

**ESTIMATION OF THE BASE HAZARD FUNCTION
BY BOOTSTRAPPING**

RIFINA ARLIN

**Master of Science
Universiti Putra Malaysia**

December 2004

**ESTIMATION OF THE BASE HAZARD FUNCTION
BY BOOTSTRAPPING**

By

RIFINA ARLIN

**Thesis Submitted to the School of Graduate Studies, Universiti Putra Malaysia,
in Fulfilment of the Requirements for the Degree of Master of Science**

December 2004

Dedicated to:

My parents
(Drs. H. Zainal Arifin and Hj.
Darlina)

and

My beloved husband,
Danil Junaidy Daulay

Abstract of thesis presented to the Senate of Universiti Putra Malaysia in
fulfilment of the requirements for the degree of the Master of Science

**ESTIMATION OF THE BASE HAZARD FUNCTION
BY BOOTSTRAPPING**

By

RIFINA ARLIN

December 2004

Chairman: Associate Professor Isa Daud, Ph.D.

Faculty: Science

This thesis examines the techniques in estimating the base hazard function by bootstrapping. The base hazard function is a crucial part of survival analysis. It is used to construct an estimate of the proportional hazard model for every individual.

As in many methods for analysing survival data, this thesis utilizes the nonparametric model of Kaplan Meier, the Cox proportional hazard regression of the parametric model and the data validation by bootstrapping.

The Cox proportional hazard regression is used to model failure time data in censored data. Bootstrapping schemes validate the models based on Efron's technique and the data samples are generated using S-Plus programme randomizer.

Assessment of this method is investigated by performing simulation study on generated data. Two simulation studies are carried out to confirm the suitability of the models. Graph obtained from the results indicated that bootstrapping provides an alternative method in constructing estimation for base hazard function. This method is good alternative for a distribution-free approach with a minimal set of data.

Abstrak tesis yang dikemukakan kepada Senat Universiti Putra Malaysia
sebagai memenuhi keperluan untuk ijazah Master Sains

**PENGANGGARAN FUNGSI BAHAYA DASAR
DENGAN BOOTSTRAPPING**

Oleh

RIFINA ARLIN

Disember 2004

Pengerus: Profesor Madya Isa Daud, Ph.D.

Fakulti: Sains

Kajian di dalam tesis ini mengkaji teknik untuk penganggaran fungsi bahaya dasar dengan bootstrapping. Fungsi bahaya dasar merupakan bahagian yang penting didalam analisis mandirian. Ianya digunakan untuk membangun suatu anggaran model bahaya berkadaran untuk setiap individu.

Sebagaimana dalam banyak kaedah untuk menganalisis data mandirian, tesis ini menggunakan Kaplan-Meier model yang tak berparameter, model regresi bahaya berkadaran Cox biasa yang berparametrik dan pengesahan data oleh bootstrapping.

Regresi bahaya berkadaran Cox biasa digunakan untuk memodelkan data masa kegagalan di dalam data tertapis. Bootstrapping mengesahkan model berdasar pada teknik Efron's dan contoh data dihasilkan dengan menggunakan program S-Plus secara rawak.

Penilaian kaedah ini diselidiki dengan melakukan kajian simulasi ke atas data yang dijana. Dua kajian simulasi dijalankan bagi menentukan kesahihan model yang digunakan. Gambarajah yang diperolehi daripada keputusan menunjukkan bootstrapping menyediakan suatu kaedah alternatif dalam membangunkan penganggaran untuk fungsi bahaya dasar. Kaedah ini adalah alternatif yang baik untuk suatu pendekatan bebas taburan dengan suatu data yang minimal.

ACKNOWLEDGEMENTS

In the name of Allah, the Most Gracious, the Most Merciful. He has given me the strength to complete my study. Peace be upon the Prophet Muhammad S.A.W., his family, and his followers.

I would like to express my great appreciation and sincere gratitude to Associate Professor Dr. Hj. Isa Daud, chairman of my supervisory committee, for providing useful suggestions to improve this work. His guidance, encouragement and support had made my dream of pursuing higher education comes true. I am also indebted to my co-supervisor Associate Professor Dr. Hjh. Noor Akma Ibrahim and Associate Professor Dr. Mat Yusof Abdullah.

