

TIME DEPENDENT LOGISTIC MODELS IN FOLLOW-UP STUDIES AND
CLINICAL TRIALS, I. BINARY DATA

Institute of Statistics Mimeo Series No. 1309

TIME DEPENDENT LOGISTIC MODELS IN FOLLOW-UP STUDIES AND
CLINICAL TRIALS, II. MULTI-RESPONSE DATA

Institute of Statistics Mimeo Series No. 1310

by

Regina C. Elandt-Johnson

Department of Biostatistics

University of North Carolina at Chapel Hill

October 1980

TIME DEPENDENT LOGISTIC MODELS
IN FOLLOW-UP CLINICAL TRIALS
I. BINARY DATA.
II. MULTI-RESPONSE DATA.

Regina C. Elandt-Johnson
Department of Biostatistics
University of North Carolina
Chapel Hill, N.C. 27514, U.S.A.

INTRODUCTION

In recent years there have been many epidemiological follow-up studies and controlled clinical trials, investigating the roles of various characteristics as possible either predictors or risk factors for the incidence of, or death from different chronic diseases. In particular, an especially controversial topic is a possible association between lipid levels and coronary heart diseases. This has drawn the attention of many medical research centers.

Two mathematical models for the analysis of this type of data are in common use: logistic linear model and Cox's proportional hazard rate model. We are considered here with the logistic linear model.

Originally, this was designed as a response model for use when the responses are independent of time, and so it could be used in retrospective type experiments. Of course, even in prospective studies, one might neglect the effect of time, asking the simple question: "What is the probability of death in a *specified* period?" If a *cohort* of N individuals is observed and D individuals die during this period, then this probability is simply estimated by the proportion $p = D/N$. If, however, there are withdrawals during this period, p would underestimate this probability. This fact is well recognized in the construction of life tables by adjusting for withdrawals. The same principle applies to more advanced models, in which death probability is a function of concomitant variables. In such cases, it is often practicable to use approximations (usually exponential) to the survival function over fairly short periods of time.

The purpose of these two articles is to present theoretical bases for constructing various kinds of likelihood functions where using exponential-logistic models as approximations to survival functions, and when there is full or partial information on time of death or withdrawal. Part I is concerned with binary data (deaths from all causes), and Part II with competing risk situations. The likelihoods are constructed for data grouped in fixed intervals, as well as for small sets of data from clinical trials, when only information on individuals in the risk sets at the times of death is taken into account. It is also emphasized that over a short interval of time, exponential-logistic models can be successfully approximated by Cox's proportional hazard rate models.

Some of the likelihood functions obtained (especially in Part I) have already appeared (though sometimes in different forms) but are scattered in the literature. But even these seem to have escaped the attention of most epidemiological researchers. Nearly all of the recently published articles and reports seem to use standard programs constructed for situations, where time is *not* taken into account.

I hope that these two papers, aimed both at reviewing some work already done and also deriving and discussing some new formulae, will be helpful for the analysis of follow-up and clinical trial data. The focus in both papers is on estimation - the problems of testing of goodness of fit are not considered here in detail. For short intervals, however, exponential-logistic or Cox's models appear to be useful tools in analysis of morality data.

TIME DEPENDENT LOGISTIC MODELS IN FOLLOW-UP
STUDIES AND CLINICAL TRIALS, I. BINARY DATA

Regina C. Elandt-Johnson^{*)}
Department of Biostatistics
University of North Carolina
Chapel Hill, N.C. 27514, U.S.A.

ABSTRACT

The 'classical' logistic model given by (1.4) is useful in analyzing data from experiments in which response does not depend on time. Its parameters can be estimated from cross-sectional- or retrospective-type data. The model is, however, inappropriate to use in analyzing data from prospective studies in which endpoint events (incidence of a disease, death) are time dependent. In such situations, a survival distribution function needs to be incorporated in the model. A motivation for this paper was to present some theoretical bases and concepts on applicability of logistic models for consideration by traditional epidemiologists and some statisticians, who may sometimes feel that the time-independent logistic function is a 'magic' device for answering the question: which factors are more beneficial and which are more harmful for prolongation of life? Since the follow-up or clinical trial data are usually censored, it would be difficult to fit a completely defined parametric distribution; instead, we consider in this paper fitting piecewise exponential distribution with hazard rates determined by a logistic linear relation. In Section 1, the exponential-logistic model defined over a specified period τ [also considered by Myers et al. (1973)] is briefly described. Two kinds of

likelihood functions when fitting exponential-logistic model to grouped data are discussed; in Section 2.1 we consider data when information on time of death or withdrawal is available to the nearest interval [see also Thompson (1977)], and in Section 2.2, exact times in each interval are taken into account. Estimation techniques based on partial (and conditional) likelihoods, applicable to small sets of data from clinical trials, are presented in Section 3. It is emphasized that exponential-logistic model is closely related to that of Cox's (1972) proportional hazard rate model.

Key Words & Phrases: Response model; Exponential distribution; Exponential -logistic model; Piecewise exponential fitting; Likelihood function (conditional, partial).

*) This work was supported by U.S. National Heart, Lung and Blood Institute contract NIH-NHLI-712243 from the National Institutes of Health.

1. BINARY RESPONSE LOGISTIC MODELS

1.1 Introduction. Time-Independent Response Models

Logistic linear models were originally designed as tools in discriminant analysis, for classifying individuals or items into one of two (or more) distinct classes on the basis of certain characteristics selected as discriminant factors and measured on each object [Fisher (1936)]. Logistic models have also been used in special types of response experiments, when the response does not involve time [Cox (1970)]. Examples in medical applications of logistic models of this kind are given by Anderson (1973), (1974).

Consider first a response experiment in which, at a given time point, the response variable, Y , say, takes only two values

$$Y = \begin{cases} 1 & \text{if response is observed} \\ 0 & \text{otherwise.} \end{cases} \quad (1.1)$$

Let $\underline{z}' = (z_1, z_2, \dots, z_s)$ represent a set of s characteristics measured on each individual, which are selected as possible predicting factors of the response.

Further, let

$$\Pr\{Y = 1 | \underline{z}\} = q(\underline{z}), \quad \Pr\{Y = 0 | \underline{z}\} = p(\underline{z}), \quad (1.2)$$

with

$$q(\underline{z}) + p(\underline{z}) = 1, \quad (1.3)$$

be the *posterior* probabilities that an individual with characteristics \underline{z} is a responder or a non-responder, respectively.

The logistic linear model is defined by the formula

$$\log \frac{q(\underline{z})}{p(\underline{z})} = \gamma + \sum_{u=1}^s \beta_u z_u = \gamma + \underline{\beta}' \underline{z}, \quad (1.4)$$

where β_u is a parameter determining the contribution of variable z_u

to the response, and γ represents an 'average' response of all characteristics, including also an 'underlying' (unknown) contribution to the response.

Clearly, the z 's do not need to be direct measurements; they can represent some functions (e.g., logarithms, powers, etc.) of original observations, or possibly products of these; in the last cases the appropriate β 's represent interaction effects among covariables.

From (1.3) and (1.4), we obtain

$$p(z) = [1 + \exp(\gamma + \beta'z)]^{-1}, \quad (1.5)$$

and

$$q(z) = \exp(\gamma + \beta'z) [1 + \exp(\gamma + \beta'z)]^{-1}. \quad (1.6)$$

Likelihood function

Suppose that N individuals were in the trial and N_1 were observed to respond. Let

$$\delta_j = \begin{cases} 1 & \text{if individual (j) responds} \\ 0 & \text{otherwise,} \end{cases} \quad (1.7)$$

for $j = 1, 2, \dots, N$.

The likelihood function is

$$\begin{aligned} L(\gamma, \beta) &= \prod_{j=1}^N [q(z_j)]^{\delta_j} [p(z_j)]^{1-\delta_j} = \prod_{j=1}^N \left[\frac{q(z_j)}{p(z_j)} \right]^{\delta_j} p(z_j) \\ &= \prod_{j=1}^N \exp[\delta_j(\gamma + \beta'z_j)] [1 - \exp(\gamma + \beta'z_j)]^{-1} \end{aligned} \quad (1.8)$$

[cf., Cox (1970), formula (2.17)]. Hence

$$\log L(\gamma, \beta) = \sum_{j=1}^N \delta_j(\gamma + \beta'z_j) - \sum_{j=1}^N \log[1 + \exp(\gamma + \beta'z_j)]. \quad (1.9)$$

This leads to the set of $(s+1)$ equations

$$\frac{\partial \log L}{\partial \gamma} = N_1 - \sum_{j=1}^N q(z_j) = 0 \quad (1.10)$$

and

$$\frac{\partial \log L}{\partial \beta_v} = \sum_{j=1}^N [\delta_j - q(z_j)] z_{vj} = 0, \quad v = 1, 2, \dots, s. \quad (1.11)$$

Note that the expected number of responders, E , say is

$$E = \sum_{j=1}^N \xi(Y_j = 1 | z_j) = \sum_{j=1}^N q(z_j). \quad (1.12)$$

Thus (1.10) is equivalent to the condition that the estimated expected number of responders is equal to the observed number of responders. In order to find how well the data fit the model, it might be useful to calculate the number of false positives and false negatives [Harrell (1978)].

