject  $i$  is at risk at time  $t$  or not. When the risk set is not adjusted for the gaps of the visit processes due to hospitalization, the at-risk indicator is the same as the censoring indicator.

We make the following two assumptions about the data.

- The study subjects are independent with each other.
- $N_i(\cdot)$  and  $C_i$  are independent conditional on  $Z_i(\cdot)$ .

### 2.1.3 Potential Risk Factors

Several potential risk factors (covariates) associated with the frequency of EDMH visits are identified by the PMHC data. They include demographic factors such as sex, age, pSES over time, geographic factors such as health region of residence, residential region (urban/rural), and diagnostic factors including triage level. We choose to focus on four potential risk factors, which are listed in the following:

- Age: the individual's age (in years) at the index initial EDMH visit, with values from 0 to 17.
- Sex: the indicator of male.
- pSES: the proxy of the social-economic status.

There are four categories: Others/Registration without Subsidy (O), Aboriginal Groups (A), Government Sponsored Program (S), Welfare (W). Most subjects who revisited ED for mental illness had stable pSES. That is, their social-economic status did not change over time. Thus we utilized pSES at the index initial EDMH visit. Since most subjects were with pSES of O (65.3%), we combined the other three groups (A, S, and W) when conducting the analysis. See Table 2.5. This variable is coded as an indicator of O in this analysis.

- Region: residential region with two categories, urban and rural. Region is not time-varying overtime very much either. We consider the effect of region at the index initial EDMH visit. This variable is coded as an indicator of urban in this analysis.



## 2.2 Descriptive Analysis

### 2.2.1 Summary of Descriptive Statistics

The data analyzed in this project include 27947 subjects and their relevant information over the time windows from 0 to  $C_i$ . The maximum length of the time windows is 3064 days. There are a total of 41159 EDMH visits, with an average of 1.47 visits per subject (max 52). Most subjects (20871, 74.7%) had only one EDMH visit during their observation period, while 25.3% of them had multiple EDMH visits. Table 2.4 represents the numbers of EDMH visits and the totals of the associated subjects across the fiscal year. The yearly numbers of EDMH visits and individuals visiting ED for mental illnesses increased from 4278 in 2003 to 4849 in 2011, and from 3438 in 2003 to 3773 in 2011, respectively. A fiscal year is defined as a one year period from April 1 of the previous year to March 31 of the current year. For example, fiscal year 2003 is from April 1 2002 to March 31 2003.

Table 2.4: Numbers of EDMH visits and distinct children and youth by fiscal year.

	Fiscal Year (yyyy)									Total
	2003	2004	2005	2006	2007	2008	2009	2010	2011	
EDMH Visits	4278	4258	4472	4629	4661	4584	4849	4579	4849	41159
Children/youth	3438	3443	3643	3724	3715	3663	3915	3684	3773	27947

Table 2.5: Characteristics of children and youth, according to the observations of the four risk factors at the index initial EDMH visit (Total  $n = 27947$ ).

Age		Sex		pSES		Region	
Category	Number	Category	Number	Category	Number	Category	Number
Pre-school (0 ~ 5)	579	Male	12095	O	18260	Urban	21293
Elementary (6 ~ 13)	6063	Female	15852	A	3785	Rural	6654
Teenager (14 ~ 17)	21305			S	4098		
				W	1804		

Table 2.5 summarizes the characteristics of children and youth who had EDMH visit during the period, according to the observations of the four risk factors at the index initial EDMH visit. The average age is 14.43 years old. More than half (76.2%, 21305) of the study subjects were teenagers

of age 14 ~ 17. The subject group had slightly more females than males (56.7% females, 43.3% males). The distribution of the pSES is 65.3%, 13.5%, 14.7%, and 6.5% of the four groups, O, A, S, and W, respectively. There are more children and youth from urban (76.2%).

Of the total 41159 EDMH visits, 58.7% (24160 visits) were made by females, while 41.3% (16999 visits ) were made by males. Visits made by females exceeded visits by males overall. However, Figure 2.1 shows that younger males tend to have slightly higher frequency of EDMH visits than their peer females. While females of age 14 ~ 17 contribute more EDMH visits than males in the same age group. This suggests that the association of EDMH visit frequency with sex differs overtime.

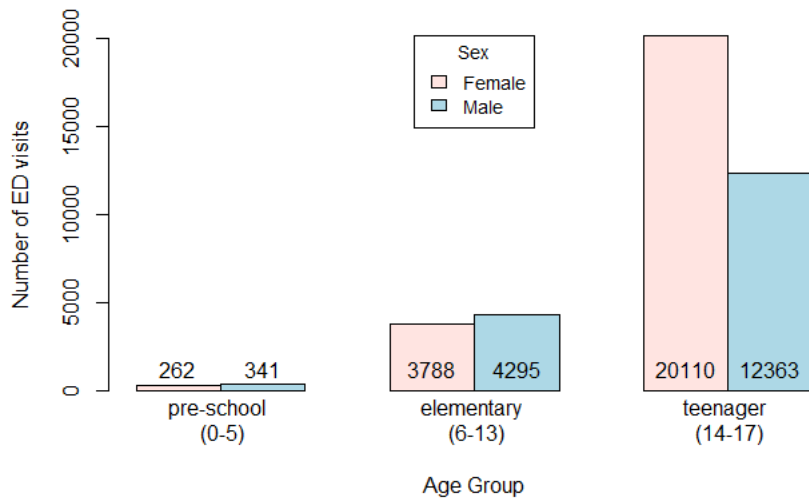


Figure 2.1: ED mental health visits by age group and sex.

We summarize the four risk factors according to the number of recurrent EDMH visits that one subject observed, by the relative frequency of each category among its corresponding risk factor. Table 2.6 reports the summary. The column with zero number of recurrent EDMH visit represents the summary at the index initial EDMH visit. The distributions of age, sex, and residential region among the study subjects with different number of recurrent EDMH visit are quite similar; while a slightly decreased pattern of pSES with category O (vs ASW), is observed with the increase of the number of recurrent EDMH visit.

Table 2.6: Summary of risk factors by the number of recurrent EDMH visits (%).

Variable	Category	Number of recurrent EDMH visit						
		0	1	2	3	4	5	>5
Age	Pre-school (0 ~ 5)	0.03	0.01	0.00	0.00	0.01	0.00	0.00
	Elementary (6 ~ 13)	0.20	0.23	0.30	0.32	0.34	0.36	0.44
	Teenager (14 ~ 17)	0.77	0.76	0.69	0.68	0.65	0.64	0.56
Sex	Male	0.45	0.41	0.37	0.37	0.38	0.40	0.28
	Female	0.55	0.59	0.63	0.63	0.62	0.60	0.72
pSES	O	0.68	0.61	0.58	0.53	0.47	0.47	0.50
	ASW	0.32	0.39	0.42	0.47	0.53	0.53	0.50
Region	Urban	0.76	0.75	0.77	0.80	0.76	0.77	0.76
	Rural	0.24	0.25	0.23	0.20	0.24	0.23	0.24

It appears that the status of pSES and residential region may change overtime, since some subjects who had multiple EDMH visits show different pSES and/or region at their recurrent EDMH visits. We checked the distributions of the two risk factors at the index initial EDMH visit, the first EDMH visit after the index initial visit, the second, and so on. Figure 2.2 shows the distribution of pSES at the index initial EDMH visit, and first to fifth EDMH visit after the index initial EDMH visit. Figure 2.3 shows the distribution of residential region. The number at the the bottom of each bar stands for the proportion of individuals with pSES in the category of O in Figure 2.2, and the proportion of individuals living in urban in Figure 2.3, respectively. From the two figures, we can see the pattern of each bar varies slightly, especially for pSES, which may indicate that some patients did change their pSES status and/or residential region during the study period. However, since we have a quite large study population, the proportion of patient whose pSES or region varied overtime is very small, 0.1% and 0.01% for pSES and region, respectively. Thus, the two risk factors are assumed as time-independent variables in this study.

From April 1 2002 and March 31 2011, 12528 hospitalization records were reported from 7868 subjects who had EDMH visits within this observation window. Among them, there is an average of 1.56 hospitalization records per individual (median 1, max 28). The average length of hospitalization per individual is 16.81 days (median 5, IQR 2 to 16, max 843 days). We found from the analysis in Chapter 2.3 that estimates from a parametric model with risk set adjusted for the duration of the hospitalization are quite similar to those with risk set not-adjusted for the duration of the hospitalization. That is because only 28% of the study subjects had observed hospitalization

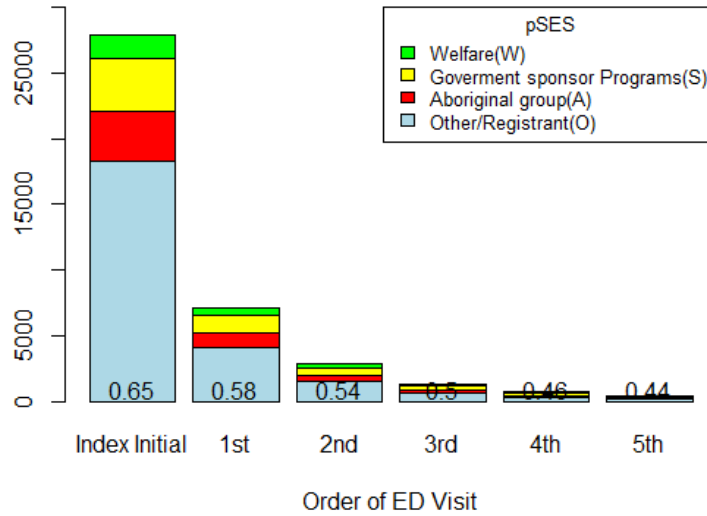


Figure 2.2: Distributions of pSES by the order of EDMH visit per subject had observed.

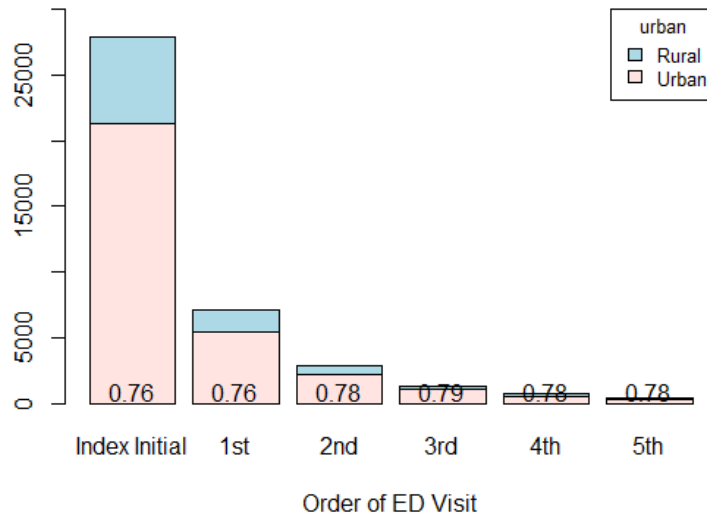


Figure 2.3: Distributions of region by the order of EDMH visit per subject had observed.

records.

## 2.2.2 Non-parametric Estimates of Cumulative Intensity Function

We group subjects according to age at the index initial EDMH visit (pre-school 0 ~ 5, elementary school 6 ~ 13, and teenager (14 ~ 17), sex (male and female), pSES at the index initial EDMH visit (O and ASW), and region at the index initial EDMH visit (urban and rural). For each of the different subgroups, we evaluate the generalized Nelson-Aalen estimator for the cumulative intensity function with adjustment for the hospital duration proposed by Hu et.al (2011), using the following formula:

$$\hat{\Lambda}_0(t) = \int_0^t \frac{\sum_{i=1}^n Y_i^C(u) dN_i(u)}{\sum_{l=1}^n Y_l^C(u) Y_l^R(u)} = \sum_{\{t_d: t_d \leq t\}} \frac{\sum_{i=1}^n Y_i^C(t_d) I(dN_i(t_d) = 1)}{\sum_{l=1}^n Y_l^C(t_d) Y_l^R(t_d)}, \quad t > 0,$$

where  $t_d$  is the distinct event time.

Figure 2.4 shows the log-transformed generalized Nelson-Aalen estimates of the cumulative intensity functions. The vertical distance between each pair of curves represents the corresponding covariate effect by controlling other covariates. The non-parametric estimates are roughly parallel with each other. It indicates that the cumulative intensity functions of the study group are likely proportional according to the four risk factors. There were no observed recurrent EDMH visits until around 6 years (2200 days) after the index initial EDMH visit, for the subgroup of subjects who are females of age 0 ~ 5 with pSES of O in rural area. Therefore, the log-transformed estimate of cumulative intensity function within the corresponding time period (from 0 to 6 years) is negative infinite and not shown in Figure 2.4(A).

Figure A.1 in Appendix A represents the generalized Nelson-Aalen estimates of the cumulative intensity functions. According to Figure A.1, males are less likely to have EDMH visits than females; the subjects with pSES in the category of O tend to have lower risk of having EDMH visits, compared to those in the pSES category of A, S, and W; the subjects living in rural area are at lower risk of recurrent EDMH visits than those living in urban area.

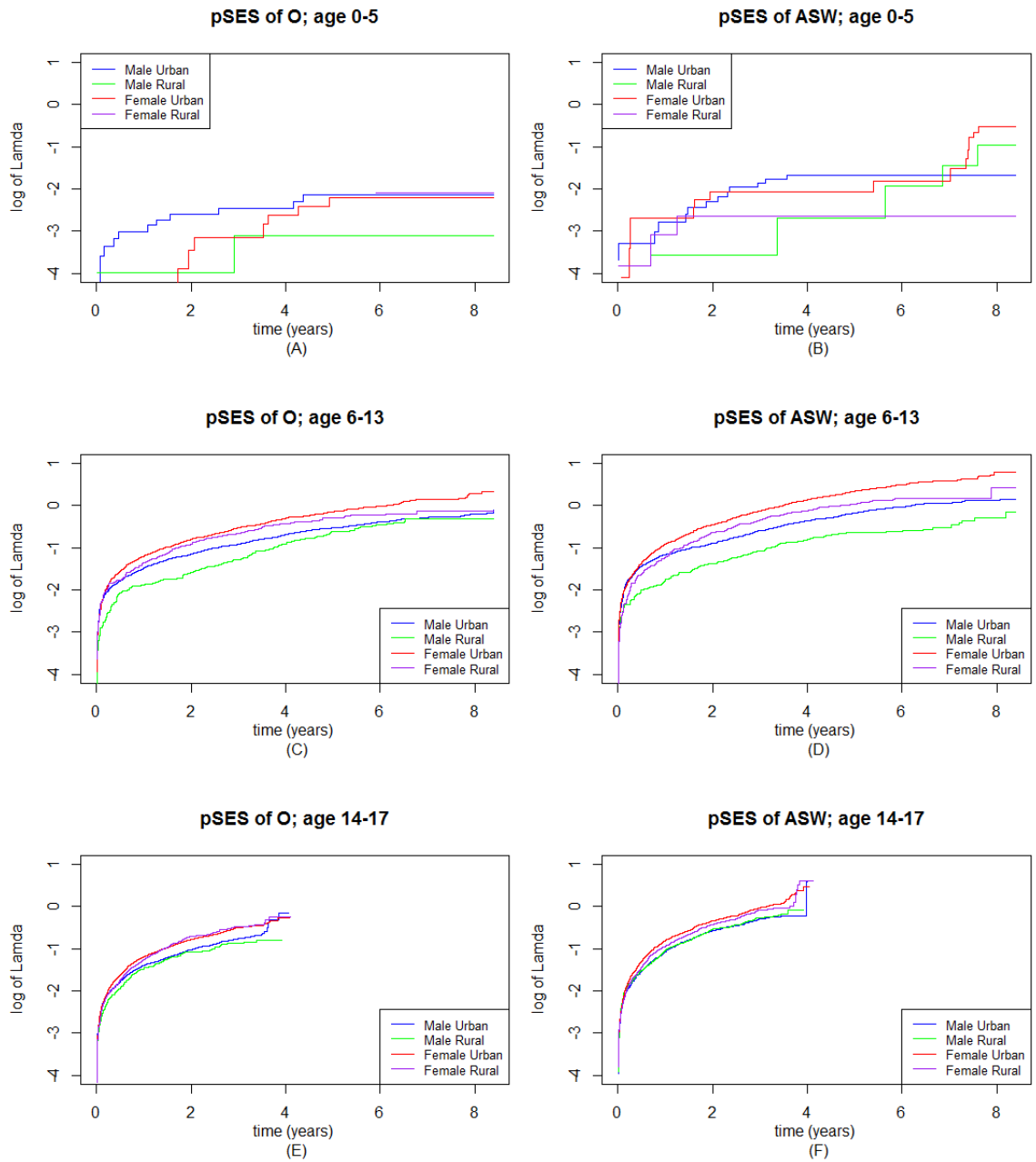


Figure 2.4: The log-transformed generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration.