Thank you to the authorities of IRPA 54064 led by Associate Prof. Dr. Isa bin Daud for granting me the financial assistance through the Graduate Assistantship scheme.

Special thanks to Mr. Kudus and Mr. Fauzi, PhD students in UPM, who have helped me in understanding S-Plus and Bootstrap. My appreciation also goes to my housemates Ms. May, Ms. Tessa, Mrs. Mirna and Mrs. Siffa for their kindness, help and assistance. Thanks to all my fellow country men Dr. Mayastri, Dr. Susila, Ms.Desi, Mrs. Tryana, Mr. Iing and Mr. Sulaiman for the encouragement.

Last, but not least, my greatest appreciation and respect are to my loving parents and my parents-in-law for their endless support and motivation throughout my study. My

sister, my brother, my brother-in-law, my sister-in-law and my unforgettable husband, who have shown endless interest in my study. “Hasbunallah wa ni’mal wakiil ni’mal maulaa wa ni’man nashiir.”

I certify that an Examination Committee met on 21st December 2004 to conduct the final examination of Rifina Arlin on her Master of Science thesis entitled “Estimation of Base Hazard Function by Bootstrapping“ in accordance with Universiti Pertanian Malaysia (Higher Degree) Act 1980 and Universiti Pertanian Malaysia (Higher Degree) Regulation 1981. The Committee recommends that the candidate be awarded the relevant degree. Members of the Examination Committee are as follows:

Habshah Midi, PhD

Associate Professor
Faculty of Science
Universiti Putra Malaysia
(Chairman)

Mohd Rizam Abu Bakar, PhD

Associate Professor
Faculty of Science
Universiti Putra Malaysia
(Member)

Kassim Haron, PhD

Associate Professor
Faculty of Science
Universiti Putra Malaysia
(Member)

Zainodin Hj Jubok, PhD

Professor
Centre for Management of Research and Conference
Universiti Malaysia Sabah
(Independent Examiner)

GULAM RUSUL RAHMAT ALI, PhD

Professor/Deputy Dean
School of Graduate Studies
Universiti Putra Malaysia

Date:

This thesis submitted to the Senate of Universiti Putra Malaysia has been accepted as fulfilment of the requirement for the degree of Master of Science. The members of the Supervisory Committee are as follows:

Isa Daud, PhD

Associate Professor
Faculty of Science
Universiti Putra Malaysia
(Chairman)

Noor Akma Ibrahim, PhD

Associate Professor
Faculty of Science
Universiti Putra Malaysia
(Member)

Mat Yusof Abdullah, PhD

Associate Professor
Faculty of Science
Universiti Putra Malaysia
(Member)

AINI IDERIS, PhD

Professor/ Dean
School of Graduate Studies
Universiti Putra Malaysia

Date:

DECLARATION

I hereby declare that the thesis is based on my original work except for quotations and citations which have been duly acknowledged. I also declare that it has not been previously or concurrently submitted for any other degree at UPM or other institutions.

RIFINA ARLIN

Date: 28st April 2004

TABLE OF CONTENTS

	Page
DEDICATION	i
ABSTRACT	ii
ABSTRAK	iv
ACKNOWLEDGMENTS	vi
APPROVAL	viii
DECLARATION	x
LISTS OF TABLES	xiii
LISTS OF FIGURES	xv
CHAPTER	
1 INTRODUCTION	1.1
1.1 Background	1.1
1.2 Objective of the Research	1.3
2 LITERATURE REVIEW	2.1
2.1 Censored Data	2.1
2.1.1 Type I Censoring	2.2
2.1.2 Type II Censoring	2.2
2.1.3 Type III Censoring	2.3
2.2 Survivor Function and Hazard Function	2.3
2.3 Kaplan-Meier Estimate	2.5
2.4 Proportional Hazards Model	2.8
2.4.1 Estimation of the Proportional Hazards Regression Parameters	2.9
2.4.2 Estimating Baseline hazard function	2.13
2.5 Bootstrapping	2.14
2.5.1 Better Bootstrap Confidence Intervals	2.15
2.5.2 The Parametric Bootstrap	2.16
3 METHODOLOGY	3.1
3.1 Kaplan-Meier Estimate of the Hazard Function	3.1
3.2 The Cox Proportional Hazards Model	3.1
3.3 Distribution Used in the Generation of Data	3.2
3.3.1 Binomial Distribution	3.3
3.3.2 Exponential Distribution	3.3
3.4 Sample Size	3.4
3.5 Validation by Bootstrapping	3.5
3.5.1 The Algorithm for Estimating Standard Error	3.6

