



UNIVERSITI PUTRA MALAYSIA

**A MODEL FOR BINARY RESPONSE VARIABLE
WITH TIME-CENSORED OBSERVATIONS**

HANAN HASSAN ALI ADLAN

FSAS 2000 2

**A MODEL FOR BINARY RESPONSE VARIABLE
WITH TIME-CENSORED OBSERVATIONS**

By

HANAN HASSAN ALI ADLAN

**Thesis Submitted in Fulfilment of the Requirements for the
Degree of Master of Science in the Faculty of
Science and Environmental Studies
Universiti Putra Malaysia**

May 2000



To my Husband,

Husam & Aala

To

my father, mother, brothers, and sisters

♥ with love ♥

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

**A MODEL FOR BINARY RESPONSE VARIABLE
WITH TIME-CENSORED OBSERVATIONS**

By

HANAN H. A. ADLAN

May 2000

Chairman : Assoc. Prof. Dr. Isa Daud

Faculty : Science and Environmental Studies

The work in this thesis is concerned with the investigation and development of a new model for analyzing censored binary observations. Binary data is always expressed in the form 0 and 1. In other words, the response variable can only take on the value 0 or 1, which represent the non-response or response of the variable under consideration. In the case of grouped binary data, the response is expressed in the form of proportions. In certain follow-up studies, especially during the study period, some variables may result in censored observations. This becomes a problem especially in the case of ungrouped binary observations, where 0 is assigned to the non-responding individuals, or those who have got a chance to respond but not within the study period. This gives rise to the probit-exponential model.



As a basis, the probit model is used to analyze grouped binary data. It is also used in comparison with the suggested probit-exponential model in the case of ungrouped data. Comparison with the logistic fit is performed for the model with grouped data. The model diagnostic is done through the Pearson residuals, and standardized residuals.

Newton-Raphson algorithm is used to find the parameter estimates. A half normal plot is employed to investigate the usefulness of the two models, and to compare the suitability of the models when the observations are censored. Validity of the probit-exponential assumptions is also investigated through simulation.

The work above is carried out within the generalized linear model framework, which accommodate the new model as a member of this family.



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

**MODEL UNTUK PEMBOLEHUBAH RESPON DEDUA
DENGAN CERAPAN TAPISAN-MASA**

Oleh

HANAN H. A. ADLAN

May 2000

Pengerusi : Profesor Madya Dr. Isa Daud

Fakulti : Sains dan Pengajian Alam Sektiar

Skop dalam thesis ini tertumpu kepada siasatan dan perkembangan suatu model baru untuk menganalisis data tertapis bagi cerapan dedua. Data dedua lazim dinyatakan dalam bentuk 0 dan 1. Dengan kata lain, pembolehubah respon hanya boleh mengambil nilai 0 atau 1, yang mewakili tiada sambutan atau ada sambutan dari pembolehubah yang dipertimbangkan. Dalam kes data dedua terkumpul, respon dinyatakan dalam bentuk kadaran. Dalam beberapa kajian susulan tertentu, terutamanya dalam tempoh jangkamasa kajian, sebahagian daripada pembolehubah mungkin menghasilkan cerapan tertapis. Ini menimbulkan masalah terutamanya

dalam kes cerapan dedua tak terkumpul, dengan 0 diambil kira sebagai tiada sambutan.

Sebagai asas, model probit digunakan untuk menganalisis data dedua terkumpul. Ia juga digunakan sebagai perbandingan dengan model probit-eksponen tersaran bagi kes data yang tidak terkumpul. Perbandingan dengan cucukan logit dilakukan untuk model dari data terkumpul. Diagnostik model dijalankan melalui reja Pearson dan reja dipiawaikan.

Untuk mencari penganggar parameter kami tumpukan kepada penggunaan algorithma Newton-Raphson. Plot normal separuh digunakan untuk menyiasat kebaikan diantara kedua model dan untuk membandingkan kesesuaian model apabila cerapan adalah tertapis.

Kesahihan andaian probit-eksponen juga disiasat melalui simulasi. Kesemua kerja diatas dijalankan dalam kerangka model linear teritlak yang mana model baru ini tersirat sebagai ahli famili model.

ACKNOWLEDGEMENTS

Firstly and foremost, praise be to God, for giving me the strength and patience to complete this work.

I am deeply indebted to my supervisor Assoc. Prof. Dr. Isa Daud for his invaluable guidance, assistance, suggestions, helpful discussions and continuous encouragement during this work.

My thanks also to the members of my supervisory committee, Assoc. Prof. Dr. Mat Yusoff Abdullah, and Dr. Noor Akma Ibrahim.

Thanks are also due to Assoc. Prof. Dr Harun Budin, for allocating me a very comfortable place, when he was the head of the department.

Finally I would like to thank University of Khartoum for granting me study leave, and the government of Sudan for the award of scholarship.