From the view point of sampling theory, the estimation techniques (e.g., maximum likelihood) just described are applicable to data collected from cross-sectional or retrospective type experiments, in which the response data represent *prevalence* of the event under consideration.

However, epidemiologists also use the logistic function (1.4) in analyzing data from *prospective* studies, in which the response is *incidence* of some endpoint events such as, for example, onset of a chronic disease or death. Logistic models designed for retrospective studies, may be inappropriate for prospective studies.

1.2 Time Dependent Response

We now consider experiments, in which response is not obtained at once, but there might be a lapse of time before the response is observed. The time 'due to respond', T , say, is a random variable. A typical example of this kind is the incidence of a chronic age dependent disease. It should, however, be realized that because time is involved, mortality also has to be taken into account. For example, if high values of z represent higher risk for a disease, it might happen that those with

such high values have already died even before the disease was diagnosed. On the other hand, if high values have a beneficial effect on individuals suffering from the disease, then those with high (and medium) values might be the survivors. In both situations, we may find the effect of covariable z "not significant" while, in fact, such an effect may exist. Therefore, in studying onset of a disease, death should be considered as just another response variable, so that multi-response models would be more appropriate (see Part II).

In this section we assume, for convenience, that the endpoint is death.

Let for $t > 0$

$$S_T(t|z) = \exp\left[-\int_0^t \lambda(y|z) dy\right] = \exp[-\Lambda(t|z)] \quad (1.13)$$

be the survival distribution function (SDF), $\lambda(t|z)$ - the hazard rate and $\Lambda(t|z)$ - the cumulative hazard rate, given set of covariables z is observed. We first assume that z 's do not depend on t .

The probability of surviving till time t is

$$\Pr\{T > t|z\} = \exp[-\Lambda(t|z)], \quad (1.14)$$

and the probability of death before time t is

$$\Pr\{T \leq t|z\} = 1 - \exp[-\Lambda(t|z)]. \quad (1.15)$$

By specifying the mathematical form of \log (odds ratio) as a function of t , we can determine the survival function (1.13). In the general case, we have

$$\log \frac{1 - \exp[-\Lambda(t|z)]}{\exp[-\Lambda(t|z)]} = \psi(t; z), \quad (1.16)$$

so that

$$S_T(t|z) = \exp[-\Lambda(t|z)] = [1 + \exp(\psi(t; z))]^{-1}. \quad (1.17)$$

For example, if

$$\psi(t; \underline{z}) = \gamma t + \underline{\beta}' \underline{z} , \quad (1.18)$$

then

$$S_T(t|\underline{z}) = [1 + \exp(\gamma t + \underline{\beta}' \underline{z})]^{-1} \quad (1.19)$$

is a multiple logistic survival distribution function.

Hankey and Mantel (1974) considered a model

$$\psi(t; \underline{z}) = \sum_{r=0}^k \gamma_r t^r + \sum_{u=1}^s \beta_u z_u , \quad (1.20)$$

In these examples, z 's are independent of t . The hazard rate function, however, depends on t and its explicit form is specified by the function $\psi(t; \underline{z})$. Further simplification can be effected by making the hazard rate dependent only on z 's, over a given period.

Exponential-logistic model

Of special interest is a model discussed by Meyer et al. (1973), in which

$$S_T(t|\underline{z}) = \exp[-\lambda(\underline{z})t] \quad (1.21)$$

is *exponential*, and the log (odds ratio) is a *linear* function of z 's over a period $(0, \tau)$, that is

$$\log \frac{q(\tau|\underline{z})}{p(\tau|\underline{z})} = \log \frac{1 - \exp[-\lambda(\underline{z})\tau]}{\exp[-\lambda(\underline{z})\tau]} = \gamma + \underline{\beta}' \underline{z} . \quad (1.22)$$

Hence

$$p(\tau|\underline{z}) = \exp[-\lambda(\underline{z})\tau] = [1 + \exp(\gamma + \underline{\beta}' \underline{z})]^{-1}, \quad (1.23)$$

so that for $0 < t \leq \tau$, we have

$$p(t|\underline{z}) = \exp[-\lambda(\underline{z})t] = [1 + \exp(\gamma + \underline{\beta}' \underline{z})]^{-t/\tau} = [p(\underline{z})]^{-t/\tau}, \quad (1.24)$$

where $p(\underline{z})$ is defined in (1.5).

From (1.24), we obtain

$$\lambda(\underline{z}) = \frac{1}{\tau} \log [1 + \exp(\gamma + \underline{\beta}' \underline{z})] = -\frac{1}{\tau} \log p(\underline{z}). \quad (1.25)$$

The period τ can be fixed in advance, or - as these authors have done - τ (in their notation W) can be considered as a parameter which has to be estimated from the data.

Two special cases might be of some interest [see also Meyer et al. (1973)].

(i) Suppose that $\tau \rightarrow 0$ (i.e., τ is small), then (1.22) becomes

$$\log \frac{1 - \exp[-\lambda(\underline{z})\tau]}{\exp[-\lambda(\underline{z})\tau]} = \log [1 - \exp(-\lambda(\underline{z})\tau)]$$

$$\doteq \log [\lambda(\underline{z})\tau] = \gamma + \beta' \underline{z} , \quad (1.26)$$

or

$$\lambda(\underline{z}) = \lambda_0 \exp(\beta' \underline{z}) , \quad (1.27)$$

where $\lambda_0 = \tau^{-1} \exp(\gamma)$.

This is a special case of Cox's (1972) model

$$\lambda(t|\underline{z}) = \lambda(t) \exp(\beta' \underline{z}) , \quad (1.28)$$

with the underlying hazard $\lambda(t) = \lambda_0$ over period $(0, \tau)$.

(ii) Suppose that $\tau \rightarrow \infty$ (i.e., τ is very large), then (1.22)

becomes

$$\log \frac{1 - \exp[-\lambda(\underline{z})\tau]}{\exp[-\lambda(\underline{z})\tau]} \doteq \log [\exp(\lambda(\underline{z})\tau)] = \lambda(\underline{z})\tau = \gamma + \beta' \underline{z} \quad (1.29)$$

or

$$\lambda(\underline{z}) = \gamma_0 + \sum_{u=1}^s \gamma_u z_u , \quad (1.30)$$

which is an *additive* model for the hazard rate with $\gamma_0 = \gamma\tau^{-1}$, and $\gamma_u = \beta_u \tau^{-1}$.

The basic question arises, whether and when we are justified in fitting a special parametric model to data from a given experiment. In prospective studies, the period of investigation is often restricted, so that even if we fit a certain distribution, our data relate, in fact, to only its left-hand tail. Often, we try to approximate the unknown

distribution by piecewise exponential, with hazard rate for the i th period of the form

$$\lambda_i(t|z) = \lambda_i \exp(\beta'z) . \quad (1.31)$$

As was mentioned above, (1.31) can provide a good approximation to 'logistic' hazard (1.25) if the i th period is short, and there is no significant difference in results whether (1.25) or (1.31) is used.

Although medical researchers often employ (1.31) in the analysis of data from clinical trials, there is a fairly well established tradition among epidemiologists of using logistic models. In view of this tradition, and because of some interesting and useful properties of logistic models, we will present a few kinds of likelihood functions in which logistic model is utilized. But first, we give some further details on likelihood function when exponential-logistic model defined in (1.22) [with hazard rate defined by (1.25)] over a period $(0, \tau)$ is fitted.

Likelihood function

Similarly as in Section 1.2, let N denote the number of participants in the study, and N_1 the number of deaths. Let t_j be the time at which individual (j) with characteristics $z_j' = (z_{1j}, \dots, z_{sj})$ was last seen in $(0, \tau)$.

