

**Regression and residual analysis
in linear models with interval censored data**

Dissertation

zur Erlangung des Grades

eines Doktors der Naturwissenschaften

der Universität Dortmund

Dem Fachbereich Statistik der Universität Dortmund

vorgelegt von

Rebekka Topp

Dortmund 2002

Gutachter:

Prof. Dr. G. Gómez

Prof. Dr. S. Schach

Prof. Dr. W. Urfer

Tag der mündlichen Prüfung: 19. Juli 2002

Introduction to interval censored data and overview of the two parts of the thesis

Interval censored data arises naturally in medical longitudinal follow-up studies in which the event of interest can not be easily observed, for instance cancer recurrence or the elevation of levels of a biomarker without noticeable symptoms. In these situations, the patients are usually examined at clinical visits that take place only in certain time intervals, and the event of interest may then occur between two consecutive clinical visits. Then, one observes only a certain time interval $[X_L, X_R]$ which is known to include the true time X of onset of the event of interest. This type of interval censoring is called interval censoring case II. As special cases it includes left censoring and right censoring for X_L equal to zero and X_R infinity, respectively. Another type of interval censoring occurs when the event is only known to be smaller or larger than an observed monitoring time. This kind of data is referred to interval censoring case I, or current status data. Finally, one speaks of doubly censored data if one observes $\min\{\max\{X, X_L\}, X_R\}$. For a more extensive review of the different types of interval censored data see Gómez et al. (2001b). In this thesis, interval censoring case II will be considered and the censoring intervals will be taken to be closed on both sides in order to account for exact observations.

An example for interval censored data is given in Betensky and Finkelstein (1999) who introduce the AIDS clinical trial group protocol 181, a natural history substudy of a comparative trial of three anti-pneumocystis drugs. The patients were monitored periodically for evidence of bacterial and viral infections, with the objective of understanding the relationship between these two events, and eventually the natural history of AIDS. Many patients

missed several of the prescheduled clinic visits, and when they returned to the hospital for examination, new laboratory indications for the two events were found. Thus, their times until occurrence of the bacterial or viral infection were censored into the time intervals between their last and their new clinic visits.

Another example is the AIDS clinical trial group protocol 359, a randomized clinical trial designed to compare six different anti-retroviral treatment regimens for HIV-infected persons who had previously failed on the protease inhibitor Indinavir (see Gulick et al., 2000). The patients were monitored periodically for their viral load levels with the aim to determine the time period these levels remained below the threshold of 500 viral copies/ml. It happened that the viral load levels climbed above the threshold between two consecutive clinic visits so that the exact time below 500 copies/ml was interval censored into the time interval $[X_1, X_2]$, where X_1 is the elapsed time between the first viral load observation below 500 copies/ml and the last observation before the viral load is subsequently observed to be above 500 copies/ml. Similarly, X_2 is the elapsed time between the visit prior to the first viral load observed below 500 copies/ml and the first visit that the viral load is subsequently observed to be above this threshold.

Methods for interval censored data have been strongly developed in the past decades. An approach for the estimation of the distribution function when the data is interval censored is found in the article by Peto (1973). Turnbull in 1976 presented a theory for nonparametrically estimating the distribution function of interval censored variables, incorporating in the estimation process the idea of self-consistency developed by Efron (1967). Turnbull's work had a strong impact on the further development of all kind of statistical methods for interval censored data, including the field of linear regression. The statistical properties of Turnbull's nonparametric maximum likelihood estimator (NPMLE) have been studied very extensively. Concerning uniqueness, consistency and asymptotic properties see for example Gentleman and Geyer (1994), Yu, Schick, Li and Wong (1998), Pan and Chappell (1999) or Yu, Li and Wong (2000). Resulting from problems in developing a distribution theory of Turnbull's NPMLE, Groeneboom and Wellner (1992) characterized the NPMLE using isotonic regression theory and thereof derived a distribution theory for it.

Some research has also been done on variance estimation of the estimated

survival function for interval censored data. Two methods for this problem are studied in Sun (2001). Since the underlying survival function can be assumed to be smooth in many applications, and the NPML as a step function does not efficiently use this information, some proposals for smooth estimation of the survival function for interval censored data have been made. See for example Li, Watkins and Yu (1997) or Pan (2000). Recently, an extension of Turnbull's NPML to the case of bivariate interval censored data was proposed by Betensky and Finkelstein (1999).