4	RESULTS AND DISCUSSION	4.1
4.1	Kaplan-Meier Type Estimate	4.1
	4.1.1 Data Intrauterine Device (IUD)	4.1
	4.1.2 Data Leukaemia	4.3
	4.1.3 Data High Cholesterol	4.4
4.2	The Cox Proportional Hazards Model	4.6
	4.2.1 Data the Stanford Heart Transplant	4.6
	4.2.2 Data Diseases of the Kidney Infection	4.10
	4.2.3 Data Myeloma	4.12
4.3	The Bootstrapping	4.14
	4.3.1 Confidence Intervals of the Hazard Function Based on the Kaplan-Meier Estimator by Bootstrapping	4.15
	4.3.2 Parameter in the Proportional Hazard Model Based on Maximum Likelihood Theory by Bootstrapping	4.22
5	CONCLUSION	5.1
5.1	Summary	5.1
5.2	Conclusion	5.2
5.3	Extension	5.3
	REFERENCES	R.1
	APPENDICES	A.1
	Appendix A. Programming in S-Plus	A.1
	Appendix B. Data	A.10
	BIODATA OF THE AUTHOR	B.1

LISTS OF TABLES

Tables	Page
4.1 Time in Week to Discontinuation of the Use of an IUD	4.1
4.2 Kaplan-Meier Type Estimate of the Hazard Function for the Data from Table 4.1	4.2
4.3 Times of Remission of Leukaemia Patients	4.3
4.4 Kaplan-Meier Type Estimate of the Hazard Function for the Data from Table 4.3	4.4
4.5 Times Data in Days of High Cholesterol	4.5
4.6 Kaplan-Meier Type Estimate of the Hazard Function for the Data from Table 4.5	4.5
4.7 Estimates of the Baseline Hazard Function for the Data the Stanford Heart Transplant	4.8
4.8 Estimates of the Baseline Hazard Function for the Data the Kidney Infection	4.11
4.9 Estimates of the Baseline Hazard Function for the Data Myeloma	4.13
4.10 Kaplan-Meier Type Estimate of the Hazard Function by Bootstrapping for the Data from Table 4.1	4.17
4.11 Kaplan-Meier Type Estimate of the Hazard Function by Bootstrapping for the Data from Table 4.3	4.19
4.12 Kaplan-Meier Type Estimate of the Hazard Function by Bootstrapping for the Data from Table 4.5	4.21
4.13 Biases and Replicates for Stanford Heart Transplant Data	4.24
4.14 Estimated Value of the Coefficients of the Explanatory Variables on Fitting a Proportional Hazard Model by Bootstrapping for Stanford Heart Transplant Data	4.26
4.15 Estimates of the Baseline Hazard Function by Bootstrapping for the Data the Stanford Heart Transplant	4.26

4.16	Biases and Replicates for Diseases of the Kidney Infection Data	4.28
4.17	Estimated Value of the Coefficients of the Explanatory Variables on Fitting a Proportional Hazard Model by Bootstrapping for Diseases of the Kidney Infection Data	4.29
4.18	Estimates of the Baseline Hazard Function by Bootstrapping for the Data Diseases of the Kidney Infection	4.30
4.19	Biases and Replicates for Myeloma Data	4.31
4.20	Estimated Value of the Coefficients of the Explanatory Variables on Fitting a Proportional Hazard Model by Bootstrapping for Myeloma Data	4.33
4.21	Estimates of the Baseline Hazard Function by Bootstrapping for the Data Myeloma	4.34