Further, let

$$\delta_j = \begin{cases} 1 & \text{if } (j) \text{ died at } t_j \\ 0 & \text{if } (j) \text{ was alive at } t_j. \end{cases} \quad (1.32)$$

Then the likelihood function is

$$L(\tau; \gamma, \beta) = \prod_{j=1}^N [\lambda(z_j)]^{\delta_j} \exp[-\lambda(z_j)t_j] , \quad (1.33)$$

or

$$\begin{aligned} \log L(\tau, \gamma, \beta) &= \sum_{j=1}^N \delta_j \log \lambda(z_j) - \sum_{j=1}^N \lambda(z_j) t_j \\ &= A(\gamma, \beta) - \frac{1}{\tau} B(\gamma, \beta) - N_1 \log \tau, \end{aligned} \quad (1.34)$$

where

$$A(\gamma, \beta) = \sum_{j=1}^N \delta_j \log \{ \log [1 + \exp(\gamma + \beta' z_j)] \}; \quad (1.35)$$

and

$$B(\gamma, \beta) = \sum_{j=1}^N t_j \log [1 + \exp(\gamma + \beta' z_j)]. \quad (1.36)$$

For τ fixed in advance, the likelihood equations are

$$- \frac{\partial \log L}{\partial \gamma} = \sum_{j=1}^N \delta_j \frac{q(z_j)}{\log p(z_j)} + \frac{1}{\tau} \sum_{j=1}^N t_j q(z_j) = 0, \quad (1.37)$$

$$- \frac{\partial \log L}{\partial \beta_v} = \sum_{j=1}^N \delta_j \frac{q(z_j) z_{vj}}{\log p(z_j)} + \frac{1}{\tau} \sum_{j=1}^N t_j q(z_j) z_{vj} = 0, \quad v = 1, 2, \dots, s, \quad (1.38)$$

where $p(z_j)$ and $q(z_j)$ are defined in (1.5) and (1.6), respectively.

When τ is considered as an unknown parameter, it is easy to see that the ML-estimator of τ is

$$\hat{\tau} = N_1^{-1} B(\gamma, \beta), \quad (1.39)$$

and is entirely a function of γ and β 's.

Substituting (1.39) into (1.33), we obtain the maximized likelihood as a function of γ and β 's alone, $L'(\gamma, \beta)$, say; the estimates of γ and β 's should be chosen to maximize $L'(\gamma, \beta)$.

The algebra is straightforward.

Likelihood functions, when the t_j 's are recorded to the nearest interval are discussed by Myers et al. (1973).

2. ESTIMATION FROM GROUPED DATA USING PIECEWISE
EXPONENTIAL-LOGISTIC APPROXIMATION.

We first discuss the practicality of the assumption that z 's do not depend on t . In fact, there are but a few kinds of covariables used in survival models which do not change with time. For example, demographic characteristics such as sex or race do not depend on time. Most of the continuous characteristics, however, are subject to changes in time, and the pattern of the changes is often unknown. In some cases, it might be just random variation depending on many uncontrolled factors; in others, there may be some functional dependence of z on time, but the form of this relation is difficult to establish. To 'stabilize' such covariables, only two (or a few) discrete values are used. For example, we may consider only "normal" and "high" blood pressure, giving the corresponding z values 0 or 1, respectively. Of course, using such dichotomized values, some information is lost, and one may wonder, whether this is a right thing to do.

But we also notice that the results in the preceding section have been, in fact, derived *conditionally* on the observed set of values z . Strictly speaking, these should represent measurements at time point t , but, in practice, they are often taken at individual's entry; their values for individual (j) are denoted by $z_j' = (z_{1j}, z_{2j}, \dots, z_{sj})$.

If the period of investigation is long enough, it may be divided into M *fixed* intervals, $[t_i, t_{i+1})$ of length $\tau_i = t_{i+1} - t_i$, $i = 0, 1, \dots, M - 1$. To allow for variation of z 's in time, it is desirable to obtain a new set of measurements in each interval: for the j th individual in the interval $[t_i, t_{i+1})$, the set $z_{ij}' = (z_{1ij}, z_{2ij}, \dots, z_{sij})$. For practical purposes, these measurements can be taken at the time of last contact with the j th individual observed during the period t_i to t_{i+1} .

We now construct likelihood function when piecewise exponential-logistic model is fitted to grouped data.

Assume exponential-logistic relation (1.22) over the period τ_i in the form

$$\log \frac{q_i(\tau_i | z)}{p_i(\tau_i | z)} = \log \frac{1 - \exp[-\lambda_i(z)\tau_i]}{\exp[-\lambda_i(z)\tau_i]} = \gamma_i + \beta'z, \quad (2.1)$$

for $i = 0, 1, \dots, M - 1$.

We have then for $t_i \leq t < t_{i+1}$

$$p_i(t|z) = [1 + \exp(\gamma_i + \beta'z)]^{-\frac{(t-t_i)}{\tau_i}} - [p_i(z)]^{-\frac{(t-t_i)}{\tau_i}} \quad (2.2)$$

where

$$p_i(z) = [1 + \exp(\gamma_i + \beta'z)]^{-1}, \quad (2.3)$$

$$q_i(z) = 1 - p_i(z) = \exp(\gamma_i + \beta'z) [1 + \exp(\gamma_i + \beta'z)]^{-1}, \quad (2.4)$$

and

$$\lambda_i(z) = \frac{1}{\tau_i} \log[1 + \exp(\gamma_i + \beta'z)] = -\frac{1}{\tau_i} \log p_i(z), \quad (2.5)$$

for $t_i \leq t < t_{i+1}$.

Let t_{ij} be the time at which the j th individual was last seen in the interval $[t_i, t_{i+1})$; $\tau_{ij} = t_{ij} - t_i$ ($0 < \tau_{ij} < \tau_i$), and z_{ij} be the set of characteristics measured on the j th individual in the interval $[t_i, t_{i+1})$. Note that in model (2.1), the β 's are assumed to be the same over all intervals, while γ_i 's change from one interval to another. Further, we also assume that the vector of observations, z_{ij} , is available for each interval. If, however, z 's are measured only once, at the time of individual's entry, z_{ij} can easily be replaced by z_j .

We now outline two kinds of likelihood functions, depending on which information on the time of last seen is available.

2.1 Actuarial-Type Approach

A similar approach to that presented below was discussed by Thompson (1977), though the likelihood equations in his paper are given in a different form.

Suppose that the information on the exact times t_{ij} is not available or not utilized; we only know the interval in which death or withdrawal occurred. Assuming that, on the average, withdrawal occurred in the middle of the interval, the probability that an individual with covariables \underline{z} , who withdraws, survives till the time of withdrawal in $[t_i, t_{i+1})$ is approximately $[p_i(\underline{z})]^{1/2}$.

Let \mathcal{D}_i , \mathcal{W}_i and \mathcal{S}_i denote the sets of individuals who died, withdrew in, or survived over the interval $[t_i, t_{i+1})$, and d_i , w_i and s_i - the corresponding numbers in each category. The likelihood function over $[t_i, t_{i+1})$ is

$$L_i^{(1)}(\gamma_i, \beta) = \prod_{j \in \mathcal{D}_i} q_i(z_{ij}) \prod_{j \in \mathcal{S}_i} p_i(z_{ij}) \prod_{j \in \mathcal{W}_i} [p_i(z_{ij})]^{1/2} \quad (2.6)$$

and the overall likelihood function is

$$L^{(1)}(\underline{\gamma}, \beta) = \prod_{i=0}^{M-1} L_i(\gamma_i, \beta), \quad (2.7)$$

where $\underline{\gamma}' = (\gamma_0, \gamma_1, \dots, \gamma_{M-1})$.

The likelihood equations are

$$\frac{\partial \log L^{(1)}}{\partial \gamma_i} = d_i - \left[\sum_{j \in \mathcal{D}_i \cup \mathcal{S}_i} q_i(z_{ij}) + \frac{1}{2} \sum_{j \in \mathcal{W}_i} q_i(z_{ij}) \right] = 0 \quad (2.8)$$

and

$$\frac{\partial \log L^{(1)}}{\partial \beta_v} = \sum_{j \in \mathcal{D}_i} z_{vij} - \left[\sum_{j \in \mathcal{D}_i \cup \mathcal{S}_i} q_i(z_{ij}) z_{vij} + \frac{1}{2} \sum_{j \in \mathcal{W}_i} q_i(z_{ij}) z_{vij} \right] = 0 \quad (2.9)$$

for $v = 1, 2, \dots, s$.

The sum in the square brackets in the left hand side of (2.8) gives the expected number of deaths, E_i , say, so that we have from (2.8) that the observed number of deaths is equal to their expected value [cf., (1.10) and (1.12)].