Concerning parameter estimation in linear models with interval censored data, Finkelstein and Wolfe in 1985 developed estimation theory for linear models when the response is interval censored. They proposed a semi-parametric approach using an EM algorithm for the maximization of the likelihood function under different parametric models for the covariate distribution, but without assuming a parametric form for the distribution of the response variable. Li and Pu (1999) applied a least squares approach to the log-linear model with interval censored response. For regression analysis with an interval censored covariates, Gómez, Espinal and Lagakos (2002) proposed a semiparametric approach by maximizing the data likelihood under the assumption of a normal distribution for the response. The covariate distribution is estimated nonparametrically via Turnbull's (1976) method. Recently, Gil, López-García, Lubiano and Montenegro (2001) considered linear relations between two interval censored variables by defining a metric for the distance between the observed values of the response and those predicted from the model.

The estimation of the regression parameters of a linear model is also considered in the first part of this thesis where a new estimation theory is presented for models with both interval censored response and covariate. Unlike Gil et al. (2001), it does not use certain distances between the observed and predicted data but is an extension of the method of Gómez et al. (2002) and considers a semiparametric maximum likelihood approach.

Closely related to linear model estimation is the field of residual analysis. In regression theory, the analysis of residuals is an integrated tool necessary to complete the process of fitting linear models. However, in connection with interval censored data, only very few research has been done. For proportional hazard models, Farrington (2000) derived interval censored counterparts to the right censored Cox-Snell, martingale, deviance, and Schoenfeld residuals.

For linear models, Gómez et al. (2002) proposed an intuitive definition of residuals coming from linear models that incorporate interval censored covariates. The second part of this thesis presents a new residual theory for regression analysis with interval censored covariates, which is shown to be superior to that proposed by Gómez et al. (2002).

Introduction

The first part of this thesis deals with linear regression analysis when both response and covariate are interval censored. Linear regression analysis is a statistical technique for investigating and modelling relationships between different variables. A statistical relation between two random variables (Y and Z , say) is defined such that one variable can be expressed in terms of a mathematical function of the other variable, for example $Y = f(Z) + \varepsilon$. In this case, Y is called the dependent variable or response, Z is the independent variable or covariate, and ε is an error term. To examine the linear relationship between Y and Z (or some more Z), an appropriate model should be chosen on the nature of the statistical relation and the variable types under consideration.

When saying a relationship between some variables is 'linear', this usually refers to linearity in the parameters. In contrast, the value of the highest power of the independent variable in the model is called the 'order' of the model. For example, $Y = \beta_0 + \beta_1 Z + \beta_2 Z^2 + \varepsilon$ is a second-order (in the covariate Z) linear (in the parameters β_i , $i = 0, 1, 2$) regression model. The ε are called 'model errors' and are a random component reflecting the inaccuracy of the relationship between the variables which can never be exact due to e.g. measurement errors in the observations.

The history of linear models can be traced back to the early 19th century where Legendre was the first to introduce a linear model. The principle for the determination of the unknown parameters β_i , $i = 0, 1, 2$, was to minimize the sum of squares of the residuals $e = Y - \beta_0 - \beta_1 Z - \beta_2 Z^2$. Among the various approaches of performing regression, the least squares method is probably the most widely used.

Applications of linear regression analysis are numerous and occur in almost every field, including engineering, physical sciences, economics, management, life and biological science, and the social sciences. In this thesis,

the main focus is on variables coming from the field of medicine, and more specifically, the interest will be on variables that are interval censored, that is, the response Y and the covariate Z are not observed directly but only known to lie in some interval $[Y_L, Y_R]$ and $[Z_L, Z_R]$, respectively.

Chapter 1 of this part of the thesis presents the statistical methods necessary for the development of the new regression theory. It contains an introduction of the theory for nonparametrically estimating the distribution function of interval censored variables, both in the one-dimensional case and the two-dimensional case. Furthermore, it introduces the regression method of Gómez et al. (2002) who proposed an approach for parameter estimation in linear models with exactly observed response and interval censored covariates. Their method will be extended in Chapter 2 when developing a new regression theory for the case that the response variable is interval censored as well. It uses a maximum likelihood approach for the estimation of the regression parameters while estimating at the same time the unknown distribution function of the interval censored covariate. The performance of the proposed method is assessed via a simulation study as described in Chapter 3. Finally, Chapter 4 contains a discussion of possible alternative approaches for the estimation of the regression parameters in the given context.

