

Probation review report

Alessandro Gasparini

2017-08-02

Contents

Introduction	5
1 Introduction to survival analysis	7
1.1 Survival data	8
1.2 Censoring	8
1.3 Terminology and notation	10
1.4 Non-parametric survival analysis	10
1.5 Parametric survival analysis	10
1.6 The Cox model	10
1.7 Advanced survival analysis	10
2 Survival models with random effects	11
3 Joint models for longitudinal and survival data	13
4 Computational challenges in survival models with random effects	15
5 Simulation study: accuracy of Gaussian quadrature	17
5.1 Aim	17
5.2 Data-generating mechanisms	17
5.3 Methods	17
5.4 Estimands	17
5.5 Performance measures	17
5.6 Results	17
6 Simulation study: impact of misspecification in survival models with shared frailty terms	19
6.1 Aim	19
6.2 Data-generating mechanisms	19

6.3	Methods	19
6.4	Estimands	19
6.5	Performance measures	19
6.6	Results	19
7	Exploring results from simulation studies interactively	21
8	Informative visiting process	23
9	Future research developments	25
10	Personal development	27
10.1	Supervisory meetings	27
10.2	Training and courses	27
10.3	Conferences	29
A	Slides	31
B	Manuscript	33

Introduction

This report presents the work I have done during my first year as a PhD student at the Department of Health Sciences, University of Leicester, under the supervision of Dr. Michael Crowther and Prof. Keith Abrams.

I will begin by briefly introducing the topic of survival analysis in Chapter 1. Second, I will introduce survival models with random effects (e.g. frailties, in the simplest form) and joint models for longitudinal and time-to-event data in Chapters 2 and 3, respectively. Computational challenges that survival models with random effects and joint models pose are presented in Chapter 4. Third, I will present the results of two simulation studies in Chapters 5 and 6; the first simulation study investigates the accuracy of quadrature methods when approximating analytically intractable terms, while the second simulation study investigates the impact of model misspecification in survival models with shared frailty terms. Fourth, I will introduce an interactive tool I have been developing to aid the dissemination of results from simulation studies in Chapter 7. Then, I will introduce the problem of informative visiting process in clinical research using healthcare consumption data in Chapter 8, and how we aim to evaluate and compare the different approaches that have been proposed and utilised in literature to tackle such problem in Chapter 9. Finally, I will briefly summarise the training and personal development activities I have participated to during the first year of my PhD in Chapter 10.

This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/), and can be accessed online at <https://ellessenne.github.io/prr/>.

Chapter 1

Introduction to survival analysis

Survival analysis is a branch of statistics in which the main outcome consists in the time until the occurrence of a given event. Time could be years, months, weeks, or any amount of calendar time or even age time; event could be death, disease occurrence or relapse, or any other experience of interest. Survival analysis is also known as reliability theory in engineering, duration analysis in economics, and event history analysis in sociology. A broad overview of survival analysis is given in [Kleinbaum and Klein \(2012\)](#).

Some examples of time to event data are:

- disease remission in leukemia patients. In this study, leukemia patients are followed over several weeks to study how long they stay in remission status;
- heart disease occurrence. In this study, healthy subjects are followed over several years until occurrence of heart disease, or end of the study;
- renal failure. In this study, individuals with kidney disease are followed until renal failure, or end of the study;
- reliability of complex technical installations. For instance, studies assessing failure rates of components such as bulbs and valves.

In this Chapter I will define survival data and its peculiarities in Section [1.1](#) and [1.2](#). Terminology and notation used throughout this report will be introduced in Section [1.3](#). I will introduce common non-parametric and parametric methods in survival analysis in Section [1.4](#) and [1.5](#), respectively. I will introduce the widely used semi-parametric Cox model in Section [1.6](#). Finally, I will provide a brief overview of modern, advanced statistical methods in survival analysis in Section [1.7](#).

1.1 Survival data

Survival data generally consists, as previously mentioned, in time until the event of interest and an event indicator. In the leukemia remission example, time to event would be how many weeks it takes before a given patient experiences disease relapse and the event indicator would signal whether the individual relapsed or not before the end of the study. Nevertheless, in certain situations, we may have some information about the survival time but the actual survival time may be unknown. This problem is known as censoring and it is presented in Section 1.2.

1.2 Censoring

Censoring is a mechanism that causes survival times to be unobserved. There are many reasons why censoring may occur; among others:

1. a person does not experience the event before the end of the study;
2. a person drops out of the study before the occurrence of the event of interest;
3. a person experiences a competing event that impedes the occurrence of the event of interest (e.g.: death, when death is not the study outcome).