When there are no covariates ($\beta = \tilde{0}$), we have

$$q_i(0) = q_i = \frac{d_i}{s_i + d_i + \frac{1}{2} w_i} . \quad (2.10)$$

This is the actuarial estimator of conditional probability of death in $[t_i, t_{i+1})$.

2.2 Likelihood When Individual Records Are Utilized

Suppose that data are grouped in fixed intervals, but we also have records of individual times at death or withdrawal, t_{ij} .

Let

$$\delta_{ij} = \begin{cases} 1 & \text{if the } j\text{th individual dies at } t_{ij}, \\ 0 & \text{otherwise .} \end{cases} \quad (2.11)$$

Replacing in (1.33), δ_j , z_j and t_j by δ_{ij} , z_{ij} and τ_{ij} , respectively, the likelihood function for the interval $[t_i, t_{i+1})$ is

$$L_i^{(2)}(\gamma_i; \beta) = \prod_{j=1}^N [\lambda_i(z_{ij})]^{\delta_{ij}} \exp[-\lambda_i(z_{ij}) \tau_{ij}] , \quad (2.12)$$

and the overall likelihood is

$$L^{(2)}(\gamma, \beta) = \prod_{i=0}^{M-1} L_i(\gamma_i; \beta) . \quad (2.13)$$

Clearly,

$$\log L^{(2)}(\gamma, \beta) = \sum_{i=0}^{M-1} A_i(\gamma_i, \beta) - \sum_{i=0}^{M-1} \frac{1}{\tau_i} B_i(\gamma_i, \beta) - \sum_{i=0}^{M-1} N_{1i} \tau_i , \quad (2.14)$$

where A_i , B_i are the analogues of A and B defined in (1.35) and (1.36), respectively, and N_{1i} is the number of deaths in the interval $[t_i, t_{i+1})$.

The likelihood equations are

$$-\frac{\partial \log L_i^{(2)}}{\partial \gamma_i} = \sum_{j=1}^N \delta_{ij} \frac{q_i(z_{ij})}{\log p_i(z_{ij})} + \frac{1}{\tau_i} \sum_{j=1}^N \tau_{ij} q_i(z_{ij}) = 0, \quad i=0, 1, \dots, M-1, \quad (2.15)$$

and

$$-\frac{\partial \log L_i^{(2)}}{\partial \beta_v} = \sum_{i=0}^{M-1} \sum_{j=1}^N \delta_{ij} \frac{q_i(z_{ij})}{\log p_i(z_{ij})} z_{vij} + \sum_{i=0}^{M-1} \tau_{ij} q_i(z_{ij}) z_{vij} = 0, \quad (2.16)$$

[c.f., (1.37) and (1.38), respectively].

3. PARTIAL AND CONDITIONAL LIKELIHOODS

Suppose that the data are results from a clinical trial in which not many deaths were observed. Grouping in fixed intervals might not be feasible.

Let

$$t'_1 < t'_2 < \dots < t'_n \quad (3.1)$$

denote the observed ordered (and distinct) times at death. We may use the observed times at death as division points and apply the method discussed in Section 2.2 over the intervals $(t'_{i-1}, t'_i]$, that is, conditional on t'_1, t'_2, \dots, t'_n . Note that since deaths are recorded "at" the time point, we use the intervals $(t'_{i-1}, t'_i]$ rather than $[t'_i, t'_{i+1})$ to conform with the usual habit of handling this type of data. The likelihood for the interval $(t'_{i-1}, t'_i]$ is essentially the same as that in (2.12) [for $[t_i, t_{i+1})$], but it could be slightly simplified because death(s) will occur only "at" the end of the interval.

3.1 Partial likelihood function

Sometimes only exact times at death are recorded, but there is no information on withdrawal times, or the investigator concentrates attention on the time points at which the event (death) was observed. This means that the data are treated as *if* withdrawals occur at the *beginning* of the period $(t'_{i-1}, t'_i]$. We first consider situations in which z 's are either independent of time, or are recorded only once, at the

time of individual's entry.

We also introduce a slightly different notation. Let $j = 1, 2, \dots, N$ denote now an individual who was, for some time, participating in the study; $j_f(i)$ - individual (j_f) who is the f th to die at time t'_i , and $R(t'_i)$ - the risk set at $t'_i - 0$, that is, the set of individuals alive just before t'_i , and $t'_i - t'_{i-1} = \tau'_i$. If there were $m_i (>1)$ deaths recorded "at t'_i " (i.e., we observe ties), then the likelihood (2.12) takes the form

$$L_i^{(2')}(\gamma_i, \beta) = \prod_{f=1}^{m_i} \lambda_i(z_{j_f(i)}) \times \prod_{\ell \in R(t'_i)} \exp[-\lambda_i(z_\ell) \tau'_i] . \quad (3.2)$$

Substituting for $\lambda_i(z)$ from (2.5), we obtain

$$L_i^{(2')}(\gamma_i, \beta) = (\tau'_i)^{-m_i} \prod_{f=1}^{m_i} \log[1 + \exp(\gamma_i + \beta' z_{j_f(i)})] \times \prod_{\ell \in R(t'_i)} [1 + \exp(\gamma_i + \beta' z_\ell)]^{-1} . \quad (3.3)$$

The general expression for the likelihood function in (3.3) is not altered if z 's are replaced by $z(t'_i)$'s.

The overall likelihood is

$$L^{(2')}(\gamma, \beta) = \prod_{i=1}^n L_i^{(2')}(\gamma_i, \beta) . \quad (3.4)$$

3.2 Conditional Partial Likelihood Function

When there are no multiple deaths (i.e., no tied observations) "at" each t'_i , it would be convenient to use a kind of conditional likelihood outlined below.

Let $\lambda_j(t' | z_j)$ denote the hazard rate for individual (j) at time t' , and let $R(t')$ be the risk set at $t' - 0$. The conditional probability, given that a death takes place at time point t' and individual (j) is the one who died, is

$$\pi_j(t') = \frac{\lambda_j(t' | z_j)}{\sum_{\ell \in R(t')} \lambda_\ell(t' | z_\ell)} . \quad (3.5)$$

Using the hazard rate of the exponential-logistic model defined in (2.1), the conditional likelihood that individual (j) is the one who died at t_i' [denoted by $j(i)$] is

$$L_i^{(3)}(\gamma_i, \beta) = \frac{\log[1 + \exp(\gamma_i + \beta'z_{j(i)})]}{\sum_{\ell \in R(t_i')} \log[1 + \exp(\gamma_\ell + \beta'z_\ell)]}, \quad (3.6)$$

and the overall likelihood function is

$$L^{(3)}(\gamma, \beta) = \prod_{i=1}^n L_i^{(3)}(\gamma_i, \beta). \quad (3.7)$$

It is easy to see that the form of the likelihood (3.6) does not change, if we replace $z_{j(i)}$ by $z_{j(i)}(t_i')$, and z_ℓ by $z_\ell(t_i')$, respectively, provided such observations are available.

We also note [from (2.1)] that

$$\exp(\gamma_i + \beta'z) = \frac{q_i(\tau_i|z)}{p_i(\tau_i|z)}. \quad (3.8)$$

Since usually $q_i(\tau_i|z) \ll p_i(\tau_i|z)$, then $\exp(\gamma_i + \beta'z) < 1$. Therefore, expanding the numerator and the denominator in (3.6) by Taylor series, we obtain

$$L_i(\gamma_i, \beta) \doteq \frac{\exp(\beta'z_{j(i)})}{\sum_{\ell \in R(t_i')} \exp(\beta'z_\ell)}, \quad (3.9)$$

which is the same as given by Cox (1972), p. 191, formula (12).

4. SUMMARY AND DISCUSSION

(i) It has been stressed throughout the text that the traditional time independent logistic model (1.4) may be appropriate for retrospective study, but is incorrect to use it in prospective studies; in the latter, the survival distribution function (SDF) has to be considered.

(ii) As an approximation to the SDF over a short interval, the commonly used exponential distribution was employed, and its hazard rate was determined from the logistic relation (2.1).

(iii) Various kinds of likelihood functions were constructed in the framework of exponential-logistic approximation, and depending on the information on time of death or withdrawal available (or utilized) from the data.

(iv) It was also emphasized through the paper a close relation of exponential-logistic model to Cox's model when it is assumed that the underlying hazard rate in Cox's model is approximately constant over a period between two consecutive deaths. Some advantage of likelihood (3.6) over (3.9) might be that all parameters can be estimated using (3.7) and so the SDF can be constructed while from Cox's (1972) model, only β 's can be estimated.