Contents

1	Methods for interval censored variables	11
1.1	Nonparametric estimation of the distribution of an interval censored variable	11
1.2	Nonparametric estimation of the distribution of two interval censored variables	14
1.3	Linear regression models with exactly observed response and interval censored covariate	16
1.3.1	Nonparametric estimation of \mathbf{w} when θ is known	17
1.3.2	Maximum likelihood estimation of θ when \mathbf{w} is known	19
2	Linear regression with interval censored response and covariate	21
2.1	Estimation procedure	22
2.2	Confidence intervals for the model parameters	28
2.3	Multiple regression	28
2.4	Model errors coming from the exponential family or Weibull distribution	29
2.4.1	Weibull distribution	32
3	Simulations	35
3.1	Simulation theory	37
3.2	Results of the simulations	38
4	Discussion of other approaches	43
4.1	Empirical approach	43
4.2	Least squares approach	45
5	Outlook	47

A	Derivation of the ML equations when the errors are normally distributed	49
B	Maple program for the calculation of approximate confidence intervals	55
C	Derivation of the MLE for the multiple regression setting	59
D	Derivation of the MLE when the errors come from the exponential family	63
E	Derivation of the MLE when the errors come from the Weibull distribution	67

Chapter 1

Methods for interval censored variables

This chapter gives an overview of the methods used in the development of the new regression theory for interval censored data. It describes density estimation in the context of interval censored random variables as introduced by Turnbull (1976) for the one-dimensional case, and generalized by Betensky and Finkelstein (1999) for the two-dimensional case. Furthermore, the regression theory for linear models with observed response and interval censored covariate as proposed by Gómez et al. (2002) is presented. Their method will be extended later to the case that both covariate and response are interval censored.

1.1 Nonparametric estimation of the distribution of an interval censored variable

Suppose X to be a continuous, interval censored random variable with distribution function F and realizations $x_i, i = 1, \dots, n$. Due to interval censoring, the x_i are not observed directly but only their respective censoring intervals $[x_{L_i}, x_{R_i}]$. These are known to include the true value x_i with probability one.

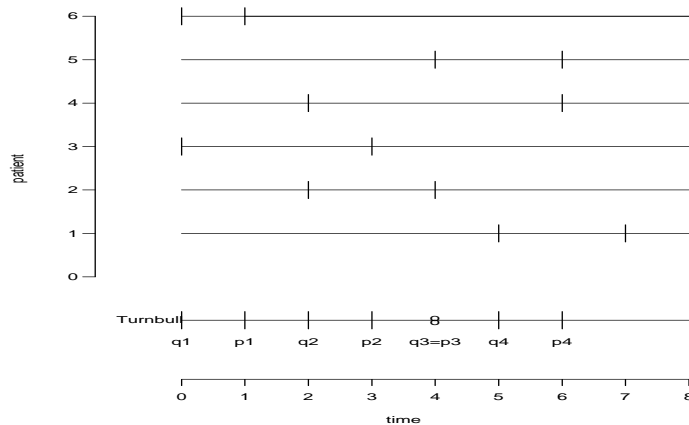
Turnbull (1976) proposed a maximum likelihood approach for determining an estimate for the distribution function F . It is a maximum likelihood approach which makes use of the equivalence between maximum likelihood estimates and self-consistent estimates as described in the following.

The construction of the likelihood for the data in the given context follows from the fact that the contribution of each individual i is $F(x_{R_i}) - F(x_{L_i})$, which results from X being interval censored. The complete likelihood accounting for all individuals is therefore given by

$$L(F) = \prod_{i=1}^n (F(x_{R_i}) - F(x_{L_i})).$$

Maximizing this likelihood with respect to F would yield the maximum likelihood estimate for the distribution function of X . Turnbull shows that this maximization problem can be reduced to a simpler one: After sorting all observed interval endpoints x_{L_i} and x_{R_i} in ascending order, one constructs a set of disjoint intervals $[q_1, p_1], \dots, [q_m, p_m]$ in the following way: Firstly, each $[q_j, p_j]$ must not contain any other member x_{L_i} or x_{R_i} except at their endpoints, and secondly, it must hold that $q_1 \leq p_1 < q_2 \leq \dots < q_m \leq p_m$. An example for the construction of the Turnbull intervals $[q_j, p_j]$, $j = 1, \dots, m$, is given in Figure 1.1. It shows six observed patient time intervals $[0,1]$, $[4,6]$, $[2,6]$, $[0,3]$, $[2,4]$, $[5,7]$ and the resulting Turnbull intervals $[0,1]$, $[2,3]$, $[4,4]$, $[5,6]$ obtained with the two construction rules given above.

Figure 1.1: Illustration of the construction of Turnbull's intervals



Turnbull proved that:

1. Any cumulative distribution function which increases outside the set $\bigcup_{j=1}^m [q_j, p_j]$ can not be a maximum likelihood estimate of F , and
2. for fixed values of $F(p_j+)$ and $F(q_j-)$, the likelihood is independent of the behavior of F within each interval $[q_j, p_j]$.