TABLE OF CONTENTS

	Page
DEDICATION	ii
ABSTRACT	iii
ABSTRAK	v
ACKNOWLEDGMENTS	vii
APPROVAL SHEETS	viii
DECLARATION FORM	x
LIST OF TABLES	xiv
LIST OF FIGURES	xvi
LIST OF ABBREVIATIONS	xvii

CHAPTER

I	INTRODUCTION	1
	General overview	1
	Chapter overview	2
	Some Key Words and Definitions	2
	Statistical Modelling	2
	Generalized Linear Models	3
	Binary Data	4
	Survival Analysis	7
	Censoring	8
	Likelihood Estimation	10
	Newton-Raphson Method	12
	Simulation	12
	Half-Normal Plots of the Residuals	12
	The Problem	13
	Purpose and Scope of the Thesis	14
 II	 GENERALIZED LINEAR MODELS, THE PROBIT MODEL AND RELATED LITREATURE	 16
	Chapter overview	16
	Generalized Linear Models	16
	Introduction	16
	Definition of Generalized Linear Models	18



Estimation in Generalized Linear Models	20
Goodness of Fit in Generalized Linear Models	21
Model Checking	23
Residuals in Generalized Linear Models	23
The Probit Model with Application to Medical Data	24
Introduction	24
The Probit Model	26
Regression Using Probits	27
Estimation of Parameters in The Probit Model	28
Assumption of The Model	28
Estimation	28
Properties of Estimates	33
Testing and Goodness of Fit	33
Goodness of Fit.....	33
Test of Hypothesis.....	34
Confidence Interval	34
Application to Coronary Heart Disease Data	34
Probit Model for Data From a Cohort Study	35
Analysis	37
Conclusions	41
III THE PROBIT-EXPONENTIAL MODEL	42
Chapter overview	42
Introduction	42
The Proposed Model	45
The Model	48
Maximum Likelihood Estimation	48
Likelihood Estimator for the Proposed Model	50
The Model with Exponential Component	50
Estimation and Inference	50
Newton-Raphson Iteration	53
Applications	58
Findings	62
Conclusions	63
IV SIMULATION	64
Chapter overview	64
Introduction	64
Random Numbers	65
Simulating Continuous Random Variable	65
The Inverse transform method	66
The Rejection Method	66
Simulating Discrete Random Variable	68
Variance Reduction Technique.....	71



Application	72
Diagnostics	74
Conclusions	76
V CONCLUSIONS AND SUGGESTIONS FOR FURTHER RESEARCH	77
Overview	77
Conclusions	77
Suggestions for Further Research	79
BIBLIOGRAPHY	82
APPENDICES	86
A Algorithm for Fitting Generalized Linear Models	87
B Software Documentation and Programs	90
B1 Generalized Linear Interactive Modeling (GLIM4)	90
B2 Programs	92
C Additional Sheets	105
BIODATA OF THE AUTHOR	123



LIST OF TABLES

Table		Page
2.1	Proportions developed CHD in a twelve year period follow-up	36
2.2	Deviance obtained from the probit fit to the data of table 2.1	38
2.3	Parameter estimates, standard errors and the t values	39
3.1	Efficiency of the mixed probit-exponential model versus complete follow-up	57
3.2	Coded variables used in the analysis of disseminated breast cancer data	59
3.3	Asymptotic variance-covariance matrix for the estimates in the analysis of disseminated breast cancer data using the probit-exponential model	59
3.4	Asymptotic variance-covariance matrix for the estimates in the analysis of disseminated breast cancer data using the probit model	60
3.5	Selected values of the efficiency ratio $\text{var}(\hat{\beta}_p)/\text{var}(\hat{\beta}_c)$ for $T = \lambda^{-1}$, for disseminated breast cancer data	60
3.6	Asymptotic variance-covariance matrix for the estimates in the analysis of HIV+ patient data using the probit-exponential model	61
3.7	Asymptotic variance-covariance matrix for the estimates in the analysis of HIV+ patient data using the probit model	61
3.8	Selected values of the efficiency ratio $\text{var}(\hat{\beta}_p)/\text{var}(\hat{\beta}_c)$ for $T = \lambda^{-1}$, for HIV+ data	62
4.1	Asymptotic variance-covariance matrix for the estimates in the analysis of the simulated data using the probit-exponential model	73



4.2	Asymptotic variance-covariance matrix for the estimates in the analysis of the simulated data using the probit model	74
4.3	Selected values of the efficiency ratio $\text{var}(\hat{\beta}_p) / \text{var}(\hat{\beta}_e)$ for $T = \lambda^{-1}$, for the simulated data	74
5.1	Selected efficiency values for the mixed probit-exponential model	78