In view of this relation and to save space, we do not illustrate our results by an example. Any data one wishes to use, should show not great difference. It is just a matter of preference (and, perhaps, availability of programs or costs) and traditional training to decide, which model - exponential-logistic or Cox's proportional hazard rate - to apply.

REFERENCES

- Anderson, J.A. (1973). Logistic discrimination with medical application. In: *Discriminant Analysis and Applications*, ed. by T. Cacoullos, Academic Press, pp. 1-15.
- Anderson, J.A. (1974). Diagnosis by logistic discriminant function: further practical problems and results. *Applied Statistics* 23, 397-404.
- Cox, D.R. (1970). *The Analysis of Binary Data*, Methuen, London.
- Cox, D.R. (1972). Regression model and life tables (with Discussion). *Journal of the Royal Statistical Society, Series B* 34, 187-220.
- Fisher, R.A. (1936). The use of multiple measurements in taxonomic problems. *Annals of Eugenics* 7, Part II, 179-188.
- Elandt-Johnson, R.C. (1981). Time dependent logistic models in follow-up studies and clinical trials. II. Multi-response data. (submitted)
- Hankey, B.F. and Mantel, N. (1974). A logistic regression analysis of response time data where the hazard function is time dependent. *Communication in Statistics - Theory and Methods* A7 (4), 333-347.
- Harrell, F. (1979). The LOGIST procedure. *Duke Medical Center Publication No.* , 83-102.
- Myers, M.H., Hankey, B.F. and Mantel, N. (1973). A logistic exponential model for use with response-time data involving regressor variables. *Biometrics* 29, 257-269.
- Thompson, W.A. (1977). On the treatment of grouped observations in life studies. *Biometrics* 33, 463-470.

TIME DEPENDENT LOGISTIC MODELS IN FOLLOW-UP
STUDIES AND CLINICAL TRIALS, II. MULTI-RESPONSE DATA

Regina C. Elandt-Johnson*)
Department of Biostatistics
University of North Carolina
Chapel Hill, N.C. 27514, U.S.A.

ABSTRACT

In Part I, we have described fitting exponential-logistic models to follow-up or clinical trial type data, when the time dependent response was death from any cause. In this paper, the results are generalized to multi-response logistic models (multiple causes of death) in the framework of competing risk theory with concomitant information on a set of characteristics z . In Section 1, we briefly discuss time independent multi-response logistic models and maximum likelihood estimation of the parameters. Basic principles and concepts of competing risk theory, a multi-response exponential-logistic model over a fixed period τ , and an appropriate likelihood function are presented in Section 2. In Section 3, fitting piecewise exponential-logistic distributions to grouped mortality data from K causes, is discussed. Methods for ungrouped data are presented in Section 4. It is first shown that when the period τ is small, the hazard rate can be approximated by a Cox's-type multiplicative exponential function, depending on parameters associated with the k th cause only. The partial and conditional likelihoods are constructed using the 'waiting time' distribution function (introduced in Section 2.1), for each cause separately.

It turns out that for estimating the parameters of a specific cause, any computer program fitting Cox's model can easily be adapted.

Key Words & Phrases: Logistic models; Multiple responses; Competing risks; Hazard rates; 'Waiting time' distributions; Likelihood functions (conditional, partial).

*)This work was supported by U.S. National Heart, Lung and Blood Institute contract NIH-NHLI-71-2243 from the National Institutes of Health.

1. MULTI-RESPONSE TIME INDEPENDENT LOGISTIC MODELS

Consider an experiment in which K types of response may take place at a specified time point. Individuals may be classified into distinct groups (classes) on the basis of their responses.

Let $\underline{z}' = (z_1, \dots, z_s)$ be a set of prediction (or discrimination) variables, and Y be the response variable such that

$$Y = \begin{cases} k & \text{if the } k\text{th response is observed} \\ 0 & \text{if there is no response.} \end{cases} \quad (1.1)$$

Further, let

$$\Pr\{Y = k | \underline{z}\} = q_k^*(\underline{z}), \quad k = 1, 2, \dots, K \quad (1.2)$$

be the *posterior* probability of the k th response, and

$$\Pr\{Y = 0 | \underline{z}\} = p(\underline{z}) \quad (1.3)$$

be the probability of no response. We must have

$$q_1^*(\underline{z}) + q_2^*(\underline{z}) + \dots + q_K^*(\underline{z}) = q(\underline{z}) = 1 - p(\underline{z}). \quad (1.4)$$

A multi-response logistic (linear) model can be defined as

$$\log \frac{q_k^*(\underline{z})}{p(\underline{z})} = \gamma_k + \sum_{u=1}^s \beta_{ku} z_u = \gamma_k + \beta_k' \underline{z}, \quad (1.5)$$

where $\beta_k' = (\beta_{k1}, \beta_{k2}, \dots, \beta_{ks})$ [see Cox (1970), Anderson (1972)].

From (1.4) and (1.5), we obtain

$$p(\underline{z}) = \left[1 + \sum_{r=1}^K \exp(\gamma_r + \beta_r' \underline{z}) \right]^{-1}, \quad (1.6)$$

and

$$\begin{aligned}
 q_k^*(z) &= [\exp(\gamma_k + \beta_k'z)]p(z) \\
 &= \exp(\gamma_k + \beta_k'z) \left[1 + \sum_{r=1}^K \exp(\gamma_r + \beta_r'z) \right]^{-1}, \quad (1.7)
 \end{aligned}$$

or

$$q_k^*(z) = c_k(z)q(z), \quad (1.8)$$

where

$$c_k(z) = \exp(\gamma_k + \beta_k'z) \left[\sum_{r=1}^K \exp(\gamma_r + \beta_r'z) \right]^{-1}, \quad (1.9)$$

for $k = 1, 2, \dots, K$, and with

$$c_1(z) + c_2(z) + \dots + c_K(z) = 1. \quad (1.10)$$

Likelihood function

Let N be the number of individuals in the sample; N_k - the number of responders in the k th category, $k = 1, 2, \dots, K$; N_{K+1} - the number of nonresponders, and $\gamma' = (\gamma_1, \dots, \gamma_K)$, $\beta' = (\beta_1, \dots, \beta_K)$. Let $\tilde{z}_j' = (z_{1j}, z_{2j}, \dots, z_{sj})$ be the vector of the z 's observed on individual (j), and

$$\delta_{kj} = \begin{cases} 1 & \text{if the } k\text{th response is observed for individual } (j); \\ 0 & \text{otherwise.} \end{cases} \quad (1.11)$$

The likelihood function is

$$\begin{aligned}
 L^{(1)}(\gamma, \beta) &= \prod_{j=1}^N \prod_{k=1}^K [q_k^*(z_j)]^{\delta_{kj}} [p(z_j)]^{1 - \sum_{k=1}^K \delta_{kj}} \\
 &= \left\{ \prod_{j=1}^N \prod_{k=1}^K [q_k^*(z_j)/p(z_j)]^{\delta_{kj}} \right\} \times \prod_{j=1}^N p(z_j) \\
 &= \left\{ \prod_{j=1}^N \prod_{k=1}^K \exp[\delta_{kj}(\gamma_k + \beta_k'z_j)] \right\} \times \prod_{j=1}^N \left[1 + \sum_{r=1}^K \exp(\gamma_r + \beta_r'z_j) \right]^{-1}, \quad (1.12)
 \end{aligned}$$

so that

$$\log L^{(1)}(\underline{\gamma}, \underline{\beta}) = \sum_{j=1}^N \sum_{k=1}^K \delta_{kj} (\gamma_k + \beta'_k z_j) - \sum_{j=1}^N \log \left[1 + \sum_{r=1}^K \exp(\gamma_r + \beta'_r z_j) \right]. \quad (1.13)$$

It is easy to show that the set of likelihood equations is

$$\frac{\partial \log L^{(1)}}{\partial \gamma_h} = \sum_{j=1}^N [\delta_{hj} - q_h^*(z_j)] = 0 \quad (1.14)$$

and

$$\frac{\partial \log L^{(1)}}{\partial \beta_{hv}} = \sum_{j=1}^N [\delta_{hj} - q_h^*(z_j)] z_{vj} = 0 \quad (1.15)$$

for $h = 1, 2, \dots, K$; $v = 1, 2, \dots, s$.