## 2.3 Analysis under Parametric Model

From Figure A.1, we can see that the cumulative intensity function increases similar to a power function of time. As an attempt to analyze the EDMH data, we specify the baseline intensity into  $\alpha t^{\alpha-1}$ , assuming all covariate effects are time-independent. The model is

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\alpha t^{\alpha-1} \exp\{\beta'Z\}. \quad (2)$$

We refer to this model as Model 2. Let  $\theta$  denote all the unknown parameters under Model 1, including the regression parameters and the baseline functions. With Model 2, the unknown parameters  $\theta$  are specified as  $(\alpha, \beta)$ . The likelihood function is

$$\begin{aligned} L(\theta|data) &\propto \prod_{i=1}^n \prod_{t \in (0, C_i]} (\lambda(t; \mathcal{H}_i(t)))^{dN_i(t)} (1 - \lambda(t; \mathcal{H}_i(t)))^{1-dN_i(t)} \\ &= \prod_{i=1}^n \prod_{t \in (0, C_i]} \left( (\alpha t^{\alpha-1} \exp\{\beta'Z_i\})^{dN_i(t)} \right) \times \exp \left\{ - \int_0^{C_i} Y_i^R(t) \alpha t^{\alpha-1} \exp\{\beta'Z_i\} dt \right\}, \end{aligned} \quad (2.1)$$

with the log-likelihood function

$$l(\theta|data) = \sum_{i=1}^n \int_0^{C_i} Y_i^C(t) \left[ \log(\alpha t^{\alpha-1} e^{\beta'Z_i}) dN_i(t) - Y_i^R(t) \alpha t^{\alpha-1} e^{\beta'Z_i} dt \right]. \quad (2.2)$$

Applying the Newton-Raphson algorithm, we attain the maximum likelihood estimate (MLE) of  $\theta$  by maximizing the log-likelihood function  $l(\theta|data)$  in (2.2). The formulas of the likelihood score function and the observed information matrix are listed in Appendix C. Table 2.7 gives the estimates of the parameters in Model 2 with the risk set adjusted and not-adjusted for hospital duration respectively. We also present the estimated standard errors in the table. The significant regression parameters are bold.

From Table 2.7, we can see that all the risk factors in Model 2 (the parametric version of Model 1) had significant effects on the recurrence of EDMH visits. The positive sign of the covariate coefficient suggests that an increasing risk of having recurrent EDMH visit is associated with subjects having index initial EDMH visits as teenagers (versus younger subjects), and people living in urban

Table 2.7: MLEs of the regression coefficients in Model 2.

Model	$\alpha$	pSES (O vs ASW)	Age (at index initial ED)	Sex (Male vs Female)	Region (Urban vs Rural)
With risk set adjusted for hospital duration					
Estimates	<b>0.205</b>	<b>-0.524</b>	<b>0.085</b>	<b>-0.474</b>	<b>0.111</b>
SE	0.006	0.017	0.003	0.019	0.020
$\log(L(\hat{\theta}))$	-104470.5				
With risk set not adjusted for hospital duration					
Estimates	<b>0.214</b>	<b>-0.531</b>	<b>0.091</b>	<b>-0.486</b>	<b>0.133</b>
SE	0.007	0.017	0.003	0.019	0.020
$\log(L(\hat{\theta}))$	-104844.1				

area are more likely to visit ED for mental illness than those in rural area. The coefficients of pSES and sex are negative, which suggest that subjects with pSES of O tend to have lower risk to the next EDMH visit than those with other pSES, and a decreased EDMH visit risk is associated with males (versus females). This is consistent with the patterns shown by the generalized Nelson-Aalen estimates in Chapter 2.2.2.

The estimates of regression parameters with adjustment for hospital duration are very similar to those without adjustment for hospital duration. As we mentioned above, it may be due to the small portion, 28%, of the study subjects had hospitalization records.

If the parametric model assumption is appropriate, the log-transformed cumulative intensity function should be a linear function of  $\log(t)$ . This provides a method to check the model assumption. We checked the pattern of the log-transformed estimator for the cumulative intensity function in each subgroup versus the log-transformed event time. Figure 2.5 represents the patterns in the three age groups consisting of males in urban area with pSES of category A, S, and W. The solid lines in Figure 2.5 represent the log-transformed generalized Nelson-Aalen estimators with the risk set not-adjusted for the hospital duration; while the dashed lines stand for the log-transformed estimators under Model 2 with the risk set not-adjusted for the hospital duration. The three dashed lines are plotted with Age=4, Age=10, and Age=16 for pre-school, elementary, and teenager respectively.

Figure 2.5 shows that the parametric assumption could be inappropriate for the baseline intensity function, especially for that in the pre-school group. That could be due to the reason that only a small number of subjects are in this group. The estimate for the pre-school group is very different



with those for the other two age groups. Although the two solid lines for elementary school group and for teenager group are very close to each other, they are not parallel. This suggests that the baseline intensity functions in different subgroups may vary. This motivates the application of the semi-parametric regression model in Chapter 3.

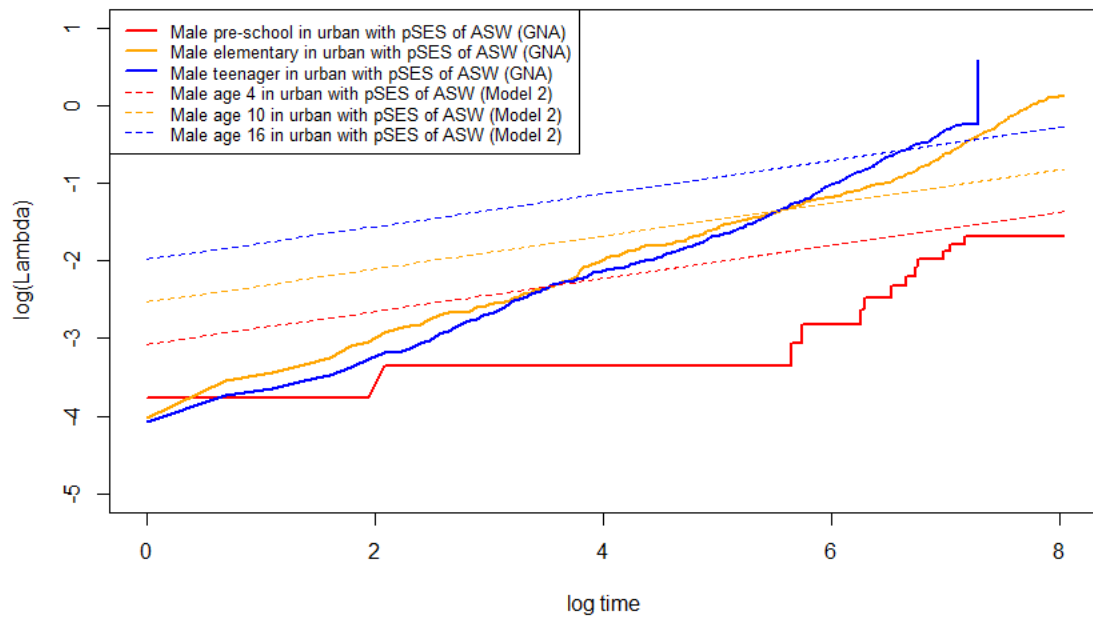


Figure 2.5: The log-transformed estimates of the cumulative intensity functions against log of event times.

## Chapter 3

# Stratified Cox Regression Models

This chapter begins with a review of a stratified Cox regression model along with the estimation procedures. The stratified Cox regression model is applied to analyze the recurrent EDMH visits described in Chapter 2. We consider special forms of the stratified Cox regression model with different stratification variables and with time-independent or time-dependent covariate effects.

### 3.1 Statistical Modeling

Assume that  $Y^R(t)$  is fully determined by  $\mathcal{H}(t)$ , the subject's history information at time  $t$ . Introduce a discrete stratification variable  $s\{\mathcal{H}(t)\}$ , determined by the history information and with all the possible values  $s = 1, \dots, S$ . We consider a specification of Model 1: for  $t > 0$  and  $s\{\mathcal{H}(t)\} = s$ ,

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\lambda_{0s}(t)\exp\{\beta(t;\beta_s)'Z(t)\}. \quad (3)$$

As refer to this model as Model 3. This model involves the stratum-specific baseline intensity function  $\lambda_{0s}(t)$ , and stratum-specific regression coefficients  $\beta(t;\beta_s)$ . Here we assume  $\beta(t;\beta_s)$  to be a known function of  $t$  up to parameters  $\beta_s$ . It is an extension of one of the two semi-parametric models considered in Prentice et al (1981), where  $\beta(t;\beta_s) = \beta_s$  is assumed and the corresponding

intensity function is

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\lambda_{0s}(t) \exp\{\beta'_s Z(t)\}. \quad (3a)$$

We refer to this model as Model 3a. This model allows the regression coefficients to differ across different strata. Prentice et al (1981) gives an example of the stratification variable as  $s\{\mathcal{H}(t)\} = N(t) + 1$ . In the current application, a subject belongs to stratum  $s\{\mathcal{H}(t)\} = N(t) + 1$ , if the number of recurrent events occurred is  $N(t)$  at time  $t$ . This specification may accommodate situations with non-Poisson processes using an appropriate stratification variable  $s\{\mathcal{H}(t)\}$ .

One may assume that risk factors have the same effect among different strata by restricting  $\beta_s = \beta$ . This assumption gives the following intensity function

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\lambda_{0s}(t) \exp\{\beta' Z(t)\}. \quad (3b)$$

We refer to this model as Model 3b. When  $\lambda_{0s}(t) = \lambda_0(t)$ , Model 3b reduces to the well-known AG model of Anderson and Gill (1982). The corresponding intensity function is then

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\lambda_0(t) \exp\{\beta' Z(t)\}. \quad (3c)$$

As refer to the AG model as Model 3c. This model assumes that the conditional intensity function is independent of the history of the counting process,  $\mathcal{N}(\cdot)$ . Thus, the counting process  $\{N(t) : t > 0\}$  is a Poisson process. That is, the number of EDMH visits in non-overlapping time interval is independent from each other given the covariates.

## 3.2 Estimation Procedures

### 3.2.1 Estimating Regression Parameters

Assume that there are no tied event times. Let  $t_{s,1} < \dots < t_{s,d_s}$  denote the ordered EDMH visit times in stratum  $s$ , where  $d_s$  is the total number of the EDMH visit times and  $s = 1, \dots, S$ .

Adapting the approach of Prentice et al (1981), a generalized version of the partial likelihood derivation given in Cox (1975), the overall likelihood function based on the data can be factored into

two terms,  $L_p(\beta)$  and  $L_1(\lambda_0(\cdot), \beta)$ , where  $L_p(\beta)$  is usually called partial likelihood function (Cox, 1975), and  $\beta$  is  $(\beta_1, \dots, \beta_S)$ , or the common values of  $\beta_s$ 's when the coefficients are the same across the strata. Maximize  $L_p(\beta)$  with respect to  $\beta$  to get the maximum partial likelihood estimator (MPLE)  $\hat{\beta}$  of  $\beta$ , then maximize  $L_1(\lambda_0(\cdot), \hat{\beta})$  with respect to  $\lambda_0(t)$ . Specifically, the partial likelihood function of the regression coefficients  $\beta_s$ 's is

$$PL(\beta) = \prod_{s=1}^S \prod_{j=1}^{d_s} \frac{\exp\{\beta(t_{s,j}^*; \beta_s)' Z_{i(t_{s,j}^*)}(t_{s,j}^*)\}}{\sum_{l:s\{\mathcal{H}_l(t_{s,j}^*)\}=s} Y_l^C(t_{s,j}^*) Y_l^R(t_{s,j}^*) \exp\{\beta(t_{s,j}^*; \beta_s)' Z_l(t_{s,j}^*)\}}, \quad (3.1)$$

where  $i(t_{s,j}^*)$  is the index of the subject who experiences the event occurrence in stratum  $s$  at time  $t_{s,j}^*$ .

Using counting process notation, the log-partial likelihood function of  $\beta$  can be expressed as  $\log\{PL(\beta)\} = C(\beta; \infty)$  with

$$C(\beta; t) = \sum_{i=1}^n \int_0^t Y_i^C(u) \sum_{s=1}^S [\beta(u; \beta_s)' Z_i(u) - \log\left\{ \sum_{l:s\{\mathcal{H}_l(u)\}=s} Y_l^C(u) Y_l^R(u) \exp\{\beta(u; \beta_s)' Z_l(u)\} \right\}] \Bigg|_{s=s\{\mathcal{H}_i(u)\}} dN_i(u). \quad (3.2)$$

The partial likelihood estimating function is then

$$U(\beta) = \frac{\partial C(\beta; \infty)}{\partial \beta} = \sum_{i=1}^n \int_0^\infty Y_i^C(u) \sum_{s=1}^S \left[ \frac{\partial \beta(u; \beta_s)'}{\partial \beta} Z_i(u) - \frac{\sum_{l:s\{\mathcal{H}_l(u)\}=s} Y_l^C(u) Y_l^R(u) \frac{\partial \beta(u; \beta_s)'}{\partial \beta} Z_l(u) \exp\{\beta(u; \beta_s)' Z_l(u)\}}{\sum_{l:s\{\mathcal{H}_l(u)\}=s} Y_l^C(u) Y_l^R(u) \exp\{\beta(u; \beta_s)' Z_l(u)\}} \right] \Bigg|_{s=s\{\mathcal{H}_i(u)\}} dN_i(u). \quad (3.3)$$

The MPLE  $\hat{\beta}$  of  $\beta$  can be attained by solving the estimating equation  $U(\beta) = 0$ . The variance of the MPLE  $\hat{\beta}$  is approximately the inverse of the second derivative of  $-\log\{PL(\beta)\}$ . We can obtain the robust variance estimator according to Hu et al (2003), which is equivalent to the so-called "sandwich" estimator given in Lin and Wei (1989). The well-developed counting process formulation for event history data analysis and the asymptotic derivation using the martingale results presented

in, for example, Fleming and Harrington (1991) and Andersen et al(1993), can be adapted to verify the consistency and asymptotic normality of MPLE  $\hat{\beta}$ , the maximum point of (3.2) at  $t = \infty$ . Furthermore, note from expression (3.2) that the counting process formulation accommodates the situations with tied event times, where are given similar treatments as to the discussion on tied failure times, for example, in Kalbfleisch and Prentice(1980, Chp 4).

A Wald-type test can be constructed based on the asymptotic normality of the MPLE on the coefficient for a covariate in Model 3 to assess the covariate effect. We may assess the goodness of fit of the model using the partial likelihood ratio test. The partial likelihood ratio test statistic, denoted by  $G$ , is calculated as the twice of the difference between the log partial likelihood of the model containing the covariates and the log partial likelihood of the model containing a subset of the covariates (the reduced model). Specifically,

$$G = 2\{C(\hat{\beta}, \infty) - C(\tilde{\beta}, \infty)\},$$

with  $\tilde{\beta}$  the MPLE of  $\beta$  under the reduced model.

### 3.2.2 Estimating Baseline Intensity Functions

With fixed  $\beta$ , the following estimating equations are unbiased:

$$\sum_{i:s\{\mathcal{H}_i(t)\}=s} Y_i^C(t) \left[ dN_i(t) - Y_i^R(t) \lambda_{0s}(t) \exp\{\beta(t; \beta_s)' Z_i(t)\} dt \right] = 0, \quad t > 0,$$

for  $s = 1, \dots, S$ . This yields the following estimation procedure. When the baseline intensity function in Model 3 varies from stratum to stratum, a consistent estimator of the cumulative baseline intensity function of stratum  $s$  is

$$\hat{\Lambda}_{0s}(t; \beta_s) = \int_0^t \frac{\sum_{i:s\{\mathcal{H}_i(u)\}=s} Y_i^C(u) dN_i(u)}{\sum_{l:s\{\mathcal{H}_l(u)\}=s} Y_l^C(u) Y_l^R(u) \exp\{\beta(u; \beta_s)' Z_l(u)\}}, \quad t > 0, \quad (3.4)$$

for  $s = 1, \dots, S$ . Here we take the convention  $0/0 = 0$ . In the situations with a single cumulative baseline intensity function, that is,  $\lambda_{0s}(t) = \lambda_0(t)$  for  $s = 1, \dots, S$ , we estimate the baseline function

using

$$\hat{\Lambda}_0(t; \beta) = \int_0^t \frac{\sum_{i=1}^n Y_i^C(u) dN_i(u)}{\sum_{s=1}^S \sum_{l:s\{\mathcal{H}_l(u)\}=s} Y_l^C(u) Y_l^R(u) \exp\{\beta(u; \beta_s)' Z_l(u)\}}, \quad t > 0. \quad (3.5)$$

With some regularity conditions and by the martingale central limit theorem, we can show that  $\hat{\Lambda}_{0s}(t; \beta)$  and  $\hat{\Lambda}_0(t; \beta)$ , after standardization, converge weakly to Gaussian processes with mean zero.

The above estimators are generalizations of the Breslow estimator for the baseline functions under Model 3 (Hu et al 2011). They can be used to estimate the intensity function of a particular group with the covariates fixed at the corresponding levels. This suggests an approach for model checking. The nonparametric versions of (3.4) and (3.5), taking  $\beta(t; \beta_s) = 0$ , are the generalized Nelson-Aalen estimator for the cumulative intensity functions with adjustment for event duration. One may group the subjects according to the covariates. With the information of each of the different subgroups, we may evaluate the generalized Nelson-Aalen estimator.