LIST OF FIGURES

Figure		Page
2.1	Maximum likelihood estimates using Fisher method of scoring	31
2.2	Proportions develop CHD by Serum Cholesterol levels	36
2.3	Data and fitted proportions by Serum Cholesterol levels	39
2.4	Pearson's residuals for the CHD data	40
2.5	Standardized residuals for the CHD data	41
3.1	Efficiency of the mixed probit-exponential model versus complete follow-up	58
4.1	The rejection method for simulating a random variable X having density function f	67
4.2	Standardized deviance residuals, with simulated envelopes (.....), and the means of the simulated values(----), for the probit-exponential model	75
4.3	Standardized deviance residuals, with simulated envelopes (.....), and the means of the simulated values(----), for the probit model.....	76



LIST OF ABBREVIATIONS

GLM	generalized linear models
GLIM4	generalized linear interactive modelling
pdf	probability density function
iid	independent and identically distributed
df	degrees of freedom
LF	likelihood function
ML	maximum likelihood
CHD	coronary heart disease
r.v.	random variable
BAN	best asymptotically normal



CHAPTER I

INTRODUCTION

General Overview

Modelling in Statistics in the recent years becomes the key for building systems in all aspects of sciences. The models are built and checked using modern methods for computing parameter estimates, and estimating errors. Fitting a model to data set enables the understanding of the essential features of the data. Of course, the model fitted will not completely describe reality but it can be thought of as an approximation to the true situation. The so called asymptotic theory is employed in looking for optimal approximate solutions. This makes the study of the analysis become the most important area of research for many years. Also the emergence of the analysis of incomplete data sets refreshes this area, especially the analysis of censored data which has become a popular field in recent researches. Old methods for building models are improved and re-evaluated so that it can deal with censoring.

In this thesis, our concern is to develop and build a probit-exponential regression model for censored data, although the work on the probit model has an older application and a lot of work has been done on it, but its application to censored data taking into account the time factor is of new interest. The problem is to use the transformed probits in building a survival model to censored observation, in



order to implement a method for analysing binary variables within the generalised linear models framework.

Chapter Overview

In this chapter we highlight some background for the scope of the research. Some important key words are defined. General description of the formulation of the problem, and the purpose of the research are outlined.

Key Words and Definitions

Statistical Modelling

The term statistical modelling referred to modelling observations of characteristics of individuals or groups, which are assumed to have random elements.

The set of observations will be a sample representing some population. The variable characteristics are assumed to respond to mechanisms involving the various circumstances of the different individuals. This is referred to as the response variables. The aim of modelling is to explain how the variation in the observed values of the response variable is explained by differences in circumstances between individual cases or groups of cases. These differences may result from nature (gender, age), nurture (social environment, education), exposure (diet, pollution, media) or treatment (drug, fertilisers, educational techniques), and are termed explanatory variables.

GLM as a single class. GLM has been introduced by Nelder and Wedderburn (1972) as a unifying family of models for non-normal cross-sectional regression analysis with non-normal responses. Their further development had a major influence on statistical modelling in a wider sense (Fahrmeir and Tutz ,1994; Francis , Green, and Payne,1994).

Binary Data

In many areas of application of statistical principles and procedures, one encounters observations made on individual experimental units that take one of two possible forms. For example, a patient in a clinical trial to compare alternative forms of treatment may or may not experience relief from symptoms; an insect in an insecticidal trial may survive or die when exposed to a particular dose of an insecticide; a seed may germinate or fail to germinate under certain experimental conditions; an integrated circuit may be defective or non-defective. Such data are said to be binary. The two possible forms are always referred to as success and failure.

The term grouped binary data refer to the number of successes out of the total number of units exposed to a particular set of experimental conditions. Binary data are assumed to have Bernoulli distribution, while data in the form of proportions can be modelled using the binomial distribution. Special form of regression model is required for analyzing such data. For example, an experiment in which it is noted that there is an increasing proportions with increasing dose leads to an S-shaped

curve relating the proportions to the dose. It would be possible to transform the response variable and fit the resulting variable using one of the statistical techniques.

Modelling of proportions involves the modern method of Probit analysis (Bliss, 1935) remains popular in this field of applications, also the logistic regression (Cox 1972) is of wider application (D.Collett, 1991; Francis, Green, and Payne, 1994).

Consider the case of ungrouped binary responses, coded by zero and one. Given the covariate x , a binary variable Y is completely determined by its response probability

$$E(Y / x) = P(Y = 1 / x) = \pi \quad (1.1)$$

implying
$$\text{var}(Y / x) = \pi(1 - \pi), \quad (1.2)$$

for the case of grouped data. Let \bar{y} denotes the relative frequency of observed 1's for the say m independent binary observations – with the same covariate vector x .