We also notice that the expected number of responders in the k th category is

$$E_k = \sum_{j=1}^N q_k^*(z_j), \quad k = 1, 2, \dots, K, \quad (1.16)$$

and the total expected number of responders is

$$E = \sum_{k=1}^K E_k = \sum_{j=1}^N \sum_{k=1}^K q_k^*(z_j). \quad (1.17)$$

Since $\sum_{j=1}^N \delta_{kj} = N$, we obtain from (1.14)

$$N_k = E_k \quad \text{for } k = 1, 2, \dots, K, \quad (1.18)$$

that is, the observed number of responders is equal to its expected number, in each category.

2. MULTIPLE TIME DEPENDENT RESPONSES

2.1 Basic Functions in Competing Risk Models

Consider K mutually exclusive time dependent responses and, for convenience, suppose that these are deaths from K different causes,

C_1, C_2, \dots, C_K .

We first briefly introduce some basic concepts in competing risk theory.

Let $\lambda_k(t|\underline{z})$ be the hazard at time $t (> 0)$ for cause C_k , acting *in the presence* of all other causes, and $\lambda(t|\underline{z})$ be the overall (i.e. from any cause) hazard rate for an individual with characteristics \underline{z} . We assume that \underline{z} 's do not depend on t .

We have

$$\lambda(t|\underline{z}) = \lambda_1(t|\underline{z}) + \lambda_2(t|\underline{z}) + \dots + \lambda_K(t|\underline{z}). \quad (2.1)$$

Let

$$S_T(t|\underline{z}) = \Pr\{T > t|\underline{z}\} = P(t|\underline{z}) = \exp\left[-\int_0^t \lambda(u|\underline{z}) du\right] \quad (2.2)$$

be the overall (from any cause) survival function, and

$$Q(t|\underline{z}) = 1 - P(t|\underline{z}) \quad (2.3)$$

be the failure distribution.

Further, let

$$Q_k^*(t|\underline{z}) = \int_0^t \lambda_k(u|\underline{z}) S_T(u|\underline{z}) du \quad (2.4)$$

denote the failure probability function for cause C_k *in the presence* of all other causes. We have

$$Q_1^*(t|\underline{z}) + Q_2^*(t|\underline{z}) + \dots + Q_K^*(t|\underline{z}) = Q(t|\underline{z}). \quad (2.5)$$

We also have [from (2.1) and (2.2)]

$$S_T(t|\underline{z}) = P(t|\underline{z}) = \prod_{k=1}^K \exp\left[-\int_0^t \lambda_k(u|\underline{z}) du\right] = \prod_{k=1}^K G_k(t|\underline{z}), \quad (2.6)$$

where

$$G_k(t|z) = \exp\left[-\int_0^t \lambda_k(u|z) du\right]. \quad (2.7)$$

The survival function in (2.7) can be considered as a kind of 'waiting time' probability function for cause C_k alone, for an individual with covariables z , and with hazard rate $\lambda_k(t|z)$.

When the hazard rates are proportional, that is

$$\lambda_k(t|z) = \pi_k(z)\lambda(t|z), \quad (2.8)$$

then

$$Q_k^*(t|z) = \pi_k(z)Q(t|z), \quad (2.9)$$

and

$$G_k(t|z) = [S_T(t|z)]^{\pi_k(z)}, \quad (2.10)$$

with

$$\pi_1(z) + \pi_2(z) + \dots + \pi_k(z) = 1. \quad (2.11)$$

2.2 Multi-Response Exponential-Logistic Distributions Over a Fixed Period

We now consider a specified period $(0, \tau)$ and denote the probabilities corresponding to $Q_k^*(t|z)$, $Q(t|z)$ and $P(t|z)$ by $q_k^*(t|z)$, $q(t|z)$ and $p(t|z)$ for $0 < t \leq \tau$, respectively.

We wish to approximate $q_k^*(t|z)$ for $0 < t \leq \tau$, by an exponential probability function of the form

$$q_k^*(t|z) = \pi_k(z)q(t|z) = \pi_k(z)\{1 - \exp[-\lambda(z)t]\}, \quad (2.12)$$

so that

$$p(t|z) = \exp[-\lambda(z)t] \quad \text{for } 0 < t \leq \tau \quad (2.13)$$

is the left-hand tail of exponential distribution.

We also have

$$q_1^*(t|z) + q_2^*(t|z) + \dots + q_K^*(t|z) = q(t|z) = 1 - p(t|z). \quad (2.14)$$

Note that the forms of $\pi_k(z)$'s and $\lambda(z)$ are not yet specified.

Suppose that it is reasonable to assume an exponential-logistic model over the period $(0, \tau)$, that is

$$\begin{aligned} \log \frac{q_k^*(\tau|z)}{p(\tau|z)} &= \log \frac{\pi_k(z) \{1 - \exp[-\lambda(z)\tau]\}}{\exp[-\lambda(z)\tau]} \\ &= \gamma_k + \sum_{u=1}^S \beta_{ku} z_u = \gamma_k + \beta'_k z \end{aligned} \quad (2.15)$$

for $k = 1, 2, \dots, K$.

We obtain [from (2.11), (2.14) and (2.15)]

$$\pi_k(z) = \exp(\gamma_k + \beta'_k z) \left[\sum_{r=1}^K \exp(\gamma_r + \beta'_r z) \right]^{-1} = c_k(z), \quad (2.16)$$

where $c_k(z)$ is defined in (1.9).

$$p(\tau|z) = \exp[-\lambda(z)\tau] = \left[1 + \sum_{r=1}^K \exp(\gamma_r + \beta'_r z) \right]^{-1}, \quad (2.17)$$

so that for $0 < t \leq \tau$, we have

$$p(t|z) = \exp[-\lambda(z)t] = \left[1 + \sum_{r=1}^K \exp(\gamma_r + \beta'_r z) \right]^{-t/\tau} = [p(z)]^{-t/\tau} \quad (2.18)$$

and

$$\lambda(z) = \frac{1}{\tau} \log \left[1 + \sum_{r=1}^K \exp(\gamma_r + \beta'_r z) \right] = -\frac{1}{\tau} \log p(z), \quad (2.19)$$

where $p(z)$ is defined in (1.6). (Compare also Part I, Section 1.2.)

Note that

$$q(t|z) = \int_0^t \lambda(z) p(u|z) du, \quad (2.20)$$

and

$$q_k^*(t|z) = \int_0^t \lambda_k(z) p(u|z) du. \quad (2.21)$$

But since [from (2.12) and (2.16)], we have

$$q_k^*(t|z) = c_k(z) q(t|z), \quad (2.22)$$

then

$$\lambda_k(t|z) = \lambda_k(z) = c_k(z) \lambda(z). \quad (2.23)$$

Therefore, we also have for $0 < t \leq \tau$

$$\begin{aligned} p(t|z) &= \prod_{k=1}^K \exp[-c_k(z) \lambda(z) t] \\ &= \prod_{k=1}^K \left\{ \exp[-\lambda(z) t] \right\}^{c_k(z)}. \end{aligned} \quad (2.24)$$

[Compare (2.6) and (2.10)]

Likelihood function

(i) Let \mathcal{D}_k denote the set of individuals who died from cause C_k , $\mathcal{D} = \bigcup_k \mathcal{D}_k$ - the set of individuals who died from any cause, S - the set of survivors, and t_j - the time at which individual (j) was last seen in $(0, \tau)$.

The likelihood function can be expressed in the form

$$\begin{aligned} L^{(2)}(\underline{Y}, \underline{\beta}) &= \prod_{k=1}^K \left\{ \prod_{j \in \mathcal{D}_k} c_k(z_j) \lambda(z_j) \exp[-\lambda(z_j) t_j] \right\} \times \prod_{j \in S} \exp[-\lambda(z_j) t_j] \\ &= \left[\prod_{k=1}^K \prod_{j \in \mathcal{D}_k} c_k(z_j) \right] \times \prod_{j \in \mathcal{D}} \lambda(z_j) \times \prod_{j=1}^N \exp[-\lambda(z_j) t_j] \end{aligned} \quad (2.25)$$

(ii) Equivalently, we may express the likelihood function (2.25) in a form using indicator variables.