We obtain estimators for the baseline intensity functions  $\Lambda_{0s}(\cdot)$  and  $\Lambda_0(\cdot)$  conventionally by plugging in the corresponding MPLE of the unknown coefficients in (3.4) and (3.5), respectively. Given the continuity of  $\hat{\Lambda}_{0s}(t; \beta)$  and  $\hat{\Lambda}_0(t; \beta)$  as functions of  $\beta$ , we can prove that the resulting estimators  $\hat{\Lambda}_{0s}(\cdot; \hat{\beta}_s)$  and  $\hat{\Lambda}_0(\cdot; \hat{\beta})$  are consistent and, after standardization, weakly converge to mean zero Gaussian processes.

### 3.3 Analysis Results

#### 3.3.1 With Time-independent Regression Coefficients

We obtain the MPLEs of the regression coefficients in the three special cases of Model 3, assuming all the covariate effects are time-independent. We consider two stratification variables, age at the index initial EDMH visit (pre-school 0-5, elementary school 6-13, and teenager 14-17), and season of the EDMH visits (fall, spring, summer, and winter). Table 3.1 represents the MPLEs of the regression coefficients in those models. The estimated robust standard errors are shown in the brackets. The log-partial likelihood functions are evaluated at the estimates of the regression coefficients under different models. The coefficient estimates with significant effect are bold in the table. Table B.1

in Appendix B provides corresponding estimates with the non-parametric bootstrap estimates of the standard errors, with bootstrap sizes  $B=1000$ ,  $B=2000$ , and  $B=5000$ , respectively. The estimates for the standard errors are similar when  $B=2000$  and  $B=5000$ . Thus, we use  $B=2000$  for all later analyses.

Model  $3a_{\beta_s, \lambda_{0s}}^{s:age}$ ,  $3b_{\beta_s, \lambda_{0s}}^{s:age}$ ,  $3a_{\beta_s, \lambda_{0s}}^{s:age-}$ , and  $3b_{\beta_s, \lambda_{0s}}^{s:age-}$  use age at the index initial EDMH visit as the stratification variable, while Model  $3a_{\beta_s, \lambda_{0s}}^{s:seasons}$  and  $3b_{\beta_s, \lambda_{0s}}^{s:seasons}$  use seasons of EDMH visits as the stratification variable.

From Table 3.1, we can see that the four risk factors have significant effects on the recurrence of EDMH visits. Lower risk of recurrent EDMH visits is associated with people having pSES in the category of O and those who are females; while older patients at index time and those who are living in urban area tend to have higher risk of recurrent EDMH visits.

Proceeding across columns in Table 3.1 under each set of models, we can see the estimated effects of the risk factors versus their baseline according to different strata. The results indicate that the covariate effects vary from stratum to stratum: The partial likelihood ratio tests for comparing the fit of the two models, denoted by  $3a_{\beta_s, \lambda_{0s}}^{s:age}$  and  $3b_{\beta_s, \lambda_{0s}}^{s:age}$ ,  $3a_{\beta_s, \lambda_{0s}}^{s:age-}$  and  $3b_{\beta_s, \lambda_{0s}}^{s:age-}$ ,  $3a_{\beta_s, \lambda_{0s}}^{s:seasons}$  and  $3b_{\beta_s, \lambda_{0s}}^{s:seasons}$ , have test statistics of value 177.2 with  $df = 8$  (p-value<0.001), 69.6 with  $df = 6$  (p-value<0.001), and 140.2 with  $df = 12$  (p-value<0.001), respectively. This indicates the effects of risk factors vary across different strata.

The effects of pSES are significant under all the different models in Table 3.1, and negative across all strata. The standard error for pSES in the first row of stratum 1 (pre-school stratum), under Model  $3a_{\beta_s, \lambda_{0s}}^{s:age}$ , is quite larger than that in other two strata. It may be due to the small number of subjects in pre-school stratum.

All the significant coefficients of sex, under different models, are negative but have different magnitudes. It appears that the difference between males and females are non-statistically significant in the pre-school group but significant in the elementary group. The coefficient  $-0.379$  indicates females are at higher risk to repeated EDMH visits than males in the elementary group. However, Figure 2.1 shows that a little more visits were made by males in the elementary group

The effects of residential region are significant and have positive sign in fall, spring, and summer; while in winter, the effect are non-significant and have opposite sign. It represents that subjects

living in urban area have more EDMH visits in the three seasons, but have no significant difference in winter with those living in rural area.

We use age at the index initial EDMH visit, classified into three age groups, as the stratification variable in Model  $3a_{\beta_s, \lambda_{0s}}^{s:age}$ ,  $3b_{\beta, \lambda_{0s}}^{s:age}$ ,  $3a_{\beta_s, \lambda_{0s}}^{s:age-}$ , and  $3b_{\beta, \lambda_{0s}}^{s:age-}$ . In order to explore the effect of age within each of the three strata, Model  $3a_{\beta_s, \lambda_{0s}}^{s:age}$  and  $3b_{\beta, \lambda_{0s}}^{s:age}$  in Table 3.1 include age at the index initial EDMH visit as the covariate. We can see from Model  $3a_{\beta_s, \lambda_{0s}}^{s:age}$  that, the estimated age effect steadily decreases from stratum 1 to stratum 3, and changes sign in stratum 3. However, the estimate in stratum 3 is not significant. The average covariate effects under model  $3b_{\beta, \lambda_{0s}}^{s:age}$  are in agreement with the overall covariate effects without stratification under model 3c. The partial likelihood ratio tests for comparing the fit of  $3a_{\beta_s, \lambda_{0s}}^{s:age}$  vs  $3a_{\beta_s, \lambda_{0s}}^{s:age-}$ , and  $3b_{\beta, \lambda_{0s}}^{s:age}$  vs  $3b_{\beta, \lambda_{0s}}^{s:age-}$ , have test statistics of value 282.2 with  $df = 3$  (p-value < 0.001), and 174.6 with  $df = 1$  (p-value < 0.001), respectively. It indicates that under the assumption that baseline intensity function have variations across the age groups, there is still difference among subjects who are in the same age group but with different age at index initial EDMH visits. Moreover, the effect of residential region shows significance in stratum 1 in Model  $3a_{\beta_s, \lambda_{0s}}^{s:age-}$ , compared with Model  $3a_{\beta_s, \lambda_{0s}}^{s:age}$ . This indicates that there might be interaction between age and residential region. This finding motivated us to consider models with two factor interactions in Chapter 3.3.3.

We also estimate the cumulative intensity function with the MPLEs of the parameters under Model 3c in Table 3.1. Figure A.2 in Appendix A gives the estimate of the cumulative intensity function for each subgroup. The dashed lines represent the Breslow estimates under Model 3c in Table 3.1, while the solid lines show the generalized Nelson-Aalen estimates. The dashed lines are plotted with Age=3, Age=11, and Age=16 for pre-school, elementary, and teenager respectively.

From Figure A.2, it is clear to see that the Breslow estimates show the closeness to the generalized Nelson-Aalen estimates in some subgroups, with exceptions in the subgroups with pSES in the category of ASW, age of 6~13 (Figure A.2.c). Nevertheless, the model assumption seems to be valid overall.



Table 3.1: MPLEs of the regression coefficients and robust estimates for the standard errors for the standard errors in Model 3a/3b/3c with time-independent covariate effects.

Model	Stratification Variable 1*			Stratification Variable 2**				AG Model
	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$	$\beta(t; \beta_s) = \beta$
	$3a^{stage}$	$3b^{stage}$	$3c^{stage}$	$3a^{seasons}$	$3b^{seasons}$	$3c^{seasons}$		
	$\beta_s, \lambda_{0s}$	$\beta_s, \lambda_{0s}$	$\beta_s, \lambda_{0s}$	$\beta_s, \lambda_{0s}$	$\beta_s, \lambda_{0s}$	$\beta_s, \lambda_{0s}$		
	$\beta(= \beta_s)$	$\beta(= \beta_s)$	$\beta(= \beta_s)$	$\beta(= \beta_s)$	$\beta(= \beta_s)$	$\beta(= \beta_s)$		
$\log(PL(\hat{\beta}))$	-112370.2	-112458.8	-112458.8	-102307.9	-102378.0	-102378.0		-119461.02
pSES	<b>-0.768</b> (0.264)	<b>-0.375</b> (0.032)	<b>-0.418</b> (0.022)	<b>-0.451</b> (0.035)	<b>-0.364</b> (0.035)	<b>-0.520</b> (0.040)	<b>-0.371</b> (0.038)	<b>-0.420</b> (0.018)
Age	<b>0.662</b> (0.115)	<b>0.137</b> (0.009)	<b>-0.009</b> (0.012)	<b>0.033</b> (0.007)	<b>0.122</b> (0.008)	<b>0.110</b> (0.009)	<b>0.094</b> (0.008)	<b>0.085</b> (0.004)
Sex	-0.265 (0.262)	<b>-0.379</b> (0.033)	<b>-0.223</b> (0.023)	<b>-0.254</b> (0.036)	<b>-0.399</b> (0.036)	<b>-0.260</b> (0.041)	<b>-0.194</b> (0.039)	<b>-0.282</b> (0.019)
Region	0.562 (0.336)	<b>0.279</b> (0.041)	0.044 (0.026)	<b>0.178</b> (0.042)	<b>0.129</b> (0.041)	<b>0.235</b> (0.048)	<b>-0.030</b> (0.043)	<b>0.126</b> (0.022)
	$3a^{stage}$	$3b^{stage}$	$3c^{stage}$					
	$\beta_s, \lambda_{0s}$	$\beta_s, \lambda_{0s}$	$\beta_s, \lambda_{0s}$					
	$\beta(= \beta_s)$	$\beta(= \beta_s)$	$\beta(= \beta_s)$					
$\log(PL(\hat{\beta}))$	-112511.3	-112546.1	-112546.1					
pSES	<b>-0.847</b> (0.263)	<b>-0.360</b> (0.032)	<b>-0.418</b> (0.022)					
Sex	0.009 (0.258)	<b>-0.460</b> (0.032)	<b>-0.224</b> (0.023)					
Region	<b>0.746</b> (0.334)	<b>0.273</b> (0.041)	0.045 (0.026)					

\* Age at the index initial EDMH visit as the stratification variable (pre-school 0-5, elementary school 6-13, and teenager 14-17).

\*\* Seasons of the EDMH visits as the stratification variable: fall, spring, summer, and winter, respectively.

pSES as the indicator of O, Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban.

Estimated robust standard error in brackets.

Significant effect with p-value  $\leq 0.05$  in boldface.

Table 3.2: MPLEs of the regression coefficients and robust estimates for the standard errors in Model 3a/3b/3c with time-dependent covariate effects.

Model	Stratification Variable 1*			Stratification Variable 2**				AG Model
	$\beta_1$	$3a_{\beta_s, \lambda_{0s}}^{stage}$ $\beta_2$	$\beta_3$	$\beta_1$	$3a_{\beta_s, \lambda_{0s}}^{seasons}$ $\beta_2$	$\beta_3$	$\beta_4$	
$\log(PL(\hat{\beta}))$	-112311.3	-112389.7	-112389.7	-0.032	-0.085	-0.211	-0.053	-102277.6
pSES	-0.421 (0.763)	-0.088 (0.104)	-0.024 (0.065)	(0.101)	(0.103)	(0.118)	(0.109)	-0.088 (0.054)
Age	0.503 (0.319)	0.036 (0.028)	-0.018 (0.031)	-0.021 (0.016)	-0.003 (0.020)	-0.001 (0.022)	-0.017 (0.020)	-0.011 (0.010)
Sex	1.390 (0.950)	0.051 (0.104)	0.024 (0.064)	0.022 (0.098)	0.040 (0.102)	0.188 (0.117)	-0.021 (0.108)	0.054 (0.053)
Region	0.359 (0.938)	<b>0.344</b> (0.133)	0.125 (0.074)	0.061 (0.117)	0.099 (0.121)	<b>0.358</b> (0.145)	0.087 (0.125)	<b>0.131</b> (0.063)
pSES $\times$ ln(t)	-0.140 (0.284)	<b>-0.116</b> (0.040)	<b>-0.189</b> (0.029)	<b>-0.191</b> (0.043)	<b>-0.123</b> (0.043)	<b>-0.136</b> (0.049)	<b>-0.146</b> (0.047)	<b>-0.150</b> (0.023)
Age $\times$ ln(t)	0.066 (0.120)	<b>0.042</b> (0.011)	0.004 (0.016)	<b>0.025</b> (0.007)	<b>0.058</b> (0.009)	<b>0.050</b> (0.010)	<b>0.053</b> (0.009)	<b>0.045</b> (0.004)
Sex $\times$ ln(t)	-0.641 (0.343)	<b>-0.174</b> (0.040)	<b>-0.121</b> (0.029)	<b>-0.127</b> (0.042)	<b>-0.195</b> (0.043)	<b>-0.198</b> (0.049)	<b>-0.077</b> (0.047)	<b>-0.151</b> (0.023)
Region $\times$ ln(t)	0.074 (0.354)	-0.026 (0.051)	-0.039 (0.033)	0.053 (0.050)	0.012 (0.051)	-0.055 (0.061)	-0.055 (0.053)	-0.004 (0.027)

\* Age at the index initial EDMH visit as the stratification variable (pre-school 0-5, elementary school 6-13, and teenager 14-17).

\*\* Seasons of the EDMH visits as the stratification variable: fall, spring, summer, and winter, respectively.

pSES as the indicator of O, Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban. Estimated robust standard error in brackets.

Significant effect with p-value  $\leq 0.05$  in boldface.

Table 3.3: MPLEs of the regression coefficients and robust estimates for the standard errors in Model 3a/3b/3c with all pairs of two factor interactions and time-independent covariate effects.

Model	Stratification Variable 1*			Stratification Variable 2**			AG Model
	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_1$	$\beta_2$	$\beta_3$	
$\log(PL(\hat{\beta}))$	-112304.1	-112441.5	-102258.9	-102258.9	-102357.8	-119440.9	-119440.9
pSES	<b>-9.949</b> (3.345)	<b>-0.278</b> (0.088)	<b>-0.392</b> (0.050)	<b>-0.319</b> (0.077)	<b>-0.304</b> (0.090)	<b>-0.376</b> (0.040)	<b>-0.367</b> (0.040)
Age	1.151	<b>0.818</b> (0.803)	<b>-0.224</b> (0.074)	<b>0.375</b> (0.050)	<b>0.415</b> (0.058)	<b>0.389</b> (0.025)	<b>0.334</b> (0.025)
Sex	2.099	<b>-0.924</b> (2.905)	<b>-0.293</b> (0.055)	<b>-0.368</b> (0.080)	<b>-0.445</b> (0.096)	<b>-0.385</b> (0.086)	<b>-0.321</b> (0.043)
Region	<b>8.810</b> (3.153)	0.147 (0.087)	0.036 (0.044)	<b>0.275</b> (0.064)	<b>0.246</b> (0.074)	<b>-0.093</b> (0.068)	<b>0.152</b> (0.034)
pSES×Age	<b>-1.978</b> (0.786)	-0.040 (0.053)	0.024 (0.063)	-0.054 (0.034)	-0.037 (0.041)	0.014 (0.044)	-0.031 (0.020)
pSES×Sex	<b>1.346</b> (0.596)	<b>0.143</b> (0.067)	-0.014 (0.048)	-0.064 (0.072)	-0.042 (0.082)	0.081 (0.080)	0.056 (0.038)
pSES×Region	1.636 (0.949)	<b>-0.239</b> (0.082)	-0.039 (0.052)	<b>-0.168</b> (0.084)	<b>-0.259</b> (0.097)	-0.026 (0.087)	<b>-0.099</b> (0.043)
Age×Sex	0.488 (0.743)	<b>-0.419</b> (0.053)	<b>0.292</b> (0.065)	0.062 (0.034)	<b>0.179</b> (0.041)	-0.047 (0.043)	0.038 (0.020)
Age×Region	<b>2.159</b> (0.834)	<b>-0.223</b> (0.072)	0.111 (0.073)	<b>-0.118</b> (0.043)	<b>-0.102</b> (0.049)	<b>-0.138</b> (0.051)	<b>-0.117</b> (0.024)
Sex×Region	-1.457 (0.812)	0.042 (0.087)	-0.017 (0.055)	-0.080 (0.086)	<b>-0.193</b> (0.086)	<b>0.254</b> (0.101)	0.018 (0.045)

\* Age at the index initial EDMH visit as the stratification variable (pre-school 0-5, elementary school 6-13, and teenager 14-17).

\*\* Seasons of the EDMH visits as the stratification variable: fall, spring, summer, and winter, respectively.

pSES as the indicator of O, Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban. Estimated robust standard error in brackets.

Significant effect with p-value  $\leq 0.05$  in boldface.

Table 3.4: MPLEs of the regression coefficients and robust estimates for the standard errors in Model 3a/3b/3c with all pairs of two factor interactions and time-dependent covariate effects.