This means that it suffices to consider only those distribution functions which increase in some or all of the intervals $[q_j, p_j]$ and are constant outside these intervals. Furthermore, the behavior of the distribution function inside these intervals is not defined but can be imagined to be arbitrary. Thus, the problem of maximizing $L(F)$ reduces to that of maximizing

$$L(s_1, \dots, s_m) = \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} s_j,$$

where $s_j = F(p_j+) - F(q_j-)$ with $\sum_{j=1}^m s_j = 1$, and $\alpha_{ij} = 1$ if $[q_j, p_j] \subseteq [x_{L_i}, x_{R_i}]$ and 0 otherwise. The meaning of the indicator α_{ij} is that only those individuals contribute to the likelihood, whose observed censoring intervals contain one or more Turnbull intervals. The estimate of the density of X is given through the weight vector $\mathbf{s} = (s_1, \dots, s_m)$.

In order to determine the maximum likelihood estimate of \mathbf{s} , Turnbull proposed to apply an algorithm which is based on the equivalence between the maximum likelihood estimates and the self-consistent estimates and is described in the following. For details on the self-consistency equations see Efron, 1967.

Define $I_{ij}=1$ if $x_i \in [q_j, p_j]$ and 0 otherwise. Because of censoring the value of I_{ij} is not known, but its expectation is given by

$$E_s(I_{ij}) = \alpha_{ij} s_j = \mu_{ij}(\mathbf{s}).$$

That is, $\mu_{ij}(\mathbf{s})$ represents the probability that the i -th observation lies in $[q_j, p_j]$. Furthermore, the proportion of observations in the interval $[q_j, p_j]$ is

$$\sum_{i=1}^n \mu_{ij}(\mathbf{s}) / M(\mathbf{s}) = \pi_j(\mathbf{s}),$$

where

$$M(\mathbf{s}) = \sum_{i=1}^n \sum_{j=1}^m \mu_{ij}(\mathbf{s}).$$

The self-consistent estimate of the s_j is then defined to be any solution of the simultaneous equation

$$s_j = \pi_j(s_1, \dots, s_m).$$

Turnbull incorporates these formulas in an iterative procedure in order to derive the nonparametric estimate for the s_j :

Step 1: Chose initial estimates s_j^0 , $j = 1, \dots, m$. This can be any set of positive numbers summing to unity, e.g. $s_j = \frac{1}{m}$ for all j .

Step 2: Evaluate $\mu_{ij}(\mathbf{s}^0)$, $M(\mathbf{s}^0)$ and $\pi_j(\mathbf{s}^0)$ using the formulas given above.

Step 3: Obtain the improved estimates s_j^1 by setting $s_j^1 = \pi_j(\mathbf{s}^0)$.

Step 4: Return to Step B replacing \mathbf{s}^0 by \mathbf{s}^1 .

Step 5: Stop when the values of \mathbf{s}^1 and \mathbf{s}^0 do not differ anymore.

Turnbull shows that the algorithm converges monotonely for those initial vectors \mathbf{s}^0 that are close to the true density vector \mathbf{s} . Gentleman and Geyer (1994) provide easily verifiable conditions for the self-consistent estimator to be a maximum likelihood estimator and for checking whether the maximum likelihood estimate is unique.

1.2 Nonparametric estimation of the distribution of two interval censored variables

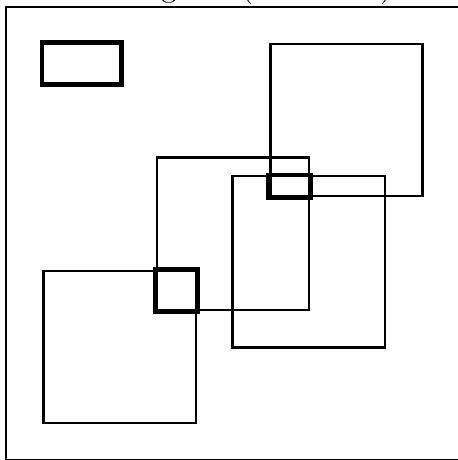
Betensky and Finkelstein (1999) generalized Turnbull's estimation procedure to bivariate discrete interval censored data. Unlike Turnbull, the likelihood function is not maximized using a self-consistent algorithm, but an extension of the method of Gentleman and Geier (1994) is applied.