Let

$$\delta_{kj} = \begin{cases} 1 & \text{if individual (j) dies from cause } C_k \text{ at time } t_j \\ 0 & \text{otherwise,} \end{cases} \quad (2.26)$$

We also have

$$\delta_{\cdot j} = \sum_{k=1}^K \delta_{kj} = \begin{cases} 1 & \text{if individual (j) dies at time } t_j \\ 0 & \text{if individual (j) is alive at time } t_j. \end{cases} \quad (2.27)$$

The likelihood function is

$$\begin{aligned} L^{(2)}(\underline{\gamma}, \underline{\beta}) &= \prod_{j=1}^N \left\{ \prod_{k=1}^K [c_k(z_j) \lambda(z_j)]^{\delta_{kj}} \exp[-\lambda(z_j) t_j] \right\} \\ &= \prod_{j=1}^N \prod_{k=1}^K [c_k(z_j)]^{\delta_{kj}} \times \prod_{j=1}^N [\lambda(z_j)]^{\delta_{\cdot j}} \times \prod_{j=1}^N \exp[-\lambda(z_j) t_j], \end{aligned} \quad (2.28)$$

so that

$$\log L^{(2)}(\underline{\gamma}, \underline{\beta}) = \sum_{j=1}^N \sum_{k=1}^K \delta_{kj} \log c_k(z_j) + \sum_{j=1}^N \delta_{\cdot j} \log \lambda(z_j) - \sum_{j=1}^N \lambda(z_j) t_j. \quad (2.29)$$

Substituting for $c_k(z_j)$ [from (1.9)] and for $\lambda(z_j)$ [from (2.19)], we obtain

$$\log L^{(2)}(\underline{\gamma}, \underline{\beta}) = D(\underline{\gamma}, \underline{\beta}) + A(\underline{\gamma}, \underline{\beta}) - \frac{1}{\tau} B(\underline{\gamma}, \underline{\beta}) - (N - N_{k+1}) \log \tau, \quad (2.30)$$

where

$$\begin{aligned} D(\underline{\gamma}, \underline{\beta}) &= \sum_{j=1}^N \sum_{k=1}^K \delta_{kj} (\gamma_k + \beta'_k z_j) - \sum_{j=1}^N \delta_{\cdot j} \log \left[\sum_{r=1}^K \exp(\gamma_r + \beta'_r z_j) \right]; \\ A(\underline{\gamma}, \underline{\beta}) &= \sum_{j=1}^N \delta_{\cdot j} \log \left\{ \log \left[1 + \sum_{r=1}^K \exp(\gamma_r + \beta'_r z_j) \right] \right\}; \\ B(\underline{\gamma}, \underline{\beta}) &= \sum_{j=1}^N t_j \log \left[1 + \sum_{r=1}^K \exp(\gamma_r + \beta'_r z_j) \right]. \end{aligned} \quad (2.31)$$

[Compare also Part I, Section 1.2.]

After some straightforward algebra the likelihood equations can be written in the form:

$$\frac{\partial \log L^{(2)}}{\partial \gamma_h} = N_h - \sum_{j=1}^N \delta_{\cdot j} \left[\frac{q_h^*(z_j)}{\log p(z_j)} + c_h(z_j) \right] - \frac{1}{\tau} \sum_{j=1}^N t_j q_h^*(z_j) = 0 \quad (2.32)$$

and

$$\begin{aligned} \frac{\partial \log L^{(2)}}{\partial \beta_{hv}} &= \sum_{j=1}^N \delta_{hj} z_{vj} - \sum_{j=1}^N \delta_{\cdot j} \left[\frac{q_h^*(z_j)}{\log p(z_j)} + c_h(z_j) \right] z_{vj} \\ &\quad - \frac{1}{\tau} \sum_{j=1}^N t_j q_h^*(z_j) z_{vj} = 0 \end{aligned} \quad (2.33)$$

for $h = 1, 2, \dots, K$; $v = 1, 2, \dots, s$; $0 < t_j \leq \tau$.

3. ESTIMATION FROM GROUPED DATA

If the period of investigation is fairly short and there are *no withdrawals*, then the likelihood function, $L^{(1)}(\gamma, \beta)$, defined in (1.12) may be approximately used in estimating β 's and γ 's. If there are *withdrawals* and the exact times of last seen, t_{ij} 's, are recorded, we may use $L^{(2)}(\gamma, \beta)$. If, however, the period of follow up is sufficiently long, an exponential approximation over the whole period might be inappropriate. If, in addition, the data represent a rather large set, we may group them into M fixed intervals and apply piecewise fitting of an exponential-logistic model.

For the interval $[t_i, t_{i+1})$ of the length $t_{i+1} - t_i = \tau_i$, we define [analogously to (2.15)]

$$\log \frac{q_{ki}^*(\tau_i | z)}{p_i(\tau_i | z)} = \log \frac{c_{ki}(z) \{1 - \exp[-\lambda_{\cdot i}(z) \tau_i]\}}{\exp[-\lambda_{\cdot i}(z) \tau_i]} = \gamma_{ki} + \beta'_{kz}. \quad (3.1)$$

Note that the γ_{ki} 's depend not only on the cause, but they also change from interval to interval, while the β_k 's are cause specific, but remain the same over all intervals. Consequently, the overall hazard rate in $[t_i, t_{i+1})$, $\lambda_{\cdot i}(z)$, depends on the interval.

Now, we may release the assumption that z 's do not change with time. Let $\tilde{z}_{ij} = (z_{1ij}, z_{2ij}, \dots, z_{s_{ij}})$ be the vector of z 's measured on the j th individual observed in the interval $[t_i, t_{i+1})$; let t_{ij} be the time at which the j th individual was last seen in $[t_i, t_{i+1})$, and $\tau_{ij} = t_{ij} - t_i$ denote the length of his exposure in this interval.

Similarly, as in Part I, Section 2, we may consider two kinds of likelihood function, depending whether the time at which the j th individual was last seen is recorded (or utilized) to the nearest interval or is 'exact' value, t_{ij} .

3.1. Times Last Seen Are Recorded to the Nearest Interval

For the interval $[t_i, t_{i+1})$, let \mathcal{D}_{ki} denote the set of individuals who died from cause C_k , and let $\mathcal{D}_i = \bigcup_k \mathcal{D}_{ki}$; further, let \mathcal{W}_i denote the set of withdrawals and \mathcal{S}_i - the set of survivors; finally, let $\gamma_{\cdot i} = (\gamma_{1i}, \gamma_{2i}, \dots, \gamma_{ki})$ for $i = 0, 1, \dots, M-1$.

Assuming that, on the average, the withdrawal time is in the middle of the interval, the contribution to the likelihood in the interval $[t_i, t_{i+1})$ is

$$L_{\cdot i}^{(3)}(\gamma_{\cdot i}, \beta) = \prod_{k=1}^K \left[\prod_{j \in \mathcal{D}_{ki}} q_{ki}^*(z_{ij}) \right] \times \prod_{j \in \mathcal{S}_i} p_i(z_{ij}) \times \prod_{j \in \mathcal{W}_i} [p_i(z_{ij})]^{1/2}, \quad (3.2)$$

where

$$p_i(z_{ij}) = \left[1 + \sum_{r=1}^K \exp(\gamma_{ri} + \beta'_r z_{r-ij}) \right]^{-1}, \quad (3.3)$$

and

$$q_{ki}^*(z_{ij}) = \exp(\gamma_{ki} + \beta'_k z_{ij}) \left[1 + \sum_{r=1}^K \exp(\gamma_{ri} + \beta'_r z_{ij}) \right]^{-1}, \quad (3.4)$$

for $k = 1, 2, \dots, K$; $i = 0, 1, \dots, M - 1$. [Compare (1.6) and (1.7), respectively.]

The overall likelihood function

$$L^{(3)}(\underline{\gamma}, \underline{\beta}) = \prod_{i=0}^{M-1} L_{\cdot i}^{(3)}(\underline{\gamma}_{\cdot i}, \underline{\beta}), \quad (3.5)$$

where $\underline{\gamma} = (\underline{\gamma}_{\cdot 0}, \underline{\gamma}_{\cdot 1}, \dots, \underline{\gamma}_{\cdot M-1})$.

Taking logarithms of both sides of (3.5) and substituting for $q_{ki}^*(z_{ij})$ and $p_i(z_{ij})$, the explicit form of the likelihood equations can be obtained in a straightforward manner. (Compare Part I, Section 2.1.)