LISTS OF FIGURES

Figures	Page
2.1 Construction of Intervals Used in the Derivation of the Kaplan-Meier Estimate	2.7
3.1 Probability Density Function of Exponential Distribution	3.4
4.2 Graph Hazard Function of Kaplan-Meier for the Data from Table 4.1	4.2
4.3 Graph Hazard Function of Kaplan-Meier for the Data from Table 4.3	4.4
4.4 Graph Hazard Function of Kaplan-Meier for the Data from Table 4.5	4.6
4.5 Graph Base Hazard Function Cox Proportional Hazard for the Data the Stanford Heart Transplant	4.9
4.6 Graph Base Hazard Function Cox Proportional Hazard for the Data Diseases of the Kidney Infection	4.11
4.7 Graph Base Hazard Function Cox Proportional Hazard for the Data Myeloma	4.14
4.8 Upper and Lower $100(1 - \alpha)$ % Confidence Interval Efron's Percentile Bootstrap of the Hazard Function for the Data from Table 4.1	4.16
4.9 Graph Hazard Function of Kaplan-Meier with Bootstrap for the Data from Table 4.1	4.18
4.10 Upper and Lower $100(1 - \alpha)$ % Confidence Interval Efron's Percentile Bootstrap of the Hazard Function for the Data from Table 4.3	4.19
4.11 Graph Hazard Function of Kaplan-Meier with Bootstrap for the Data from Table 4.3	4.20
4.12 Graph Hazard Function of Kaplan-Meier with Bootstrap for the Data from Table 4.5	4.21
4.13 Plot of Bias with Replicate from 1000 Simulation for Stanford Heart Transplant Data	4.25
4.14 Graph Base Hazard Function Cox Proportional Hazard by Bootstrapping for the Data the Stanford Heart Transplant	4.27

4.15	Plot of Bias with Replicate from 1000 Simulation for Diseases of the Kidney Infection Data	4.29
4.16	Graph Base Hazard Function Cox Proportional Hazard by Bootstrapping for the Data Diseases of the Kidney Infection	4.31
4.17	Plot of Bias with Replicate from 1000 Simulation for Myeloma Data	4.33
4.18	Graph Base Hazard Function Cox Proportional Hazard by Bootstrapping for the Data Myeloma	4.35

CHAPTER 1

INTRODUCTION

Survival analysis deals with the modeling and analysis of data that has as a principal an end point to time until an event occurs. The failure time or survival time is the time elapsed between the entry of a subject into the study and the occurrence of an event thought to be related to the concerned (Anderson *et al.* 1980).

1.1 Background

A fundamental problem in survival analysis is estimating survival probabilities, as either a point estimate or a confidence interval. If a parametric model is assumed, such as the assumption of exponential or lognormal distributions for failure times, then estimation of the model parameters can be done using maximum-likelihood methods, modified to deal with the problem of censored data.

In this research, the data of interest are sets of survival times. These times represent a measured period from some origin t_0 to a particular observation. For example, time from recruiting an individual to a clinical trial to the time of the individual's death.

In some cases, the individual will be lost to the study before a failure time is observed. For example in a clinical trial, the individual may move from the area, or still be alive at the end of the trial. In such cases, the time is noted and referred to as a censoring time. Otherwise, the observation is a failure time or death time.

Deals with methods of estimating these functions from sample of survival data and are said to be non-parametric or distribution-free, since they do not require specific assumptions to be made about the underlying distribution of the survival times. Kaplan and Meier (1958) proposed a nonparametric method, which has been accepted as the standard estimator for such probabilities.

The proportional hazards model was proposed by Cox (1972) and has also come to be known as the Cox regression model. Although the model is based on the assumption of proportional hazards, no particular form probability distribution is assumed for the survival times.

Cox proportional hazards regression model or Cox model has been widely used as a well-known method for the analysis of clinical data (Farewell, 1979). The Cox model has been generalized to include designs, which include multiple failures for the same individual (Andersen and Gill, 1982; Wei, Lin, and Weissfeld, 1989) and cohort studies (Prentice, 1986).