Model	Stratification Variable 1*			Stratification Variable 2**			AG Model
	$3a^{\text{stage}}$ $\beta_1$	$\beta_2$	$\beta_3$	$3a^{\text{seasons}}$ $\beta_2$	$\beta_3$	$\beta_4$	
$\log(PL(\hat{\beta}))$	-112247.2	-112367.4	-112367.4	-102148.9	-102253.5	-119335.6	-119335.6
pSES	<b>-10.136</b> (3.500)	0.011 (0.132)	0.099 (0.088)	0.012 (0.124)	-0.017 (0.127)	0.081 (0.149)	0.008 (0.066)
Age	0.824 (1.312)	<b>0.593</b> (0.112)	-0.127 (0.116)	0.075 (0.066)	0.052 (0.077)	<b>0.208</b> (0.089)	<b>0.117</b> (0.038)
Sex	3.882 (3.182)	<b>-0.470</b> (0.142)	-0.083 (0.092)	0.124 (0.125)	-0.044 (0.131)	0.106 (0.153)	0.021 (0.067)
Region	<b>8.818</b> (3.397)	0.239 (0.156)	0.102 (0.095)	0.127 (0.136)	<b>0.351</b> (0.140)	<b>0.462</b> (0.169)	<b>0.258</b> (0.073)
pSES×Age	<b>-2.029</b> (0.791)	-0.043 (0.053)	-0.132 (0.068)	<b>-0.094</b> (0.036)	-0.069 (0.043)	-0.083 (0.047)	<b>-0.065</b> (0.021)
pSES×Sex	<b>1.334</b> (0.614)	0.127 (0.068)	-0.027 (0.048)	-0.086 (0.073)	<b>0.225</b> (0.073)	-0.076 (0.083)	0.034 (0.038)
pSES×Region	1.647 (0.951)	<b>-0.242</b> (0.082)	-0.043 (0.052)	0.041 (0.085)	<b>-0.165</b> (0.083)	<b>-0.264</b> (0.097)	<b>-0.102</b> (0.044)
Age×Sex	0.546 (0.766)	<b>-0.434</b> (0.053)	<b>0.226</b> (0.069)	0.039 (0.035)	<b>0.142</b> (0.044)	<b>-0.107</b> (0.047)	0.001 (0.021)
Age×Region	<b>2.166</b> (0.844)	<b>-0.224</b> (0.072)	0.092 (0.078)	<b>-0.121</b> (0.045)	<b>-0.114</b> (0.053)	-0.108 (0.060)	<b>-0.135</b> (0.026)
Sex×Region	-1.433 (0.831)	0.041 (0.087)	-0.020 (0.055)	-0.080 (0.087)	<b>-0.192</b> (0.087)	<b>0.240</b> (0.102)	0.014 (0.046)
pSES×ln(t)	0.001 (0.317)	<b>-0.115</b> (0.040)	<b>-0.211</b> (0.031)	<b>-0.221</b> (0.045)	<b>-0.133</b> (0.045)	<b>-0.166</b> (0.052)	<b>-0.167</b> (0.023)
Age×ln(t)	0.132 (0.380)	<b>0.099</b> (0.032)	0.012 (0.045)	<b>0.068</b> (0.021)	<b>0.170</b> (0.025)	<b>0.133</b> (0.028)	<b>0.125</b> (0.012)
Sex×ln(t)	-0.615 (0.346)	<b>-0.186</b> (0.041)	<b>-0.088</b> (0.031)	<b>-0.118</b> (0.044)	<b>-0.140</b> (0.045)	<b>-0.237</b> (0.052)	<b>-0.149</b> (0.024)
Region×ln(t)	-0.005 (0.395)	-0.037 (0.051)	-0.028 (0.036)	0.019 (0.053)	-0.037 (0.054)	-0.095 (0.065)	-0.048 (0.028)

\* Age at the index initial EDMH visit as the stratification variable (pre-school 0-5, elementary school 6-13, and teenager 14-17).

\*\* Seasons of the EDMH visits as the stratification variable: fall, spring, summer, and winter, respectively.

pSES as the indicator of O. Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban.

Estimated robust standard error in brackets.

Significant effect with p-value  $\leq 0.05$  in boldface.

### 3.3.2 With Time-dependent Regression Coefficients

To understand the variations of the covariate effects overtime, we consider models with time-dependent effects of the risk factors by assuming the regression coefficients in Model 3 as a linear function of event time. That is,  $\beta(t; \beta_s) = \beta_{s0} + \beta_{s1} \ln(t)$ , for  $s \in \{\mathcal{H}(t)\} = s$ . We start with models assuming that all risk factors have time-dependent effects. Table 3.2 summarizes the MPLEs of the regression coefficients in three special cases of Model 3, assuming time-dependent covariate effects. The estimated robust standard errors are shown in the brackets. The log-partial likelihood functions are evaluated at the estimates of the regression coefficients under model assumptions. The estimates of the coefficients with significant effect are bold in the table. Table B.2 in Appendix B shows corresponding estimates with the non-parametric bootstrap estimates of the standard errors with bootstrap size  $B=2000$ .

The partial likelihood ratio tests for comparing the goodness of fit of the models in Table 3.2 with their corresponding models in Table 3.1 have  $p\text{-value} < 0.001$ . For example, the partial likelihood ratio test for comparing the fit of the two models, denoted by Model 3c in Table 3.2 and Model 3c in Table 3.1, has a test statistic of value 202.6 with  $df = 4$ , which provides a  $p\text{-value}$  less than 0.001 with the chi-square distribution approximate. This result shows that there exists strong evidence that the effects of the risk factors are time-dependent. We remove the non-statistically significant term, the slope of the time-varying effect of region in Model 3c in Table 3.2, and refit the model. The partial likelihood ratio test gives the value of the test statistic of 0.01 with  $df = 1$ . Thus, there is no need for model reduction.

The results in Table 3.2 indicate that there are statistical significant time trends for the effect of pSES, age, and gender. The overall covariate effect without stratification, denoted by Model 3c in Table 3.2, are in agreement with the average covariate effects under Model 3b<sup>s:age</sup> <sub>$\beta, \lambda_{0s}$</sub>  except for the intercept term of the time-varying effect of region: under the assumption that the baseline intensity functions are different across strata (three age groups), the effect of region becomes non-statistically significant. It also indicates that there might be an interaction between age and residential region.

We also estimate the cumulative intensity function with the MPLEs of the regression parameters under Model 3c in Table 3.2. Figure A.3 in Appendix A gives the estimate of the cumulative intensity function for each subgroup. The dashed lines represent the Breslow estimates under Model 3c in

Table 3.2, while the solid lines show the generalized Nelson-Aalen estimates. Again, we plot the dashed lines with Age=3, Age=11, and Age=16 for pre-school, elementary, and teenager respectively.

From Figure A.3, we can see that the Breslow estimates under Model 3c in Table 3.2 become closer to the generalized Nelson-Aalen estimates in the subgroups with pSES in the category of ASW, age of 6~13 (Figure A.3.c), compared with the estimates under Model 3c in Table 3.1.

### 3.3.3 Main Effects together with Two Factor Interactions I

Results from the previous analysis indicate that there might be an interaction between age and residential region. This section consider the effects of two-factor interactions on the recurrence of EDMH visits. We start with models assuming all pairs of two-factor interactions. Table 3.3 summarizes the MPLEs of the regression coefficients in the corresponding three special cases of Model 3. The estimated robust standard errors are shown in the brackets. The log-partial likelihood functions are evaluated at the estimates of the parameters under model assumptions. The estimates of the coefficients with significant effect are bold in the table. Table B.3 in Appendix B shows corresponding estimates with the non-parametric bootstrap estimators of the standard errors ( $B=2000$ ).

We see from Table 3.3 that there exists strong evidence that the interaction terms are nonzero. For example, the partial likelihood ratio statistic for testing the null hypothesis that Model 3c in Table 3.1 holds against Model 3c in Table 3.3 gives a value of 40.2 with  $df = 6$ , which provides a p-value less than 0.001 with the chi-square distribution approximate. Therefore, Model 3c considering the two-factor interactions in Table 3.3, significantly improves the fit.

We made a model deduction for Model 3c in Table 3.3 by removing the most non-statistically significant interaction between gender and region. The partial likelihood ratio test for comparing the two model has a p-value of 0.685 with a degree of freedom 1. Repeating the step, we ended with the Model in Table 3.5. Note that the tests for interactions between pSES and region, age and gender are significant using the robust standard error estimates, but non-statistically significant when applying the non-parametric bootstrap standard error estimates. The result indicates a statistically significant effect only from the interaction between age and region.

We estimated the cumulative intensity function with MPLEs of the regression parameters under

Model 3c in Table 3.5. Figure A.4 in Appendix A represent the estimate of the cumulative intensity function for each subgroup. The dashed lines represent the Breslow estimates under Model 3c in Table 3.5, are plotted with Age=3, Age=11, and Age=16 for pre-school, elementary, and teenager respectively.

From Figure A.4, we can clearly see that the Breslow estimate under Model 3c in Table 3.5 shows closeness to the corresponding generalized Nelson-Aalen estimate in each of different subgroups. It indicates Model 3c fits reasonably well.

Table 3.5: MPLEs of the regression coefficients and the estimates of the standard errors under the fitted AG model considering two factor interactions and time-independent covariate effects.

Model	Variable	$\hat{\beta}$	SE1	p-value1	SE2*	p-value2*
3c	pSES	<b>-0.345</b>	0.038	0.000	0.054	0.000
	Age	<b>0.321</b>	0.024	0.000	0.032	0.000
	Sex	<b>-0.276</b>	0.019	0.000	0.026	0.000
	Region	<b>0.155</b>	0.030	0.000	0.047	0.001
	pSES×Region	-0.093	0.043	0.031	0.063	0.133
	Age×Sex	0.039	0.020	0.049	0.028	0.157
	Age×Region	<b>-0.121</b>	0.024	0.000	0.032	0.000
log(PL( $\hat{\beta}$ ))		-119443.7				
* SE2 are the non-parametric bootstrap standard errors of the parameters, and p-value2 are the corresponding p-values (B=2000). pSES as the indicator of O, Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban.						

### 3.3.4 Main Effects together with Two Factor Interactions II

This section consider the effects of two-factor interactions on the recurrence of EDMH visits and time-dependent effects of the risk factors. We start with models assuming all pairs of two-factor interactions and time-dependent effects of all the four risk factors. Table 3.4 summarizes the further analysis results with the three special cases of Model 3. Table B.4 in Appendix B shows corresponding estimates of the regression coefficients with the non-parametric bootstrap estimators of the standard errors (B=2000).

We compare the fits of the models in Table 3.4 with their corresponding models in Table 3.3 using the partial likelihood ratio test. The tests have p-value<0.001, which indicate that there exists strong evidence that the effects of risk factors are time-varying.

Again, the results indicate that the covariate effects vary from stratum to stratum: the p-values of the partial likelihood ratio tests for comparing the fit of the two models, denoted by  $3a_{\beta_s, \lambda_{0s}}^{s:age}$  and  $3b_{\beta_s, \lambda_{0s}}^{s:age}$ ,  $3a_{\beta_s, \lambda_{0s}}^{s:seasons}$  and  $3b_{\beta_s, \lambda_{0s}}^{s:seasons}$ , are less than 0.001.

The time trend effects of pSES are significant except in the pre-school group of Model  $3a_{\beta_s, \lambda_{0s}}^{s:age}$ . However, the interaction between pSES and age and the interaction between pSE and sex are significant in this group. It may indicate that the effect of pSES on the recurrence of EDMH visit is different among subjects of different ages and between males and females. No significant time trend for the effect of region is detected.

We made a model deduction for Model 3c in Table 3.4. The estimation results of the reduced model are shown in Table 3.6. Note that the tests for interactions between pSES and age, pSES and region are significant using the robust standard error estimates, but non-statistically significant when applying the non-parametric bootstrap standard error estimates.

The cumulative intensity function are estimated with MPLEs of the regression parameters under Model 3c in Table 3.6. Figure A.5 in Appendix A represent the estimates of the cumulative intensity functions. The dashed lines represent the Breslow estimates under Model 3c in Table 3.6, are plotted with Age=3, Age=11, and Age=16 for pre-school, elementary, and teenager respectively.

Figure A.5 shows that the Breslow estimates under Model 3c in Table 3.6 are closer to their corresponding generalized Nelson-Aalen estimates, compared with the Breslow estimates under the model in Table 3.5. Model 3c assuming the interactions and time-dependent coefficients improves the fit.



Table 3.6: MPLEs of the regression coefficients and the estimates of the standard errors under the fitted AG model considering two factor interactions and time-dependent covariate effects.

Model	Variable	$\hat{\beta}$	SE1	p-value1	SE2*	p-value2*
3c	pSES	0.031	0.063	0.626	0.076	0.685
	Age	<b>0.106</b>	0.036	0.003	0.045	0.019
	Sex	0.051	0.053	0.330	0.062	0.408
	Region	<b>0.154</b>	0.030	0.000	0.041	0.000
	pSES×Age	-0.068	0.021	0.001	0.035	0.059
	pSES×Region	-0.094	0.043	0.030	0.056	0.092
	Age×Region	<b>-0.124</b>	0.025	0.000	0.032	0.000
	pSES × ln(t)	<b>-0.174</b>	0.023	0.000	0.030	0.000
	Age × ln(t)	<b>0.126</b>	0.012	0.000	0.025	0.000
	Region × ln(t)	<b>-0.150</b>	0.023	0.000	0.031	0.000
log(PL( $\hat{\beta}$ ))		-119337.5				

\* SE2 are the non-parametric bootstrap standard errors of the parameters, and p-value2 are the corresponding p-values (B=2000).  
pSES as the indicator of O, Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban.

## Chapter 4

# Extended Renewal Process Model

As aforementioned in Chapter 1, the available data utilized in this project only contains records of mental health related ED visits from children and youth during April 1 2002 to March 31 2011. The starting point of one individual is from the time when he/she had the first EDMH visit during the observation window. Under this observation mechanism, different individuals have different starting points, and the starting point is related to the occurrence of EDMH visit. The interpretations of the analysis outcomes under the Cox regression models in Chapter 3, as well as under the parametric model in Chapter 2.3, are easier and meaningful only when the time origin is non-informative to the response. This motivates the renewal process model which assumes one individual is "brand-new" after the occurrence of each EDMH visit. If the renewal process model is appropriate, the corresponding analysis outcomes are easier to interpret than the outcomes under the Cox regression models. Since if individuals start from different time points, it is hard to tell the difference shown by the end.

This chapter focuses on the renewal process model, assuming that the intensity function of the event process depends only on the time since the most recent event. We first review a renewal process model with corresponding estimation procedures. Both a parametric and a semi-parametric version of the renewal process model are fitted with Alberta's PMHC data described in Chapter 2. We compare the analysis outcomes from the Cox regression models and the renewal process models by the end of this chapter.

## 4.1 Statistical Modeling

Assume that the baseline intensity function in the generalized Cox regression model (Model 1) to be  $\lambda_0\{t; \mathcal{H}(t)\} = \lambda_0(t - T_{N(t-)})$ , a function of the gap time since the most recent ED visit, and restrict the regression coefficients as  $\beta(t; \mathcal{H}(t)) = \beta(t - T_{N(t-)}; \beta)$ . We have, for  $t > 0$ ,

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\lambda_0(t - T_{N(t-)}) \exp\{\beta(t - T_{N(t-)}; \beta)'Z\}. \quad (4)$$

This model uses the gap time as the index of the baseline intensity function and regression coefficients. It assumes the shape of the intensity function to depend on the time from the immediately preceding ED visit time.

When  $\beta(t - T_{N(t-)}; \beta) = \beta$ , Model 4 assumes the effects of the risk factors are time-independent. The corresponding intensity function is

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\lambda_0(t - T_{N(t-)}) \exp\{\beta'Z\}. \quad (4a)$$

Further assuming  $Y^R(t) \equiv 1$  reduces Model 4a to the second semi-parametric model of Prentice et al (1981), labeled as formula (3) in their paper.

One may assume a parametric form for the baseline intensity function  $\lambda_0(\cdot)$  in Model 4a. Specifically, the intensity function can be

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\alpha(t - T_{N(t-)})^{\alpha-1} \exp\{\beta'Z\}. \quad (4b)$$

## 4.2 Estimation Procedures

The estimation procedures for the extended renewal process model are based on likelihood function  $L(\theta|data)$ , where  $\theta$  denotes all the unknown parameters in our model. The following presents estimation procedures with the semi-parametric model (Model 4a) and the parametric model (Model 4b).

### 4.2.1 Estimation with a Semi-parametric Model

With Model 4a, the unknown parameters  $\theta$  are specified as  $\beta, \lambda_0$ . The likelihood function becomes

$$L(\theta|data) = \prod_{i=1}^n \left\{ \prod_{k=1}^{K_i} \left[ \lambda_0(g_{ik}) e^{\beta' Z_i} \exp\left\{-\int_0^{g_{ik}} Y_i^R(u + t_{i,k-1}) \lambda_0(u) e^{\beta' Z_i} du\right\}\right] \right\} \\ \times \exp\left\{-\int_0^{C_i - t_{iK_i}} Y_i^R(u + t_{iK_i}) \lambda_0(u) e^{\beta' Z_i} du\right\}, \quad (4.1)$$

with the log-likelihood function

$$l(\theta|data) = \sum_{i=1}^n \left\{ \sum_{k=1}^{K_i} \left[ \log(\lambda_0(g_{ik})) + \beta' Z_i - \int_0^{g_{ik}} Y_i^R(u + t_{i,k-1}) \lambda_0(u) e^{\beta' Z_i} du \right] \right\} \\ - \int_0^{C_i - t_{iK_i}} Y_i^R(u + t_{iK_i}) \lambda_0(u) e^{\beta' Z_i} du, \quad (4.2)$$

where  $g_{ik} = t_{ik} - t_{i,k-1}$  with  $t_{i0} = 0$  are the gap times between the  $(k-1)$ th and  $k$ th EDMH visits for all subjects  $i = 1, \dots, n$ .