3.2. Exact Times Last Seen Are Recorded

In this situation, a likelihood function of the kind (2.25) or equivalently (2.28) is appropriate to apply in each interval. Using, for instance, formula (2.28), the contribution to the likelihood in the interval $[t_i, t_{i+1})$ is

$$L_{\cdot i}^{(4)}(\underline{\gamma}_{\cdot i}, \underline{\beta}) = \prod_{j=1}^N \prod_{k=1}^K [c_{ki}(z_{ij})]^{\delta_{kij}} \times \prod_{j=1}^N [\lambda_{\cdot i}(z_{ij})]^{\delta_{\cdot ij}} \times \prod_{j=1}^N \exp[-\lambda_{\cdot i}(z_{ij})\tau_{ij}], \quad (3.6)$$

where $c_{ki}(z)$ is an analogue of $c_k(z)$ defined in (1.9), and $\lambda_{\cdot i}(z)$ is an analogue of $\lambda(z)$ defined in (2.19), but in both γ_k is replaced by γ_{ki} . Also δ_{kij} and $\delta_{\cdot ij}$ are analogues of δ_{kj} [see (2.26)], and $\delta_{\cdot j}$ [see (2.27)], respectively, for the interval $[t_i, t_{i+1})$.

The overall likelihood function is

$$L^{(4)}(\underline{\gamma}, \underline{\beta}) = \prod_{i=0}^{M-1} L_{\cdot i}^{(4)}(\underline{\gamma}_{\cdot i}, \underline{\beta}). \quad (3.7)$$

The likelihood equations can be derived in a straightforward manner. (Compare Part I, Section 2.2.)

4. ESTIMATION FROM UNGROUPED DATA: APPROXIMATION BY COX'S MODEL

When data are of relatively small extent, and also not many deaths are observed over a restricted period of investigation, the methods described in Section 3 are, in general, not very useful. First, we recall from Section 2.1 that we have introduced two kinds of mortality functions for cause C_k . The first is $Q_k^*(t|z)$ defined in (2.4) - that is the probability of death from cause C_k before time t in the presence of all causes acting simultaneously; these probabilities were employed in construction of likelihoods when the mortality from all causes is considered jointly (Sections 2 and 3). The second is the 'survival' function $G_k(t|z)$ defined in (2.7) which is a kind of 'waiting time' distribution function for cause C_k alone, assuming that the hazard rate is $\lambda_k(t|z)$. This function is more convenient to use for analysis of the clinical trial type data mentioned at the beginning of this Section.

For practical purposes, it seems to be sufficient to use only partial information on an event (death from cause C_k) which occurs at a given time point. This implies that only time at death from a particular cause, and the risk set at this time point supply the information for the contribution to likelihood. However, the exponential-logistic model is defined over a specified period or interval, and it would be difficult to decide what value of τ_i should be used for this type of data. Rather, we consider an approximation to this model when τ_i is small.

4.1. Approximation by Cox's Model

Consider exponential-logistic model defined in (2.15), with

$\pi_k(z) = c_k(z)$, and assume that τ is small.

Then

$$\begin{aligned} \log \frac{c_k(z)\{1 - \exp[-\lambda(z)\tau]\}}{\exp[-\lambda(z)\tau]} &\doteq \log\{c_k(z)[1 - \exp(-\lambda(z)\tau)]\} \\ &\doteq \log[c_k(z)\lambda(z)\tau] = \gamma_k + \beta_k'z, \end{aligned} \quad (4.1)$$

or

$$\lambda(z) = [c_k(z)\tau]^{-1} = \exp(\gamma_k + \beta_k'z), \quad (4.2)$$

Substituting for $c_k(z)$ [from (1.9)], we obtain

$$\lambda(z) = \sum_{k=1}^K \lambda_k \exp(\beta_k'z), \quad (4.3)$$

where $\lambda_k = \tau^{-1} \exp(\gamma_k)$ is a constant.

Since

$$\lambda(z) = \sum_{k=1}^K \lambda_k(z),$$

we obtain

$$\lambda_k(z) = \lambda_k \exp(\beta_k'z), \quad k = 1, 2, \dots, K, \quad (4.4)$$

which is a special case of Cox's (1972) model, with the underlying hazard rate $\lambda_k(t) = \lambda_k$ over a period $(0, \tau)$.

Note that now $\lambda_k(z)$, defined in (4.4), depends solely on the coefficients λ_k and β_k 's associated with cause C_k . These parameters can then be estimated from mortality data for cause C_k alone while treating the deaths from other causes as withdrawals.

4.2. Estimation of Survival Function: Likelihood Based on Density Function

We first assume that z 's do not depend on time, or are measured only once, usually at the time of individual's entry. Let

$$t'_{k1} < t'_{k2} < \dots < t'_{ki} < \dots < t'_{kn_k} \quad (4.5)$$

denote the observed ordered (and distinct) times at death and $\lambda_{ki}(z)$ denote the hazard rate for cause C_k in the interval $(t'_{k,i-1}, t'_{ki}]$ with $t'_{ki} - t'_{k,i-1} = \tau'_{ki}$. Further, let $j (= 1, 2, \dots, N)$ denote individual (j) participating for some time in the study; m_{ki} - the number of individuals who die "at" time t'_{ki} ; $j_f(ki)$ - individual (j_f) who is the f th to die at time t'_{ij} , ($f = 1, 2, \dots, m_{ki}$), and $R(t'_{ki})$ - the risk set at $t'_{ki} - 0$. Then the contribution to the likelihood at time point t'_{ki} is

$$L_{ki}^{(5)}(\gamma_{ki}, \beta_k) = \prod_{f=1}^{m_{ki}} \lambda_{ki}(z_{j_f(ki)}) \times \prod_{\ell \in R(t'_{ki})} \exp[-\lambda_{ki}(z_\ell) \tau'_{ki}]. \quad (4.6)$$

Substituting for $\lambda_{ki}(z) = \lambda_{ki} \exp(\beta'_k z)$ [cf., (4.4)], we obtain

$$L_{ki}^{(5)}(\gamma_{ki}, \beta_k) = \prod_{f=1}^{m_{ki}} \lambda_{ki} \exp(\beta'_k z_{j_f(ki)}) \times \prod_{\ell \in R(t'_{ki})} \exp[-\lambda_{ki} \exp(\beta'_k z_\ell) \tau'_{ki}], \quad (4.7)$$

and

$$L_{k \cdot}^{(5)}(\gamma_{k \cdot}, \beta_k) = \prod_{i=1}^{n_k} L_{ki}^{(5)}(\gamma_{ki}, \beta_k). \quad (4.8)$$

4.5. Conditional Partial Likelihood Function

When there are no tied observations, and we are primarily interested in estimating β_k 's, the conditional partial likelihood is similar to that described in Section 3.2 of Part I.

Let $j(ki)$ denote individual (j) who dies at time t'_{ki} . The conditional probability that individual (j) is the one who died at time t'_{ki} given a death from cause C_k takes place, is

$$\omega_j(t'_{ki}) = \frac{\exp(\beta'_k z_{j(ki)})}{\sum_{\ell \in R(t'_{ki})} \exp(\beta'_k z_\ell)},$$

and

$$L_{k^*}^{(6)}(\beta) = \prod_{i=1}^{n_k} \frac{\exp(\beta' z_{k^*j}(ki))}{\sum_{\lambda \in R(t'_{ki})} \exp(\beta' z_{k^*l})}, \quad (4.9)$$

which is, of course, the likelihood function for Cox's model [see e.g. Holt (1978)].

First, note that this does not depend on the underlying hazard rates, λ_{ki} . These parameters can be estimated from likelihood $L_{k^*}^{(5)}$. Second, z 's can be replaced by $z(t'_{ki})$'s, so that use of time dependent covariables is possible. Finally, we notice that the approach presented in this section - using only mortality data from cause C_k alone - is possible in this context, because we ignore the other events (deaths from other causes and withdrawals) between $t'_{k,i-1}$ and t'_{ki} .

5. CONCLUDING REMARKS

1. A general exponential-logistic model over a specified period of time was defined, and applied in construction of likelihood functions of mortality data of large extent from K different causes (Sections 2 & 3).

2. When the data arise from clinical trials and/or grouping is inappropriate, there are some difficulties with using a logistic model. It has been shown that sometimes Cox's model can be used, instead, as a fair approximation. This simplifies the calculations considerably, as well as facilitating them, since already existing programs for fitting Cox's model can be easily adapted.

REFERENCES

- Anderson, J. A. (1972). Separate logistic discrimination. *Biometrika* 59, 19-35.
- Cox, D. R. (1970). *The Analysis of Binary Data*, Methuen, London.
- Cox, D. R. (1972). Regression model and life tables (with Discussion). *Journal of the Royal Statistical Society, Series B* 34, 187-220.
- Elandt-Johnson, R. C. (1981). Time dependent logistic models in follow-up studies and clinical trials. I. Binary data. (submitted).
- Holt, J. D. (1978). Competing risk analyses with special reference to matched pair experiments. *Biometrika* 65, 159-165.