I simulated survival data for illustration purposes: I assumed a clinical trial with 10 individuals enrolled during a recruitment window of 5 years, and followed for up to 15 years. Not all individuals experience the event of interest during the study period, and are therefore censored at 15 years. The observation time for each individual is depicted in Figure 1.1 with a solid dark grey line, a cross represents the occurrence of the study event, and a circle represents administrative censoring. Individuals *C*, *E*, *G*, *H*, and *J* all have censored survival time: I know that they were still event-free at the end of follow-up, i.e. their real survival time is greater than the observed one, but the former is unknown.

This example represents a particular form of censoring: *right censoring*. The defining characteristic of right-censored data is that it is censored (or incomplete) at the right side of the follow-up time, hence the true survival time is greater than the observed time. This example represents *administrative censoring* as well, as individuals are censored at the end of the study to artificially restrict follow-up time (e.g. for financial reasons).

It is also possible to encounter data that is *left censored* or *interval censored*. In the former case, the true survival time is shorter than the observed one, e.g. I know that the event

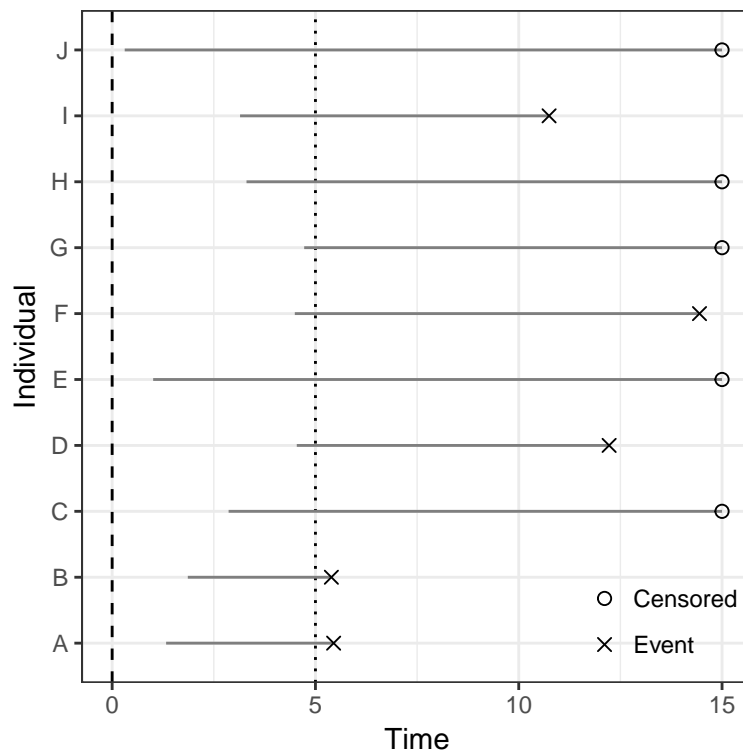


Figure 1.1: Simulated right censored survival data.

occurred before the observation time, but I do not know when - imagine onset of a viral infection, which can be detected only at a visit time. In the latter, I know that the event occurred within a certain interval of time but I do not know when; using the same example of infection onset, if infection was detected at a visit date but the individual was known to be infection-free at the previous visit, the true infection onset time is unknown and the event time is said to be interval censored.

Finally, another important concept related to right censoring is that of *left truncation* (or *delayed entry*). Left truncation occurs when an individual enrolls in the study some time after the inclusion criteria are satisfied; individuals that die (or emigrate, ...) before the start of observation time will never enter the study, and inclusion time may differ between individuals. Data arising from such phenomenon is therefore said to be left truncated.

1.3 Terminology and notation

1.4 Non-parametric survival analysis

1.5 Parametric survival analysis

1.6 The Cox model

1.7 Advanced survival analysis

Chapter 2

Survival models with random effects

Chapter 3

Joint models for longitudinal and survival data

Chapter 4

Computational challenges in survival models with random effects

Chapter 5

Simulation study: accuracy of Gaussian quadrature

5.1 Aim

5.2 Data-generating mechanisms

5.3 Methods

5.4 Estimands

5.5 Performance measures

5.6 Results

Chapter 6

Simulation study: impact of misspecification in survival models with shared frailty terms

6.1 Aim

6.2 Data-generating mechanisms

6.3 Methods

6.4 Estimands

6.5 Performance measures

6.6 Results

Chapter 7

Exploring results from simulation studies interactively

Chapter 8

Informative visiting process

Chapter 9

Future research developments

Chapter 10

Personal development

In this chapter I will introduce and briefly discuss the personal development activities I carried out during the first year of my PhD. In particular, I will present the supervisory meetings, training courses, and conferences I attended.