Let  $0 < g_1 < \dots < g_J$  be the distinct values of the gap times  $\{g_{ik} : k = 1, \dots, K_i; i = 1, \dots, n\}$ . Following Breslow (1972), we attain the MLE of  $\beta, \lambda_0$  by maximizing log-likelihood function  $l(\theta|data)$  in (4.2), viewing  $\lambda_0(g) = 0$  except for  $g = g_j, j = 1, \dots, J$ . This motivates the likelihood estimating equations

$$\partial l(\theta|data) / \partial \lambda_0(g_j) = 0, \quad j = 1, \dots, J; \quad \partial l(\theta|data) / \partial \beta = 0.$$

Let  $g_j + t_{i,k-1} = t_{ik}^{(j)}$  for  $k = 1, \dots, K_i$ . The first set of estimating equations can be written as, for  $j = 1, \dots, J$ ,

$$\lambda_0(g_j) = \frac{\sum_{i=1}^n \sum_{k=1}^{K_i} I(g_{ik} = g_j)}{\sum_{i=1}^n \sum_{k: g_j \in (0, g_{ik}] } Y_i^R(t_{ik}^{(j)}) e^{\beta' Z_i} + \sum_{i: g_j \in (0, C_i - t_{iK_i}] } Y_i^R(t_{i, K_i+1}^{(j)}) e^{\beta' Z_i}}. \quad (4.3)$$

Plugging it in the second estimating equations leads to

$$\begin{aligned} & \sum_{i=1}^n \sum_{k=1}^{K_i} \frac{\partial \beta'}{\partial \beta} Z_i \\ = & \sum_{j=1}^J \frac{\partial \beta'}{\partial \beta} \left\{ \frac{\sum_{i=1}^n \sum_{k: g_j \in (0, g_{ik}]} Y_i^R(t_{ik}^{(j)}) Z_i e^{\beta' Z_i} + \sum_{i: g_j \in (0, C_i - t_{iK_i}]} Y_i^R(t_{i, K_i+1}^{(j)}) Z_i e^{\beta' Z_i}}{\sum_{i=1}^n \sum_{k: g_j \in (0, g_{ik}]} Y_i^R(t_{ik}^{(j)}) e^{\beta' Z_i} + \sum_{i: g_j \in (0, C_i - t_{iK_i}]} Y_i^R(t_{i, K_i+1}^{(j)}) e^{\beta' Z_i}} \right\} \end{aligned} \quad (4.4)$$

A natural algorithm to attain the MLE of  $\beta$ ,  $\Lambda_0(t) = \int_0^t \lambda_0(u) du$  is as follows:

- (i) Obtain the MLE of  $\beta$ , denoted by  $\hat{\beta}$ , by solving the equations (4.4).
- (ii) Plug  $\hat{\beta}$  in the right-hand-sides of the equations (4.3), resulting in an estimator of  $\lambda_0(g_j)$  (denoted by  $\hat{\lambda}_0(g_j; \hat{\beta})$  in the following), and obtain the MLE of  $\Lambda_0(\cdot)$  as  $\hat{\Lambda}_0(t) = \sum_{j: g_j \leq t} \hat{\lambda}_0(g_j; \hat{\beta})$  for  $t > 0$ .

## 4.2.2 Estimation with a Parametric Model

With Model 4b, the unknown parameters  $\theta$  are specified as  $(\alpha, \beta)$ . The likelihood function becomes

$$\begin{aligned} L(\theta|data) &= \prod_{i=1}^n \prod_{k=1}^{K_i} (\alpha (g_{ik})^{\alpha-1} e^{\beta' Z_i}) \exp\left\{-\int_0^{g_{ik}} Y_i^R(t_{i, k-1} + u) \alpha u^{\alpha-1} e^{\beta' Z_i} du\right\} \\ &\quad \exp\left\{-\int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u) \alpha u^{\alpha-1} e^{\beta' Z_i} du\right\}, \end{aligned} \quad (4.5)$$

with the log-likelihood function

$$\begin{aligned} l(\theta|data) &= \sum_{i=1}^n \sum_{k=1}^{K_i} \log(\alpha) + (\alpha - 1) \log(g_{ik}) + \beta' Z_i - \int_0^{g_{ik}} Y_i^R(t_{i, k-1} + u) \alpha u^{\alpha-1} e^{\beta' Z_i} du \\ &\quad - \int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u) \alpha u^{\alpha-1} e^{\beta' Z_i} du. \end{aligned} \quad (4.6)$$

We attain the MLE of  $\theta$  by maximizing log-likelihood function  $l(\theta|data)$  in (4.6). The formulas of the likelihood score function and the observed information matrix are listed in Appendix C.

### 4.3 Analysis Results

Applying the Newton-Raphson algorithm, we obtained the MLE of  $\theta$  with the four time-independent covariates, with risk set adjusted for hospital duration, for Model 4a and Model 4b, respectively.

Table 4.1 represents the MPLEs of the regression coefficients in Model 4a (the extended semi-parametric renewal process model). The estimated standard errors are given below the estimates. The coefficient estimates with significant effect are bold in the table. We can see from Table 4.1 that there exist statistically significant effects of pSES, age, and sex, on the recurrence of EDMH visit, but no statistically significant difference is detected for region under Model 4a. Although the signs of the regression coefficients in Model 4a are the same with Model 3c in Table 3.1 (the AG model with time-independent covariate effects), the interpretations for the two models are different. The results under Model 4a indicate that a decreased risk of an EDMH visit between successive recurrent EDMH visits, is associated with pSES of O vs pSES of ASW, and males vs females, no matter when and how many repeated EDMH visits those subjects had; while older subjects at index initial visit tend to have higher risk to repeated EDMH visits than younger subjects. There is no statistically significant difference between those in urban and rural regions in the risk to the next EDMH visit from the time of most recent EDMH visit.

Table 4.1: MLEs of the regression coefficients in Model 4a.

Model	pSES (O vs ASW)	Age (at index initial EDMH)	Sex (Male vs Female)	Region (Urban vs Rural)
Estimates	<b>-0.368</b>	<b>0.194</b>	<b>-0.062</b>	0.107
SE	0.067	0.035	0.028	0.187

Under Model 4a, we estimate the cumulative intensity function for a subgroup of subjects who have pSES in the category of O, are of age 11 at the index initial EDMH visit, and live in urban area. See Figure 4.1. Figure 4.1(A) is for individuals who had EDMH visit at time  $T_1 = 1000, T_2 = 2000, T_3 = 2500$ . That is, the gap times between two successive EDMH visits are  $g_1 = 1000, g_2 = 1000, g_3 = 500$ . Figure 4.1(B) is for individuals who had EDMH visit at time  $T_1 = 500, T_2 = 1000, T_3 = 2000$ . The gap times between two successive EDMH visits for the individual are  $g_1 = 500, g_2 = 500, g_3 = 1000$ . We can see that the pattern is quite different with the one under the Cox

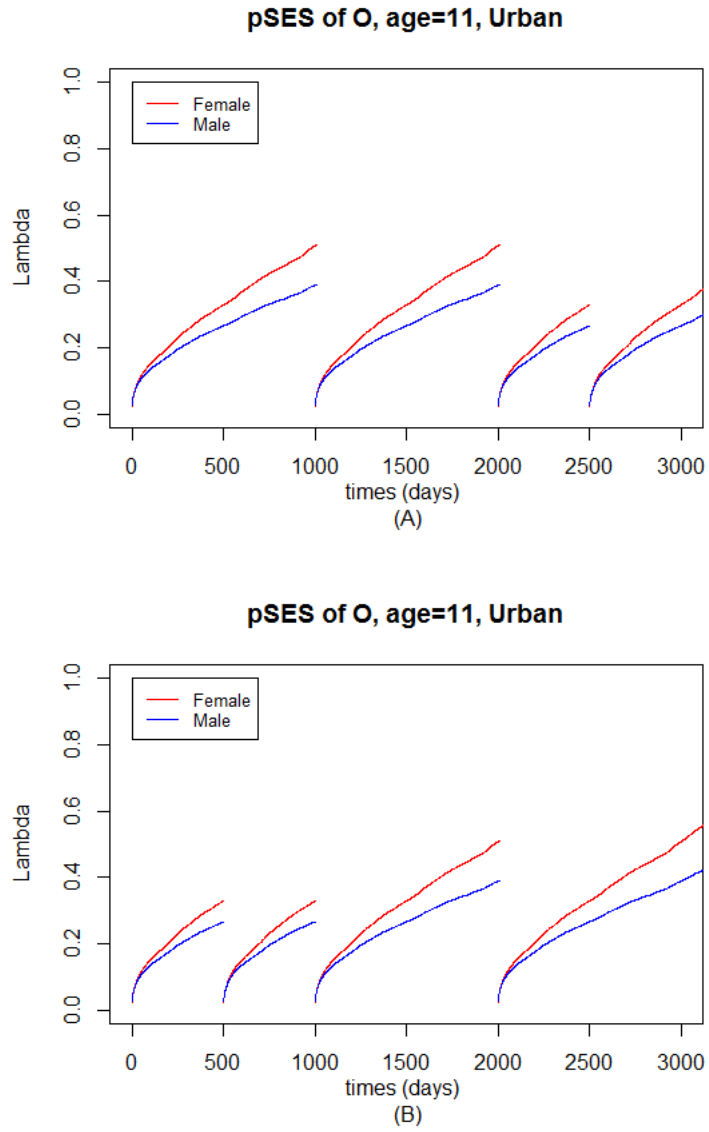


Figure 4.1: The estimates of the cumulative intensity function for an individual who had EDMH visit at time (A)  $T_1 = 1000, T_2 = 2000, T_3 = 2500$ , and (B)  $T_1 = 500, T_2 = 1000, T_3 = 2000$ .

regression model. Under the renewal process model, one individual returns back to the status at the index initial EDMH visit immediately after the occurrence of each EDMH visit.

Table 4.2 represents the MPLEs of the regression coefficients in the extended parametric renewal process model (Model 4b). The estimated standard errors are given below the estimates. The coefficient estimates with significant effect are bold in the table. We see from Table 4.2 that covariate effects with the parametric renewal process model (Model 4b) are in agreement with the covariate effects under the semi-parametric renewal process model (Model 4a).

Table 4.2: MLEs of the regression coefficients in Model 4b.

Model	$\alpha$	pSES (O vs ASW)	Age (at index initial EDMH)	Sex (Male vs Female)	Region (Urban vs Rural)
Estimates	<b>0.413</b>	<b>-0.206</b>	<b>0.104</b>	<b>-0.041</b>	0.062
SE	0.011	0.021	0.013	0.011	0.23

We estimate the cumulative rate of risk to the next EDMH visit since the index initial EDMH visit under the two renewal process models. Figure 4.2 shows the corresponding estimates within different subgroups. Dashed lines represent Model 4a, solid lines stand for Model 4b. We can see from Figure 4.2 that the estimates under the two renewal process models show closeness in some subgroups, such as pre-school males with pSES of O in urban; while in some subgroups, they show discrepancy. The parametric assumption for the baseline intensity function might be inappropriate.



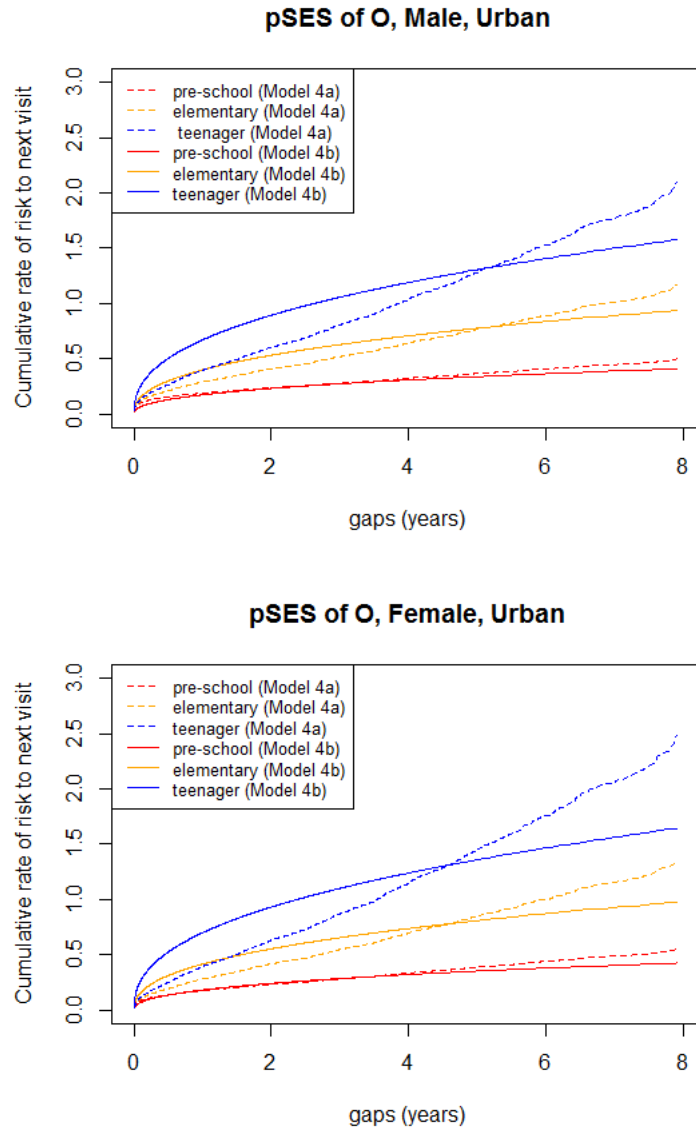


Figure 4.2: The estimates of the cumulative rates of risk to the next EDMH visit since the index initial EDMH visit under Model 4a and Model 4b. The lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively.

#### **4.4 Comparison with the Stratified Cox Regression Model**

The main difference between the two classes of models, the stratified Cox regression model and the renewal process model, is the index time for the baseline intensity function. The Stratified Cox regression model uses the time since the index initial EDMH visit, while the renewal process model uses the time since the most recent EDMH visit. As aforementioned in Chapter 1, the available data utilized in this project only contains records of mental health related ED visits from children and youth during April 1 2002 to March 31 2011. The starting point of one individual is from the time when he/she had the first EDMH visit during the observation window. The first EDMH visit observed within this observation window may not be the first EDMH visit of the whole life of one subject. With the renewal process model, we can more easily to interpret the analysis outcome. However, the objective of this project is to evaluate the frequency of EDMH visits and identifying the risk factors. It is worth examining the data through a variety of other models.

## Chapter 5

# Final Remarks

### 5.1 Summary

This project aims to evaluate the frequency of children and youth ED visits for mental illnesses and identify the associated risk factors. Based on the PMHC data, we focus on four main potential risk factors: social-economic status, gender, age at the index initial EDMH visit, and the region of residence. Adapting the counting process framework, we conduct data analyses under two classes of the generalized Cox regression model which accommodates the duration of hospitalization using the adjusted risk sets: the stratified Cox regression model and the renewal process model. The data analysis shows that all four risk factors have significant effects on the recurrence of EDMH visits. We found that an increasing risk of having recurrent EDMH visit is associated with subjects having index initial visit as teenagers versus at younger ages; children and youth who live in urban area are more likely to have multiple EDMH visits for mental illness than those in rural area; while those with pSES of category O tend to have lower risk of the next EDMH visit than those with other pSES, and a decreased multiple EDMH visit risk is associated with males versus females.

With a parametric version of the generalized Cox regression model, we present the maximum likelihood based estimation procedures in Chapter 2 and obtain the MLEs of the regression parameters with the risk set adjusted and not-adjusted for hospital duration. The estimation results under the two at-risk settings are very similar, which may due to a small proportion of children and

youth having observed hospitalization records. We evaluate the cumulative intensity function with the MLEs of the regression coefficients, and plot the log-transformed estimates against log of event time to check the model assumption. The result indicates that it is inappropriate to assume the baseline intensity function as a power function of event time. The differences of the Nelson-Aalen estimates for the intensity functions within different subgroups motivate the analysis in Chapter 3 and Chapter 4.

In Chapter 3, we consider the stratified Cox regression models and obtain the MPLEs of the regression coefficients. The robust estimates for the standard errors are very small. We also obtain the non-parametric bootstrap estimates for the standard errors, to make a comparison. The analysis results show that the effects of the four risk factors are different across different strata and there are statistical significant time trend for the effect of pSES, age, and gender. The statistically significant effect from the interaction between age and region was also detected.

In Chapter 4, we fit a semi-parametric and a parametric renewal process model based on PMHC data, respectively. We obtain the MLEs of the regression parameters. The interpretations under the renewal process models are different with those under the Cox regression models. The analysis results show that only three risk factors, pSES, age at the index initial EDMH visit, and sex, have statistically significant effects on the occurrence of the next EDMH visit since the most recent visit. The results also indicate that the parametric assumption for the baseline intensity function might be inappropriate.