The absolute frequencies $m\bar{y}$ are binomially distributed with

$$E(m\bar{Y} / x) = m\pi \quad (1.3)$$

$$\text{var}(m\bar{Y} / x) = m\pi(1 - \pi). \quad (1.4)$$

The relative frequencies are scaled binomial, i.e. they take the corresponding values $0, 1/m, 2/m, \dots, 1$ with the same binomial probabilities as $m\bar{y}$ and

$$E(\bar{Y} / x) = \pi. \quad (1.5)$$

$$\text{var}(Y / x) = \frac{\pi(1 - \pi)}{m}. \quad (1.6)$$

Thus for grouped binary data, response functions are the same as for individual binary responses, with the variance function has to be divided by m i.e. $w = m$. When

the individual response (given \mathbf{x}) is binomially distributed with $Y \sim B(m, \pi)$, the scaled binomial or relative frequency y/m will be considered as response, since y/m can act as an average of individual independent binary observations, and can be treated as occurring from grouped observations. Then the variance

$$\text{var}(Y / \mathbf{x}) = \frac{\pi(1 - \pi)}{m} \quad (1.7)$$

is the same as for grouped binary observations.

Models for binary and binomial responses are determined by relating the response probability π to a linear predictor $\eta = \beta X$ via some response function $\pi = h(\eta)$, link function $g(\pi) = \eta$. The most common models in this field are linear probability model, the probit, and the logit model.

Our concern in this thesis based on implementation of the probit model, which can be defined by $\pi = \Phi(\eta) = \Phi(\beta X)$ where Φ is the standard normal distribution function, no restrictions on η .

The Logit model corresponds to the natural link function

$$g(\pi) = \log\left(\frac{\pi}{1 - \pi}\right) = \eta \quad (1.8)$$

with the logistic distribution function $\pi = h(\eta) = \frac{\exp(\eta)}{1 + \exp(\eta)}$, (1.9)

as the resulting response function. The logistic distribution function also has support on the entire real axis and is symmetric, but it has some heavier tails than the standard normal. Apart from π values near zero or one which correspond to the tails, fits by probit or logit models are generally quite similar (Fahrmeir and Tutz, 1994).

Survival Analysis:

Survival analysis is concerned with studying the time between entry to a study and a subsequent event. Originally the analysis was concerned with time from treatment until death, hence the name, but survival analysis is applicable to many areas as well as mortality. Recent examples include time to discontinuation of a contraceptive, maximum dose of bronchoconstrictor required to reduce a patient's lung function to 80% of baseline, time taken to exercise to maximum tolerance, time that a transdermal patch can be left in place, time for a leg fracture to heal.

When the outcome of a study is the time between one event and another, a number of problems can occur.

1. The times are most unlikely to be Normally distributed.
2. We cannot afford to wait until events have happened to all the subjects, for example until all are dead. Some patients might have left the study early - they are lost to follow up. Thus the only information we have about some patients is that they were still alive at the last follow up. These are termed censored observations.

Let the survival time, T for individuals in a population have density function $f(t)$. The corresponding distribution function $F(t) = \int_{-\infty}^t f(s)ds$ is the function of the population dying by time t . The complementary function $1 - F(t)$ (survivor function) is the fraction still surviving at time t . The hazard function $\lambda(t)$ measures the instantaneous risk, and it means that $\lambda(t)\delta_t$ is the probability of dying in the next small interval δ_t given survived to time t . These functions can be expressed as

$$\Pr(\text{survival to } t + \delta_t) = \Pr(\text{survival to } t) \Pr(\text{survival for } \delta_t / \text{survival to } t)$$

$$1 - F(t + \delta_t) = (1 - F(t))(1 - \lambda(t)\delta_t) \quad (1.10)$$

hence $\delta_t F'(t) = (1-F(t))\lambda(t)\delta_t$ and we have $\lambda(t) = f(t)/(1-F(t))$ (1.11)

(Cox, and Oakes, 1984; Miller, 1981; McCullagh, and Nelder, 1989).

Censoring

Censoring represent a mechanism which is considered a special feature of survival data. It occurs when some subjects in the study have not experienced the event of interest. At the end of the study, the exact survival times of these subjects are unknown. And they are referred to as censored observations or censored times. In other words in an experiment in which subjects are followed over time until the event of interest (such as death or other type of failure) occurs, it is not always possible to follow every subject until the event is observed. Subjects may drop out of the study and be lost to follow-up, or be deliberately withdrawn, or the end of the data collection period may arrive before the event is observed to happen. For such a subject, all that is known is that the time to the event is at least as long as the time to when the subject was last observed. The observed time to the event under such circumstances is censored. As mentioned before censoring may occur from the right that is the observation stops before the event is observed, as in censorship for survival analysis, or from the left when the observations does not begin until after the event has occurred.

The types of right censoring

1- type I censoring

Subjects are entered the study for a fixed period in the calendar. At the end of the study, subjects that don't experience the event of interest or lost to follow-up are censored.