10.1 Supervisory meetings

I have been having frequent meetings with my supervisors, formally and informally. Formal supervisory meetings, recorded on PROSE (<https://prose.le.ac.uk>), have been held on average every other week, with summaries produced and shared between us. A comprehensive list is available on PROSE. Additionally, we held informal meetings to discuss developments and more urgent matters more often, whenever it was needed and every week on average.

10.2 Training and courses

I have attended a wide variety of courses during my first year, both externally and internally to the University of Leicester. The external courses I attended are:

- *Efficient R Programming*, on November 8th 2016, organised by the Royal Statistical Society in London. The instructor was Dr. Colin Gillespie, from the University of Newcastle, United Kingdom, and Jumping Rivers. The course covered how to program efficiently with R; in particular, it covered common pitfalls when writing R code, code profiling, RCpp, and parallel programming. General hints and tips were provided.

- *Introduction to causal inference*, on April 25th and 26th 2017, organised by the Biostatistics Research Group at the University of Leicester and delivered by Dr. Arvid Sjölander from Karolinska Institutet, Stockholm, Sweden. The course provided foundational concepts of causal inference such as the difference between association and causation, the counterfactual framework, exchangeability, directed acyclic graphs, methods for estimating a causal effect, etc. Additionally, it provided an introduction to more advanced methods such as instrumental variables and Mendelian randomisation.
- *Using simulation studies to evaluate statistical methods*, on May 22nd 2017, organised by University College London. The course was delivered by Dr. Tim Morris, Prof. Ian White and Dr. Michael Crowther, and it covered the rationale for using simulation studies, important concepts to keep in mind when planning a simulation study, computational tools, estimates of uncertainty, and tools for improving reporting and dissemination.
- Workshop on *Joint modelling of longitudinal and time-to-event data with R*, on July 5th, 2017, organised by the Department of Biostatistics of the University of Liverpool. The course was delivered by Dr. Graeme Hickey, and provided an introduction to joint models of longitudinal and survival data, including extensions to incorporate competing risks and multiple longitudinal processes and a practical session using R.

I have attended a few courses within the University and not offered on PROSE; specifically, I attended a course on *Time series analysis with R* (November 10th, 2016), a course on *Data visualisation* (November 15th, 2016), and a course on *High performance computing at Leicester* (February 8th, 2017). The latter was particularly important, as it allowed me to make better use of the high-performance computing facilities offered by the University. I also attended the *Preparing to teach in higher education* workshop, strand A (July 24th and 27th 2017).

Additionally, I have attended the following PROSE training sessions to develop personal and communication skills in research settings. These are listed below:

- *Planning your literature search*, October 21st 2016;
- *Conducting your literature search*, October 25th 2016 ;
- *Assertiveness*, November 14th 2016;
- *Introduction to critical thinking*, December 15th 2016;
- *Presentations A: Fundamentals of an effective presentation*, January 30th 2017;

- *Communication in research and other work settings*, January 31st 2017;
- *Enhancing your digital profile*, February 2nd 2017;
- *Saying it with your abstract*, February 10th 2017;
- *Designing a poster*, February 27th 2017;
- *Leadership in research and other work environments*, February 28th 2017;
- *Preparing for the probation review (Physical natural and medical sciences)*, May 30th 2017.

10.3 Conferences

I have attended a number of conferences during this year, in which I delivered the following oral presentations:

- Survival Analysis for Junior Researchers conference, held in Leicester, UK, on April 5th and 6th 2017. I delivered a talk titled *Direct likelihood maximisation using numerical quadrature to approximate intractable terms*;
- Statistical Analysis of Multi-Outcome Data (SAM) conference, held in Liverpool, UK, on July 3rd and 4th 2017. I delivered a talk titled *Impact of model misspecification in survival models with frailties*;
- Annual Conference of the International Society for Clinical Biostatistics conference, held in Vigo, Spain, on July 9th to July 13th 2017. I delivered two talks: a titled *Impact of model misspecification in survival models with frailties* during the main conference, and a talk titled *Exploring results from simulation studies interactively* during the Students' Day organised on July 13th.

Additionally, I delivered an oral presentation on previous work external to my PhD project during the 54th ERA-EDTA Congress held in Madrid, Spain, between June 3rd and June 6th. The ERA-EDTA Congress is the main conference in the field of Nephrology in Europe, with approximately 10,000 participants in 2017. I delivered my presentation, titled *Inappropriate prescription of nephrotoxic drugs to individuals with chronic kidney disease*, to an audience of clinicians, epidemiologists, clinical researchers, and other stakeholders.

Appendix A

Slides

Appendix B

Manuscript

Bibliography

Kleinbaum, D. G. and Klein, M. (2012). *Survival analysis: A self-learning text*. Springer-Verlag New York, 3 edition.