## 5.2 Future Investigation

There are many issues to consider. Some interesting points of future investigations are listed below.

- The estimation procedures in this project require the censoring times  $C_i$  for all study subjects. The censoring times depend on the unavailable birth dates of the subjects. We can use different methods to estimate the birth dates of the study subjects.
- There is also an interest in the consideration of time-varying covariates. When time-varying covariates are of interest, the estimation procedures need observations of the covariates

throughout the whole study period. In the current application, we only have the covariate information at subject's each own EDMH visit times. One approach to deal with this issue is to assume the time-dependent covariates of one subject have only jumps at his or her observed EDMH visit times.

- We are also interested in considering different stratification variables. We explore season of the EDMH visits as a stratification variable in this project. This stratification variable may not be informative, as winter in Calgary and Edmonton, for example, may start in different month. It is interesting to consider temperature at EDMH visit time as the stratification variable, if data is available.
- Different methods can be used to accommodate the informative time origin problem. For example, one may use age in years as the time scale and set age of zero as the starting point of each subject.
- The PMHC data were extracted from four population-based administrative database. The information is only available for children and youth in Alberta who had ED records for mental health during the data collection period from April 01, 2002 to March 31, 2011. In addition, we know that those individuals who are eligible to be recorded but have no information about ED visits recorded in the database must not have any ED visit during the data collection period. Utilizing available demographic information on the whole population of children and youth in Alberta who are eligible to be recorded in the database to supplement the available ED visit data, we can evaluate the ED visit frequency of the whole population of children and youth in Alberta.
- The ED visit data and non-ED physician visit data allow us to study the relationship of ED and non-ED visits. We can jointly model the ED and non-ED visits together by assuming different forms of the intensity functions. It would also be interesting to jointly model the time of ED visit (time-to-event data) and diagnosis or severity of ED visit (longitudinal data).

# Bibliography

- [1] P.K. Anderson, O. Borgan, R.D. Gill, and N. Keiding. *Statistical Models Based on Counting Processes*. Springer, New York, 1993.
- [2] P.K. Anderson and R.D. Gill. Cox's regression model for counting processes: A large sample study. *Annals Statistics*, 10:1100–1120, 1982.
- [3] N.E. Breslow. Discussion following "regression models and life-tables" by d.r. cox. *Journal of the Royal Statistical Society*, 1972.
- [4] Deborah Burr. A comparison of certain bootstrap confidence intervals in the cox model. *Journal of the American Statistical Association*, 89:1290–1302, 1994.
- [5] R.J. Cook and J. Lawless. *The Statistical Analysis of Recurrent Events*. Springer, 2007.
- [6] D.R. Cox. Regression models and life-tables. *Journal of the Royal Statistical Society*, 1972.
- [7] D.R. Cox. Partial likelihood. *Biometrika*, 62:269–276, 1975.
- [8] T.R. Fleming and D.P. Harrington. *Counting Processes and Survival Analysis*. John Wiley and Sons, Inc., 1991.
- [9] D.W. Hosmer, S. Lemeshow, and S. May. *Applied Survival Analysis*. John Wiley and Sons, Inc., 2008.
- [10] X.J. Hu, M. Lorenzi, J.J. Spinelli, S.C. Ying, and M.L. McBride. Analysis of recurrent events with non-negligible event duration, with application to assessing hospital utilization. *Lifetime Data Analysis*, 17:215–233, 2011.
- [11] X.J. Hu, J. Sun, and L. Wei. Regression parameter estimation from panel counts. *Scandinavian Journal of Statistics*, 30:25–43, 2003.
- [12] J.D. Kalbfleisch and R.L. Prentice. *The Statistical Analysis of Failure Time Data*. John Wiley and Sons, Inc., 2002.
- [13] J.F. Lawless. *Statistical Models and Methods for Lifetime Data*. John Wiley and Sons, Inc., 2003.
- [14] K.K. Leitch. Reaching for the top: A report by the advisor on healthy children and youth: Ottawa, Ontario: Health Canada, 2007.
- [15] D.Y. Lin and L.J. Wei. The robust inference for the cox proportional hazards model. *Journal of the American Statistical Association*, 84:1074–1078, 1989.

- [16] A.S. Newton, S. Ali, D.W. Johnson, C. Haines, R.J. Rosychuk, R.A. Keaschuk, P. Jacobs, M. Cappelli, and T.P. Klassen. Who comes back? characteristics and predictors of return to emergency department services for pediatric mental health emergencies. *Academic Emergency Medicine*, 17:177–186, 2010.
- [17] A.S. Newton and R.J. Rosychuk. *The Emergency Department Compass: Children's Mental Health*. Pediatric mental health emergencies in Alberta, Canada: Emergency department visits by children and youth aged 0 to 17 years, 2002-2008. Edmonton, AB, 2011.
- [18] R.L. Prentice, B.J. Willians, and A.V. Peterson. On the regression analysis of multivariate failure time data. *Biometrika*, 68:373–379, 1981.
- [19] T.M. Therneau and P.M. Grambsch. *Modeling Survival Data: Extending the Cox Model*. Springer, New York, 2000.

## **Appendix A**

# **Estimates of the Cumulative Intensity Functions**

This section represent the estimates of the cumulative intensity functions. We group subjects according to age at the index initial EDMH visit (pre-school 0 ~ 5, elementary school 6 ~ 13, and teenager (14 ~ 17), sex (male and female), pSES at the index initial EDMH visit (O and ASW), and region at the index initial EDMH visit (urban and rural). For each of the different subgroups, we evaluate the generalized Nelson-Aalen estimator and the Breslow estimator under the Cox regression models in Chapter 3 for the cumulative intensity function.



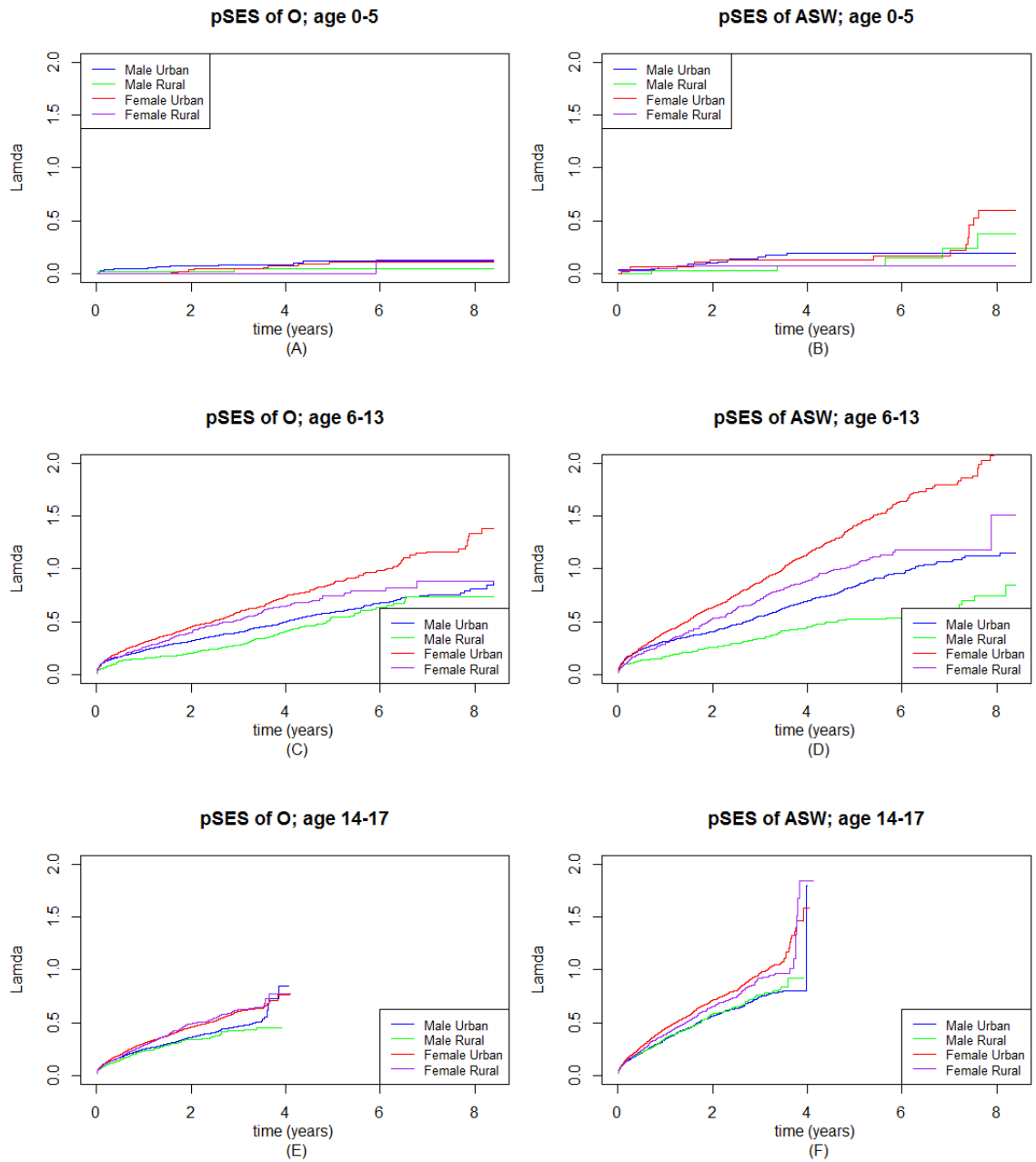


Figure A.1: The Generalized Nelson-Aelon estimates of the cumulative intensity functions with adjustment of hospital duration.

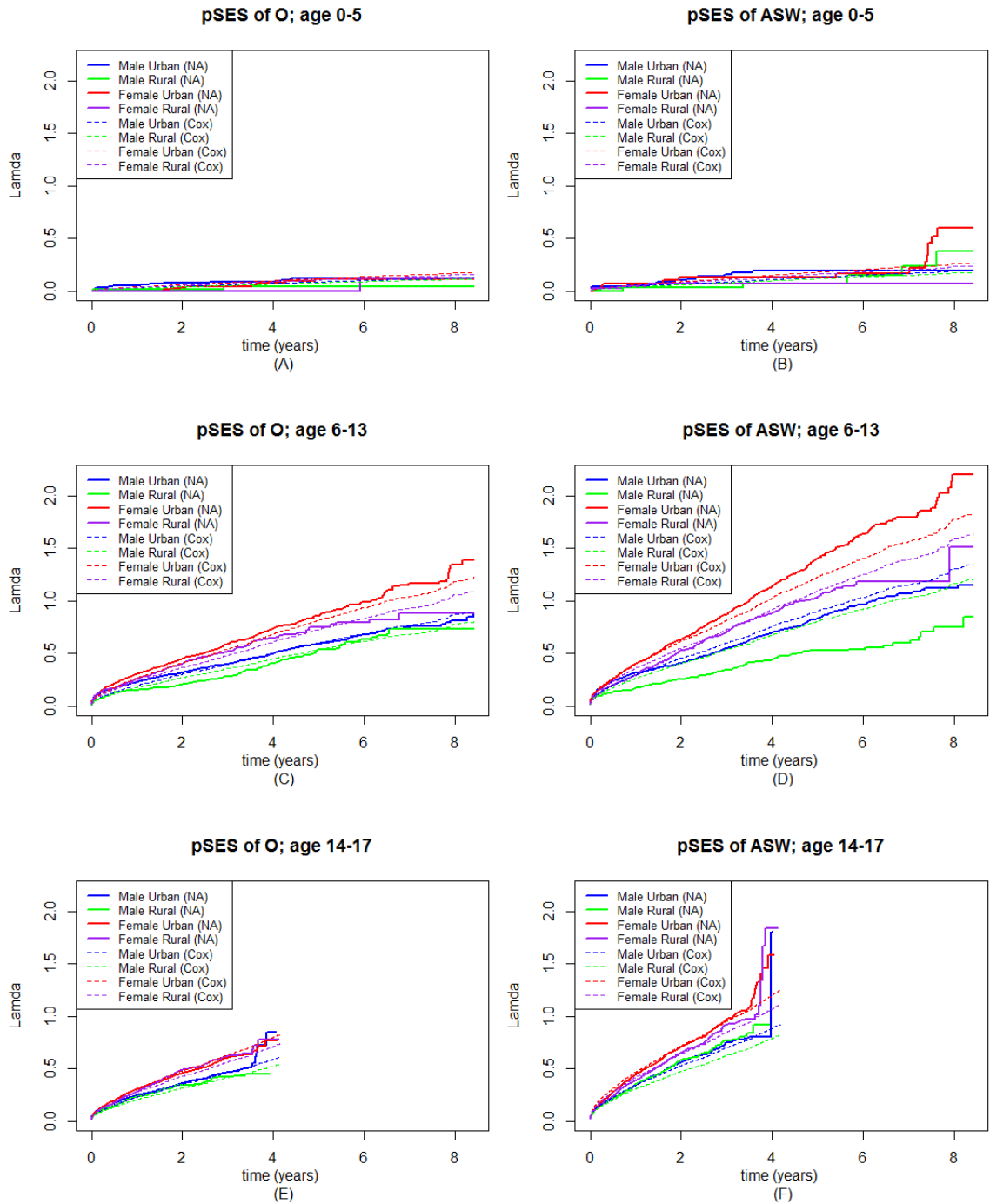


Figure A.2: The Generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration, together with the Breslow estimates under Model 3c (the model considering time-independent covariate effects) in Table 3.1. The dashed lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively.

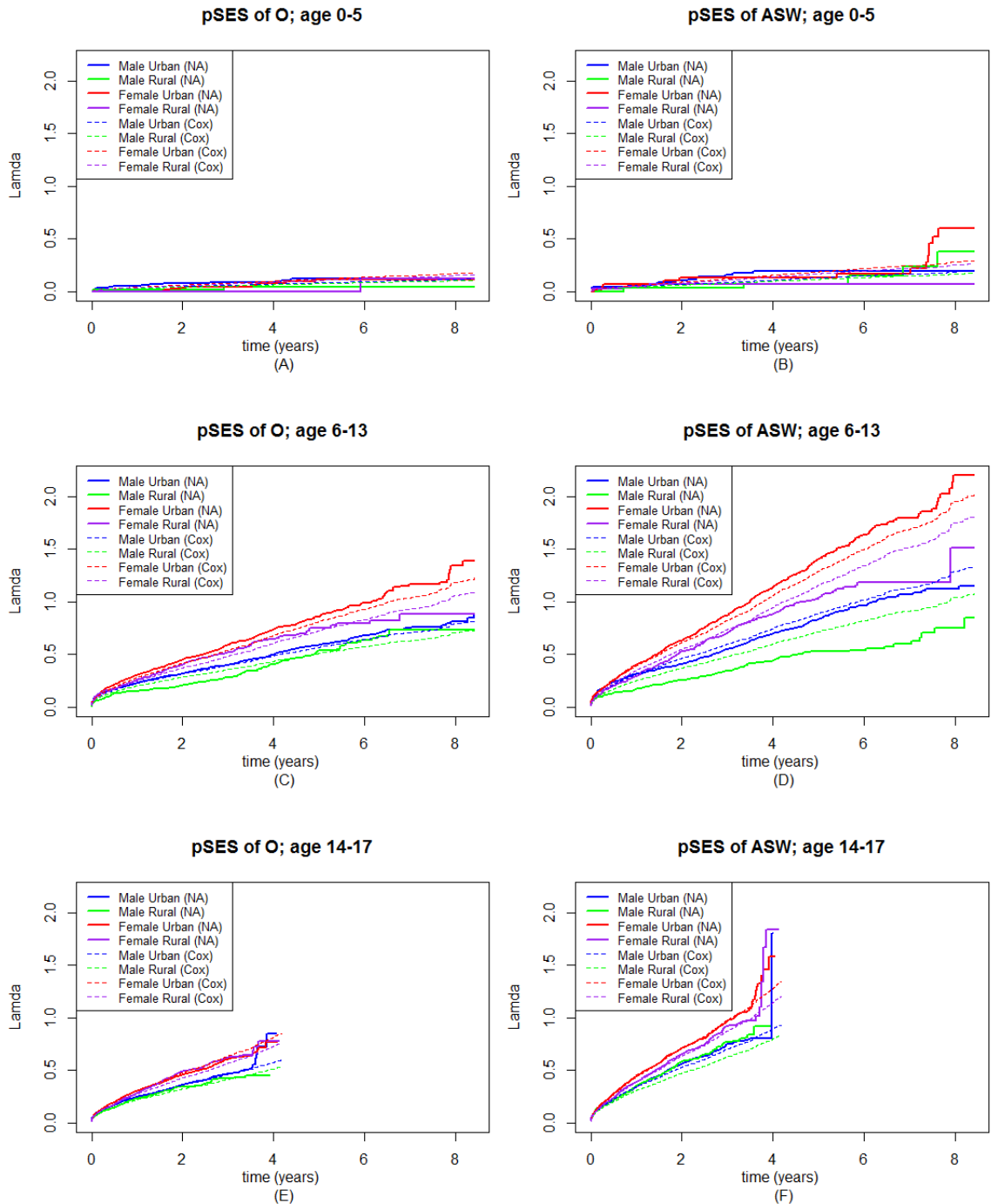


Figure A.3: The Generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration, together with the Breslow estimates under Model 3c (the model considering time-dependent covariate effects) in Table 3.2. The dashed lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively.