Bootstrap methods resembling techniques for assessing uncertainty. They are useful when inference is to be based on a complex procedure for which theoretical results are unavailable or not useful for the sample sizes met in practice. Or when a standard model is suspect but it is unclear what to replace it with or where a ‘quick and dirty’ answer is required. They can also be used to verify the usefulness of standard approximations for parametric models, and to improve them if they seem to give inadequate inferences.

An improved bootstrap method that is second-order correct in a wider class of problem concerns setting approximate confidence intervals for a real-valued parameter θ in a multi parameter family (Efron, 1987). Beran, 1985 developed a higher-order approximate confidence intervals based on Edgeworth expansions which sometimes use bootstrap methods to reduce the theoretical computations.

1.2 Objective of the Research

The general objective of the research is to estimate the hazard function of survival data. These estimates can then be used to summarize the survival experience of individuals in the study. Therefore, the main objective of the study is to estimate the base hazard function in the proportional hazards model by bootstrapping through the estimation of the parameters, whereas the specific objectives are as follows:

- a) To estimate the hazard function based on the Kaplan-Meier estimator in small sample;
- b) To determine the confidence intervals of the hazard function based on the Kaplan-Meier estimator by bootstrapping;
- c) To obtain the parameters in the proportional hazards model based on maximum likelihood theory by bootstrapping.

CHAPTER 2

LITERATURE REVIEW

Survival analysis is the study of experiments which are performed to measure the amount of time that elapses until a particular event occurs. The time event can only occur once for a given subject and is usually described as the subject's failure time.

According to Miller (1981) survival analysis is a loosely defined statistical term that encompasses a variety of statistical techniques for analyzing positive-valued random variables. Typically the value of the random variable is the time to the failure of a physical component (mechanical or electrical unit) or the time to the death of a biological unit of patient, animal, cell etc.

Ideally, each subject would be observed until they fail, however, this is not always possible. For example, some subject may not have reached their failure time when the study is terminated or some patients in a medical study may die from causes unrelated to the disease being studied. The time at which a subject ceases to be observed for some reason other than failure is called the subject's censoring time. So, survival analysis is used for data in which censoring is present.

2.1 Censored Data

A distinctive characteristic of survival data is that the event of interest may not be observed for every experimental unit. This feature is known as censoring. For example when monitoring the life span of an electric bulb, the time to failure may not be

observable due to it is unavailability before it actually fails. The time to this event can be considered as censored time. Censoring may arise due to time limits and other restrictions depending on the nature of the experiment.

Survival data can arise from many of situations, particularly in problems of reliability and life time analysis (Basu, 1984).

2.1.1 Type I censoring

In clinical and epidemiological studies, censoring is mainly caused by a time restriction, and is therefore of type I (Lagakos, 1982).

Let t_c be some fixed number that is the fixed censoring time. Instead of observing T_1, T_2, \dots, T_n (the random variable of interest) we can only observe Y_1, Y_2, \dots, Y_n where

$$Y_i = \begin{cases} T_i & \text{if } T_i \leq t_c \\ t_c & \text{if } T_i > t_c \end{cases} \quad \begin{cases} \text{(uncensored)} \\ \text{(censored)} \end{cases} \quad (2.1)$$

The distribution function of Y has positive mass $P\{T > t_c\} > 0$ at $Y = t_c$.

2.1.2 Type II censoring

Let $r < n$ be fixed, and let $T_1 < T_2 < T_3 < \dots < T_n$ be the order statistic of failure times $T_1, T_2, T_3, \dots, T_n$. Observation ceases after the r -th failure, so we can observe $T_1, T_2, T_3, \dots, T_r$. In this design, the number of failures is not a random variable.

Both types of censoring designs, or generalizations of them, may be usefully adopted, provided that n, r and T are adequately chosen and appropriate methods are adopted for the analysis (Lawless, 1982).

2.1.3 Type III censoring

Let $C_1, C_2, C_3, \dots, C_n$ be independent and identically distributed (i.i.d) each with degrees of freedom (df) G . C_i is the censoring time associated with T_i . We observe $(Y_1, \delta_1), (Y_2, \delta_2), \dots, (Y_n, \delta_n)$, where

$$Y_i = \min(T_i, C_i)$$