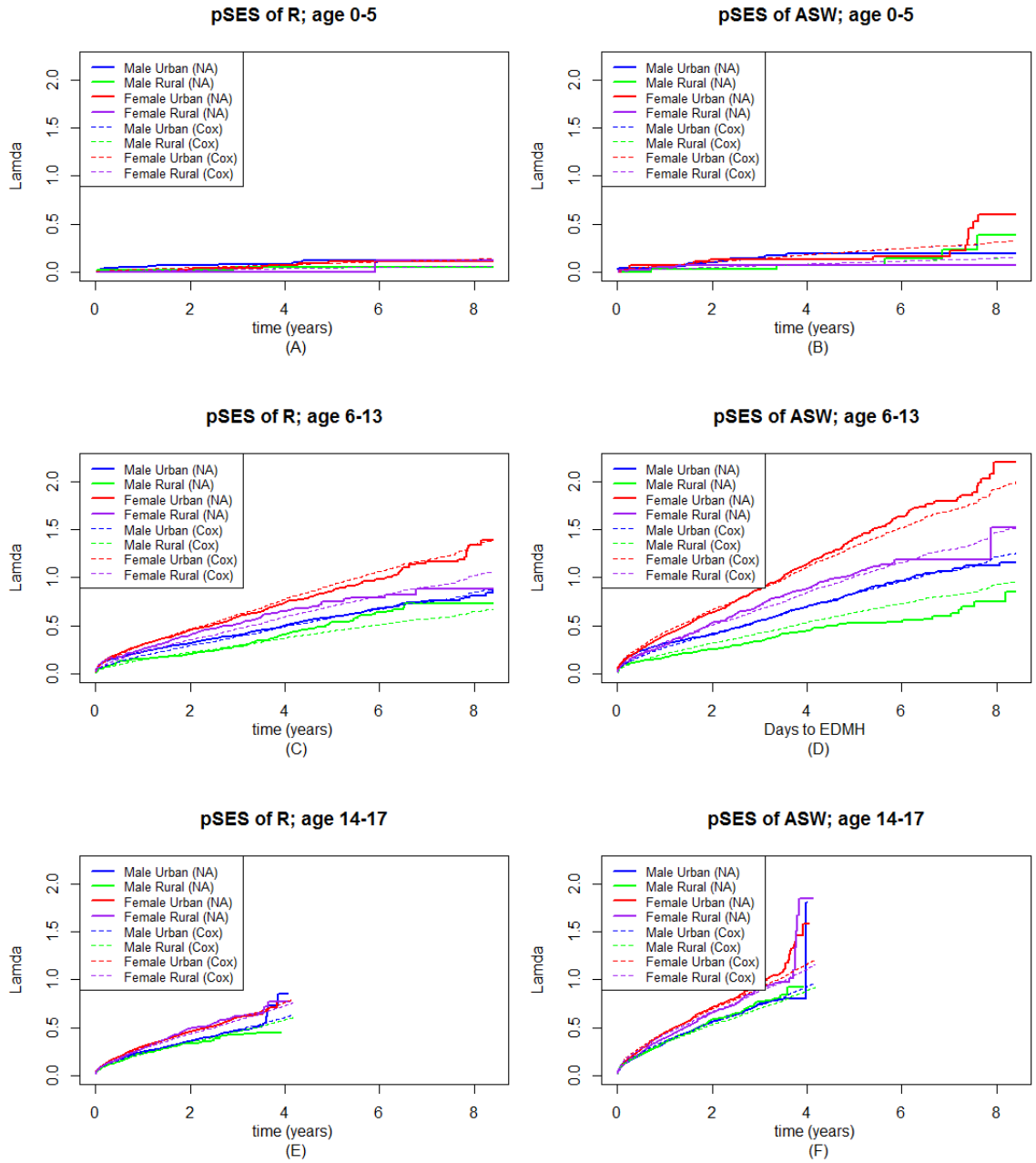


Figure A.4: The Generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration, together with the Breslow estimates under Model 3c (the two factor interaction model with time-independent covariate effects) in Table 3.5. The dashed lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively.

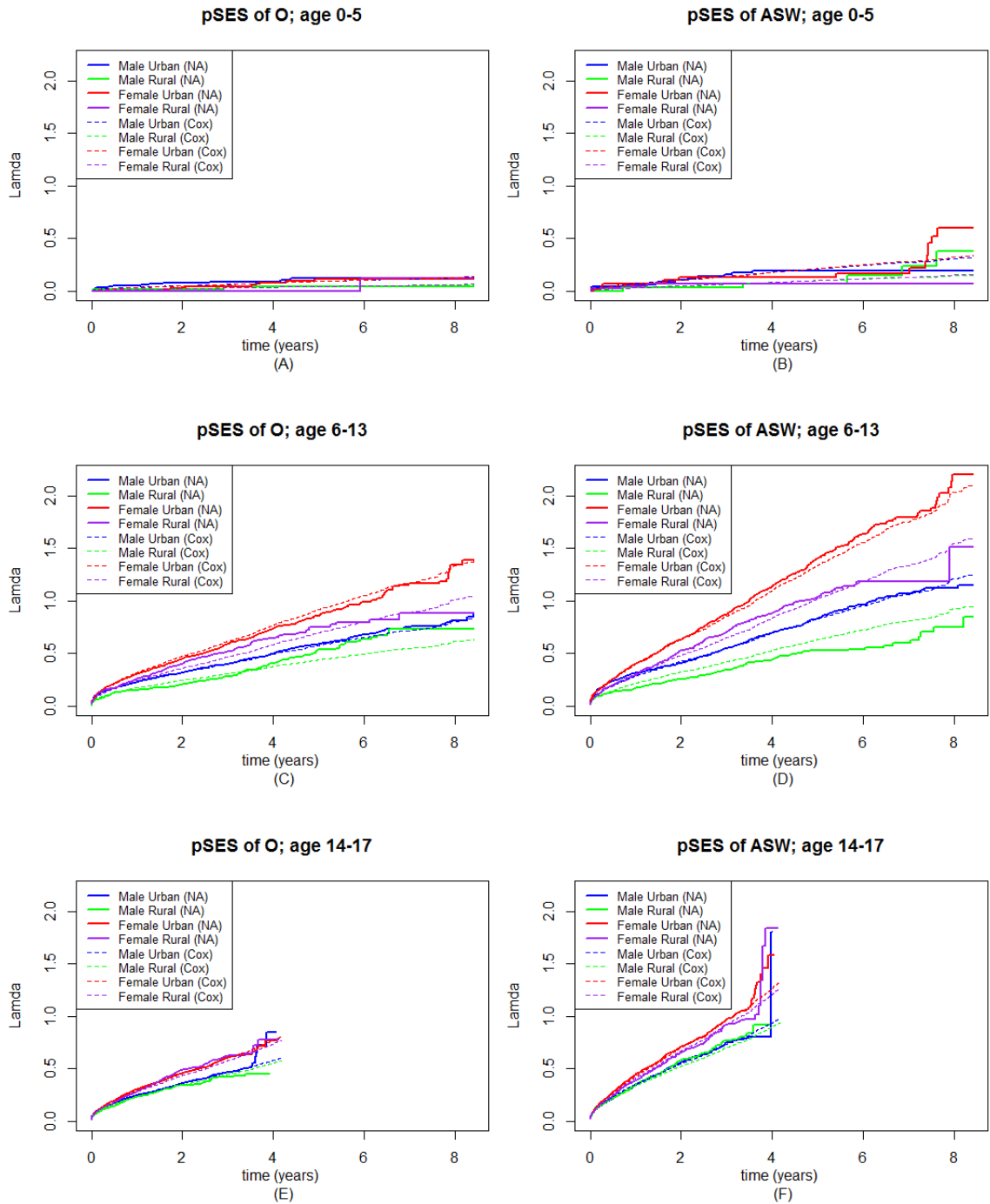


Figure A.5: The Generalized Nelson-Aalen estimates of the cumulative intensity functions with adjustment of hospital duration, together with the Breslow estimates under Model 3c (the two factor interaction model with time-dependent covariate effects) in Table 3.6. The dashed lines are plotted with Age=3, Age=11, and Age=16 for the three age group 0-5, 6-13, 14-17, respectively.

## **Appendix B**

# **Non-parametric Bootstrap Estimates of Standard Errors**

This section represent the estimations under models in Chapter 3, along with the non-parametric bootstrap estimates of the standard errors. We explored bootstrap sizes  $B=1000$ ,  $B=2000$ , and  $B=5000$ , respectively. The estimates for the standard errors are similar when  $B=2000$  and  $B=5000$ . Thus, we shows the estimation results with  $B=2000$  in this section.

Table B.1 corresponds to Table 3.1; Table B.2 can be compared with Table 3.2; while Table B.3 corresponds to Table 3.3.

Table B.1: MPLEs of the regression coefficients and nonparametric bootstrap estimates for the standard errors in Model 3a/3b/3c with time-independent covariate effects.

Model	Stratification Variable 1*			Stratification Variable 2**			AG Model
	$\beta_1$	$3a_{\beta_s, \lambda_{0s}}^{stage}$	$3b_{\beta_s, \lambda_{0s}}^{stage}$	$\beta_1$	$3a_{\beta_s, \lambda_{0s}}^{seasons}$	$3b_{\beta_s, \lambda_{0s}}^{seasons}$	
$\log(PL(\hat{\beta}))$							
		$\beta_3$	$\beta(=\beta_s)$		$\beta_4$	$\beta(=\beta_s)$	$\beta(t; \beta_s) = \beta$
pSES	<b>-0.768</b> (0.281)	<b>-0.375</b> (0.065)	<b>-0.418</b> (0.041)	<b>-0.451</b> (0.044)	<b>-0.364</b> (0.046)	<b>-0.385</b> (0.022)	<b>-0.420</b> (0.027)
Age	<b>0.662</b> (0.126)	<b>0.137</b> (0.010)	-0.009 (0.023)	<b>0.033</b> (0.008)	<b>0.122</b> (0.008)	<b>0.085</b> (0.008)	<b>0.085</b> (0.005)
Sex	-0.265 (0.289)	<b>-0.379</b> (0.043)	<b>-0.223</b> (0.031)	<b>-0.254</b> (0.057)	<b>-0.399</b> (0.044)	<b>-0.213</b> (0.025)	<b>-0.282</b> (0.026)
Region	0.562 (0.337)	<b>0.279</b> (0.048)	0.044 (0.033)	<b>0.178</b> (0.051)	<b>0.129</b> (0.046)	<b>0.113</b> (0.022)	<b>0.126</b> (0.032)
		$3a(3)$	$3b(3)$				
$\log(PL(\hat{\beta}))$							
		-112511.3	-112546.1				
pSES	<b>-0.847</b> (0.268)	<b>-0.360</b> (0.038)	<b>-0.418</b> (0.040)				
Sex	0.009 (0.261)	<b>-0.460</b> (0.041)	<b>-0.224</b> (0.027)				
Region	<b>0.746</b> (0.344)	<b>0.273</b> (0.049)	0.045 (0.031)				

\* Age at 1st EDMH visit as the stratification variable (pre-school 0-5, elementary school 6-13, and teenager 14-17).

\*\* Seasons of the EDMH visits as the stratification variable: fall, spring, summer, and winter, respectively.

pSES as the indicator of O, Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban.

Non-parametric bootstrap standard errors in brackets (B=2000);

Significant effect with p-value  $\leq 0.05$  in boldface

Table B.2: MPLEs of the regression coefficients and nonparametric bootstrap estimates for the standard errors in Model 3a/3b/3c with time-dependent covariate effects.

Model	Stratification Variable 1*			Stratification Variable 2**			AG Model 3c $\beta(t; \beta_s) = \beta$ -119359.737
	$\beta_1$	$3a_{\beta_s, \lambda_{0s}}^{stage}$ $\beta_2$	$\beta_3$	$\beta_1$	$3a_{\beta_s, \lambda_{0s}}^{seasons}$ $\beta_2$	$\beta_3$	
$\log(PL(\hat{\beta}))$	-112311.3	-112389.7	-102201.3	-102277.6	-102277.6	-102277.6	-102277.6
pSES	-0.421 (0.781)	-0.088 (0.1113)	-0.024 (0.071)	-0.032 (0.110)	-0.085 (0.131)	-0.211 (0.122)	-0.088 (0.064)
Age	0.503 (0.321)	0.036 (0.039)	-0.018 (0.038)	-0.021 (0.025)	-0.003 (0.031)	-0.001 (0.038)	-0.011 (0.013)
Sex	1.390 (0.958)	0.051 (0.124)	0.024 (0.084)	0.022 (0.101)	0.040 (0.113)	0.188 (0.132)	0.054 (0.063)
Region	0.359 (0.941)	<b>0.344</b> (0.139)	0.125 (0.074)	0.061 (0.122)	0.099 (0.120)	<b>0.358</b> (0.152)	<b>0.131</b> (0.073)
pSES $\times$ ln(t)	-0.140 (0.295)	<b>-0.116</b> (0.041)	<b>-0.189</b> (0.038)	<b>-0.191</b> (0.052)	<b>-0.123</b> (0.061)	<b>-0.136</b> (0.057)	<b>-0.150</b> (0.041)
Age $\times$ ln(t)	0.066 (0.125)	<b>0.042</b> (0.026)	0.004 (0.032)	<b>0.025</b> (0.007)	<b>0.058</b> (0.010)	<b>0.050</b> (0.015)	<b>0.045</b> (0.005)
Sex $\times$ ln(t)	-0.641 (0.351)	<b>-0.174</b> (0.052)	<b>-0.121</b> (0.033)	<b>-0.127</b> (0.042)	<b>-0.195</b> (0.048)	<b>-0.198</b> (0.052)	<b>-0.151</b> (0.032)
Region $\times$ ln(t)	0.074 (0.361)	-0.026 (0.062)	-0.039 (0.043)	0.053 (0.051)	0.012 (0.052)	-0.055 (0.060)	-0.004 (0.031)

\* Age at the index initial EDMH visit as the stratification variable (pre-school 0-5, elementary school 6-13, and teenager 14-17).

\*\* Seasons of the EDMH visits as the stratification variable: fall, spring, summer, and winter, respectively.

pSES as the indicator of O, Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban.

Non-parametric bootstrap standard error in brackets (B=2000);

Significant effect with p-value  $\leq 0.05$  in boldface



Table B.3: MPLEs of the regression coefficients and nonparametric bootstrap estimates for the standard errors in Model 3a/3b/3c with all pairs of two factor interactions and time-independent covariate effects.

Model	Stratification Variable 1*				Stratification Variable 2**				AG Model
	$\beta_1$	$\beta_2$	$\beta_3$	$\beta(\equiv \beta_s)$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$	
$\log(PL(\hat{\beta}))$	-112304.1	-112441.5	-102258.9	-119440.9	-0.466	-0.319	-0.304	-0.376	-0.367
pSES	<b>-9.949</b> (3.381)	<b>-0.278</b> (0.102)	<b>-0.392</b> (0.064)	<b>-0.362</b> (0.058)	<b>-0.466</b> (0.078)	<b>-0.319</b> (0.089)	<b>-0.304</b> (0.091)	<b>-0.376</b> (0.086)	<b>-0.367</b> (0.060)
Age	1.151	<b>0.818</b>	<b>-0.224</b>	<b>0.356</b>	<b>0.187</b>	<b>0.375</b>	<b>0.415</b>	<b>0.389</b>	<b>0.333</b>
Sex	2.099	<b>-0.924</b>	<b>-0.293</b>	<b>-0.323</b>	-0.138	<b>-0.368</b>	<b>-0.445</b>	<b>-0.385</b>	<b>-0.319</b>
Region	(2.915)	(0.100)	(0.061)	(0.063)	(0.081)	(0.085)	(0.106)	(0.090)	(0.063)
pSES×Age	<b>8.810</b> (3.170)	0.147 (0.094)	0.036 (0.054)	<b>0.143</b> (0.050)	<b>0.170</b> (0.067)	<b>0.275</b> (0.068)	<b>0.246</b> (0.071)	<b>-0.093</b> (0.070)	<b>0.154</b> (0.055)
pSES×Sex	<b>-1.978</b> (0.791)	-0.040 (0.062)	0.024 (0.077)	-0.029 (0.035)	-0.054 (0.041)	-0.037 (0.045)	-0.041 (0.051)	0.014 (0.044)	-0.031 (0.029)
pSES×Region	1.636	<b>-0.239</b>	-0.039	<b>-0.092</b>	-0.064	<b>0.247</b>	-0.042	0.081	0.056
Age×Sex	0.488	<b>-0.419</b>	<b>0.292</b>	0.038	(0.084)	(0.079)	(0.089)	(0.081)	(0.058)
Age×Region	<b>2.159</b> (0.845)	<b>-0.223</b> (0.061)	0.111 (0.068)	<b>-0.120</b> (0.041)	<b>-0.118</b> (0.045)	<b>-0.102</b> (0.040)	<b>-0.074</b> (0.055)	<b>-0.138</b> (0.047)	<b>-0.117</b> (0.029)
Sex×Region	-1.457 (0.830)	0.042 (0.096)	-0.017 (0.065)	0.016 (0.055)	-0.080 (0.097)	<b>-0.193</b> (0.085)	<b>0.254</b> (0.112)	<b>0.188</b> (0.091)	0.018 (0.063)

\* Age at the index initial EDMH visit as the stratification variable (pre-school 0-5, elementary school 6-13, and teenager 14-17).  
 \*\* Seasons of the EDMH visits as the stratification variable: fall, spring, summer, and winter, respectively.  
 pSES as the indicator of O, Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban.  
 Non-parametric bootstrap standard error in brackets (B=2000);  
 Significant effect with p-value  $\leq 0.05$  in boldface

Table B.4: MPLEs of the regression coefficients and nonparametric bootstrap estimates for the standard errors in Model 3a/3b/3c with all pairs of two factor interactions and time-dependent covariate effects.

Model	Stratification Variable 1*			Stratification Variable 2**			AG Model	
	$3a^{\text{stage}}$ $\beta_1, \beta_2, \beta_3$	$3b^{\text{stage}}$ $\beta(=\beta_s)$	$3c^{\text{stage}}$ $\beta(=\beta_s)$	$3a^{\text{seasons}}$ $\beta_1, \beta_2, \beta_3, \beta_4$	$3b^{\text{seasons}}$ $\beta(=\beta_s)$	$3c$ $\beta(t; \beta_s) = \beta$		
$\log(PL(\hat{\beta}))$	-112247.2	-112367.4	-112367.4	-102148.9	-102253.5	-119335.6		
pSES	<b>-10.136</b> (3.822)	0.011 (0.091)	0.099 (0.091)	0.017 (0.131)	0.081 (0.150)	-0.050 (0.138)	0.008 (0.078)	0.010 (0.078)
Age	0.824 (1.401)	<b>0.593</b> (0.122)	-0.127 (0.125)	0.073 (0.074)	<b>0.208</b> (0.098)	0.133 (0.085)	<b>0.117</b> (0.044)	<b>0.116</b> (0.043)
Sex	3.882 (3.202)	<b>-0.470</b> (0.154)	-0.083 (0.098)	0.026 (0.077)	0.106 (0.166)	-0.167 (0.142)	0.021 (0.077)	0.019 (0.078)
Region	<b>8.818</b> (3.409)	0.239 (0.165)	0.102 (0.104)	<b>0.248</b> (0.081)	<b>0.351</b> (0.149)	0.125 (0.149)	<b>0.258</b> (0.086)	<b>0.256</b> (0.085)
pSES×Age	<b>-2.029</b> (0.805)	-0.043 (0.061)	-0.132 (0.073)	<b>-0.072</b> (0.034)	-0.069 (0.055)	-0.007 (0.064)	<b>-0.065</b> (0.033)	<b>-0.065</b> (0.033)
pSES×Sex	<b>1.334</b> (0.633)	0.127 (0.073)	-0.027 (0.056)	0.025 (0.044)	<b>0.225</b> (0.084)	0.067 (0.088)	0.034 (0.042)	0.032 (0.042)
pSES×Region	1.647 (0.965)	<b>-0.242</b> (0.093)	-0.043 (0.064)	<b>-0.095</b> (0.049)	<b>-0.165</b> (0.091)	-0.036 (0.095)	<b>-0.102</b> (0.053)	<b>-0.101</b> (0.053)
Age×Sex	0.546 (0.778)	<b>-0.434</b> (0.064)	<b>0.226</b> (0.078)	-0.007 (0.031)	<b>0.142</b> (0.050)	-0.077 (0.058)	0.001 (0.028)	-0.000 (0.028)
Age×Region	<b>2.166</b> (0.854)	<b>-0.224</b> (0.080)	0.092 (0.081)	<b>-0.139</b> (0.033)	<b>-0.114</b> (0.058)	<b>-0.171</b> (0.065)	<b>-0.135</b> (0.037)	<b>-0.135</b> (0.037)
Sex×Region	-1.433 (0.843)	0.041 (0.098)	-0.020 (0.064)	0.011 (0.055)	<b>-0.192</b> (0.094)	0.177 (0.098)	0.012 (0.057)	0.014 (0.057)
pSES×ln(t)	0.001 (0.324)	<b>-0.115</b> (0.049)	<b>-0.211</b> (0.042)	<b>-0.169</b> (0.053)	<b>-0.133</b> (0.049)	<b>-0.146</b> (0.057)	<b>-0.167</b> (0.032)	<b>-0.168</b> (0.032)
Age×ln(t)	0.132 (0.389)	<b>0.099</b> (0.038)	0.012 (0.056)	<b>0.154</b> (0.034)	<b>0.170</b> (0.033)	<b>0.148</b> (0.035)	<b>0.125</b> (0.028)	<b>0.125</b> (0.027)
Sex×ln(t)	-0.615 (0.355)	<b>-0.186</b> (0.051)	<b>-0.088</b> (0.039)	<b>-0.154</b> (0.033)	<b>-0.140</b> (0.056)	-0.092 (0.060)	<b>-0.149</b> (0.033)	<b>-0.149</b> (0.034)
Region×ln(t)	-0.005 (0.408)	-0.037 (0.062)	-0.028 (0.045)	-0.047 (0.034)	-0.037 (0.070)	-0.097 (0.066)	-0.048 (0.041)	-0.048 (0.042)

\* Age at the index initial EDMH visit as the stratification variable (pre-school 0-5, elementary school 6-13, and teenager 14-17).

\*\* Seasons of the EDMH visits as the stratification variable: fall, spring, summer, and winter, respectively.

pSES as the indicator of O. Age at the index initial EDMH visit, Sex as the indicator of male, Region as the indicator of urban.

Non-parametric bootstrap standard error in brackets (B=2000);

Significant effect with p-value  $\leq 0.05$  in boldface

## Appendix C

# Estimation Formula

### The estimation formula for Model 2 in Chapter 2

Assume that the baseline intensity function is a power function of the event time  $t$ . We specify the baseline intensity into  $\alpha t^{\alpha-1}$ . The model is

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\alpha t^{\alpha-1} \exp\{\beta'Z\}. \quad (2)$$

Let  $\theta$  denote the two set of unknown parameters in Model 2. Then the likelihood function is

$$\begin{aligned} L(\theta|data) &\propto \prod_{i=1}^n \prod_{t \in (0, C_i]} (\lambda(t; \mathcal{H}_i(t)))^{dN_i(t)} (1 - \lambda(t; \mathcal{H}_i(t))dt)^{1-dN_i(t)} \\ &= \prod_{i=1}^n \prod_{t \in (0, C_i]} (\alpha t^{\alpha-1} \exp\{\beta'Z_i\})^{dN_i(t)} \times \exp \left\{ - \int_0^{C_i} Y_i^R(t) \alpha t^{\alpha-1} \exp\{\beta'Z_i\} dt \right\} \end{aligned}$$

The log likelihood function is

$$\begin{aligned} l(\theta|data) &= \sum_{i=1}^n \left[ \sum_{t \in (0, C_i]} \log(\alpha t^{\alpha-1} e^{\beta'Z_i}) dN_i(t) - \int_0^{C_i} Y_i^R(t) \alpha t^{\alpha-1} e^{\beta'Z_i} dt \right] \\ &= \sum_{i=1}^n \int_0^{\infty} Y_i^C(t) \left[ \log(\alpha t^{\alpha-1} e^{\beta'Z_i}) dN_i(t) - Y_i^R(t) \alpha t^{\alpha-1} e^{\beta'Z_i} dt \right] \end{aligned}$$

The likelihood score function of  $\boldsymbol{\theta}$  is  $U(\boldsymbol{\theta}) = (U_\alpha(\boldsymbol{\theta}), U_\beta'(\boldsymbol{\theta}))'$ , with

$$\begin{aligned} U_\alpha(\boldsymbol{\theta}) &= \sum_{i=1}^n \int_0^\infty Y_i^C(t) \left[ \left( \frac{1}{\alpha} + lnt \right) dN_i(t) - Y_i^R(t) \left( t^{\alpha-1} + \alpha t^{\alpha-1} lnt \right) e^{\beta' Z_i} dt \right] \\ &= \sum_{i=1}^n \int_0^\infty Y_i^C(t) \left( \frac{1}{\alpha} + lnt \right) \left[ dN_i(t) - Y_i^R(t) \alpha t^{\alpha-1} e^{\beta' Z_i} dt \right] \end{aligned}$$

and

$$U_\beta(\boldsymbol{\theta}) = \sum_{i=1}^n \int_0^\infty Y_i^C(t) Z_i \left[ dN_i(t) - Y_i^R(t) \alpha t^{\alpha-1} e^{\beta' Z_i} dt \right]$$

The observed information matrix is  $I(\boldsymbol{\theta})$ , with four elements

$$I_{\alpha\alpha}(\boldsymbol{\theta}) = \sum_{i=1}^n \int_0^\infty Y_i^C(t) \left[ \left( \frac{1}{\alpha^2} \right) dN_i(t) + Y_i^R(t) \left( t^{\alpha-1} (2 + \alpha lnt) lnt \right) e^{\beta' Z_i} dt \right]$$

$$I_{\alpha\beta}(\boldsymbol{\theta}) = I_{\beta\alpha}(\boldsymbol{\theta})' = \sum_{i=1}^n \int_0^\infty Y_i^C(t) Y_i^R(t) t^{\alpha-1} (1 + \alpha lnt) Z_i e^{\beta' Z_i} dt$$

$$I_{\beta\beta}(\boldsymbol{\theta}) = \sum_{i=1}^n \int_0^\infty Y_i^C(t) Y_i^R(t) \alpha t^{\alpha-1} Z_i Z_i' e^{\beta' Z_i} dt$$

## The estimation formula for Model 4b in Chapter 4

Under the assumption that the baseline intensity function is a power function of the event time  $t$ , the extended renewal process model becomes

$$\lambda(t|\mathcal{H}(t)) = Y^R(t)\alpha(t - T_{N(t-)})^{\alpha-1} \exp\{\beta'Z\}. \quad (4b)$$

Let  $\theta$  denote the two set of unknown parameters in Model 4b. The likelihood function is derived as the following

$$\begin{aligned} L(\theta|data) &= \prod_{i=1}^n L_i(\theta|data_i) \\ &\propto \prod_{i=1}^n \prod_{t \in (0, C_i]} (\lambda(t; \mathcal{H}_i(t)))^{dN_i(t)} (1 - \lambda(t; \mathcal{H}_i(t)))^{1-dN_i(t)} \\ &= \prod_{i=1}^n \prod_{t \in (0, C_i]} (Y_i^R(t)\alpha(t - T_{N(t-)})^{\alpha-1} e^{\beta'Z_i})^{dN_i(t)} (1 - Y_i^R(t)\alpha(t - T_{N(t-)})^{\alpha-1} e^{\beta'Z_i} dt)^{1-dN_i(t)} \\ &= \prod_{i=1}^n \left[ \prod_{t \in (0, C_i]} (Y_i^R(t)\alpha(t - T_{N(t-)})^{\alpha-1} e^{\beta'Z_i})^{dN_i(t)} \right] \\ &\quad \left[ \prod_{t \in (0, C_i]} (1 - Y_i^R(t)\alpha(t - T_{N(t-)})^{\alpha-1} e^{\beta'Z_i} dt)^{1-dN_i(t)} \right] \\ &= \prod_{i=1}^n \left[ \prod_{k=1}^{K_i} (\alpha(t_{ik} - t_{i,k-1})^{\alpha-1} e^{\beta'Z_i}) \right] \left[ \prod_{k=1}^{K_i} \prod_{t \in (t_{i,k-1}, t_{ik})} (1 - Y_i^R(t)\alpha(t - T_{N(t-)})^{\alpha-1} e^{\beta'Z_i} dt) \right] \\ &\quad \left[ \prod_{t \in (t_{iK_i}, C_i]} (1 - Y_i^R(t)\alpha(t - T_{N(t-)})^{\alpha-1} e^{\beta'Z_i} dt) \right] \\ &= \prod_{i=1}^n \left[ \prod_{k=1}^{K_i} (\alpha(g_{ik})^{\alpha-1} e^{\beta'Z_i}) \right] \left[ \prod_{k=1}^{K_i} \prod_{u \in (0, g_{ik})} (1 - Y_i^R(t_{i,k-1} + u)\alpha u^{\alpha-1} e^{\beta'Z_i} du) \right] \\ &\quad \left[ \prod_{u \in (0, C_i - t_{iK_i})} (1 - Y_i^R(t_{iK_i} + u)\alpha u^{\alpha-1} e^{\beta'Z_i} du) \right] \\ &= \prod_{i=1}^n \left[ \prod_{k=1}^{K_i} (\alpha(g_{ik})^{\alpha-1} e^{\beta'Z_i}) \right] \left[ \prod_{k=1}^{K_i} \exp\left\{-\int_0^{g_{ik}} Y_i^R(t_{i,k-1} + u)\alpha u^{\alpha-1} e^{\beta'Z_i} du\right\} \right. \\ &\quad \left. \exp\left\{-\int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u)\alpha u^{\alpha-1} e^{\beta'Z_i} du\right\} \right] \\ &= \prod_{i=1}^n \prod_{k=1}^{K_i} (\alpha(g_{ik})^{\alpha-1} e^{\beta'Z_i}) \exp\left\{-\int_0^{g_{ik}} Y_i^R(t_{i,k-1} + u)\alpha u^{\alpha-1} e^{\beta'Z_i} du\right\} \\ &\quad \exp\left\{-\int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u)\alpha u^{\alpha-1} e^{\beta'Z_i} du\right\}, \end{aligned}$$

where  $g_{ik} = t_{ik} - t_{i,k-1}$  with  $t_{i0} = 0$  are the gap times between the  $(k-1)$ th and  $k$ th EDMH visits for all subjects  $i = 1, \dots, n$ .

Let  $0 < g_1 < \dots < g_J$  be the distinct values of the gap times  $\{g_{ik} : k = 1, \dots, K_i; i = 1, \dots, n\}$ . Following Breslow (1972), we attain the MLE of  $\theta$  by maximizing log-transformation of  $L(\theta)$  above, viewing  $\lambda_0(g) = 0$  except for  $g = g_j, j = 1, \dots, J$ . The log likelihood function is

$$\begin{aligned} l(\theta|data) &= \sum_{i=1}^n \sum_{k=1}^{K_i} \log(\alpha) + (\alpha - 1) \log(g_{ik}) + \beta' Z_i - \int_0^{g_{ik}} Y_i^R(t_{i,k-1} + u) \alpha u^{\alpha-1} e^{\beta' Z_i} du \\ &\quad - \int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u) \alpha u^{\alpha-1} e^{\beta' Z_i} du, \end{aligned}$$

The likelihood score function of  $\theta$  is  $U(\theta) = (U_\alpha(\theta), U_\beta(\theta))'$ , with

$$\begin{aligned} U_\alpha(\theta) &= \sum_{i=1}^n \sum_{k=1}^{K_i} 1/\alpha + \log(g_{ik}) - \int_0^{g_{ik}} Y_i^R(t_{i,k-1} + u) (u^{\alpha-1} + \alpha u^{\alpha-1} \ln u) e^{\beta' Z_i} du \\ &\quad - \int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u) (u^{\alpha-1} + \alpha u^{\alpha-1} \ln u) e^{\beta' Z_i} du \end{aligned}$$

and

$$U_\beta(\theta) = \sum_{i=1}^n \sum_{k=1}^{K_i} Z_i - \int_0^{g_{ik}} Y_i^R(t_{i,k-1} + u) \alpha u^{\alpha-1} Z_i e^{\beta' Z_i} du - \int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u) \alpha u^{\alpha-1} Z_i e^{\beta' Z_i} du$$

The observed information matrix is  $I(\theta)$ , with four elements

$$\begin{aligned} I_{\alpha\alpha}(\theta) &= \sum_{i=1}^n \sum_{k=1}^{K_i} 1/\alpha^2 + \int_0^{g_{ik}} Y_i^R(t_{i,k-1} + u) (2 + \alpha \ln u) (u^{\alpha-1} \ln u) e^{\beta' Z_i} du \\ &\quad + \int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u) (2 + \alpha \ln u) (u^{\alpha-1} \ln u) e^{\beta' Z_i} du \end{aligned}$$

$$\begin{aligned} I_{\alpha\beta}(\theta) &= \sum_{i=1}^n \sum_{k=1}^{K_i} \int_0^{g_{ik}} Y_i^R(t_{i,k-1} + u) (u^{\alpha-1} + \alpha u^{\alpha-1} \ln u) Z_i e^{\beta' Z_i} du \\ &\quad + \int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u) (u^{\alpha-1} + \alpha u^{\alpha-1} \ln u) Z_i e^{\beta' Z_i} du \\ &= I_{\beta\alpha}(\theta)' \end{aligned}$$

$$I_{\beta\beta}(\theta) = \sum_{i=1}^n \sum_{k=1}^{K_i} \int_0^{g_{ik}} Y_i^R(t_{i,k-1} + u) \alpha u^{\alpha-1} Z_i Z_i' e^{\beta' Z_i} du + \int_0^{C_i - t_{iK_i}} Y_i^R(t_{iK_i} + u) \alpha u^{\alpha-1} Z_i Z_i' e^{\beta' Z_i} du$$