In the bivariate case, one observes for each individual i , $i = 1, \dots, n$, the data rectangle $[x_{L_{1i}}, x_{R_{1i}}] \times [x_{L_{2i}}, x_{R_{2i}}]$ which are known to contain the realizations of X_{1i} and X_{2i} . Denoting $F(x_1, x_2)$ the joint cumulative distribution function of X_1 and X_2 , the likelihood for the data in this setting is

$$\prod_{i=1}^n (F(x_{R_{i1}}+, x_{R_{i2}}+) - F(x_{R_{i1}}+, x_{L_{i2}}-) - F(x_{L_{i1}}-, x_{R_{i2}}+) + F(x_{L_{i1}}-, x_{L_{i2}}-)).$$

Similar to the one-dimensional case, the support of the maximum likelihood estimate of F is contained in that set of rectangles which is formed by intersecting the observed data rectangles such that no other rectangle is contained within them. This mechanism is equivalent to the one used in the construction of the Turnbull intervals explained in the previous section. Figure 1.2 gives an illustration.

Figure 1.2: Final rectangles (thick lines), resulting from intersecting the observed regions (thin lines)



Denote the final rectangles as $[r_j, s_j] \times [t_j, u_j]$, $j = 1, \dots, J$. Define furthermore the probability associated with rectangle j to be $p_j = F(s_j+, u_j+) - F(s_j+, t_j-) - F(r_j-, u_j+) + F(r_j-, t_j-)$. Then, adopting the argumentation of Turnbull (1976), the search for the maximum likelihood estimate for F can be restricted to those vectors $\mathbf{p} = (p_1, \dots, p_J)$ having strictly non-negative components and summing to one. The maximum likelihood estimate even-

tually results from maximizing

$$L(\mathbf{p}) = \prod_{i=1}^n \sum_{j=1}^J \alpha_{ij} p_j,$$

where α_{ij} equals 1 if $[r_j, s_j] \subseteq [x_{L1_i}, x_{R1_i}]$ and $[t_j, u_j] \subseteq [x_{L2_i}, x_{R2_i}]$, and 0 otherwise.

Under the constraints for the p_j given above, the authors propose to maximize the likelihood $L(\mathbf{p})$ directly by solving a concave programming problem with linear constraints as described in Gentleman and Geier (1994).

1.3 Linear regression models with exactly observed response and interval censored covariate

Gómez et al. (2002) proposed a theory for linear regression analysis with interval censored covariates. The idea of their approach is to simultaneously maximize the data likelihood and estimate the unknown distribution function of the covariate.

The authors consider a continuous response variable Y with exactly observed realizations y_i , and a discrete and interval censored covariate Z whose realizations z_i are not observed but only the corresponding covariate intervals $[z_{L_i}, z_{R_i}]$, $i = 1, \dots, n$. These intervals are known to include z_i with probability one. The model to be established is

$$Y = \alpha + \beta Z + \varepsilon, \quad \text{model 1}$$

where the error term ε is said to be independent of Z and normally distributed with expectation zero and variance σ^2 . The aim is to estimate the parameter vector $\theta = (\alpha, \beta, \sigma^2)$ from the observed data $(y_i, [z_{L_i}, z_{R_i}])$.

Since Z is taken to be a discrete random variable, the authors suppose that it assigns positive mass w_j to the points s_j , $j = 1, \dots, m$. From the normality of the model errors follows that the conditional density f of the

response Y given s_j as a realization of Z is also normally distributed, with expectation $\alpha + \beta s_j$ and variance σ^2 :

$$f(y|s_j; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \alpha - \beta s_j)^2}{2\sigma^2}\right).$$

This density is used in the construction of the data likelihood

$$L(\theta, \mathbf{w}_j) = \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} w_j f(y_i|s_j; \theta),$$

where

$$\alpha_{ij} = \begin{cases} 1 & : s_j \in [z_{L_i}, z_{R_i}] \\ 0 & : s_j \notin [z_{L_i}, z_{R_i}] \end{cases}, \quad \text{and } w_j = P(Z = s_j).$$

Due to the unknown covariate distribution $\mathbf{w} = (w_1, \dots, w_m)$, this likelihood can not be maximized directly to obtain the maximum likelihood estimates for the model parameters. Therefore, the authors maximize L simultaneously for θ and \mathbf{w} using a two-step algorithm which first maximizes L with respect to \mathbf{w} for fixed θ , and then resolves the maximization problem for θ with \mathbf{w} known. These two steps are described in detail below.

1.3.1 Nonparametric estimation of \mathbf{w} when θ is known

Assuming that the value for θ is known, the maximization of the likelihood L reduces to the problem of finding a vector \mathbf{w} that maximizes

$$L^*(\mathbf{w}) = \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} w_j f(y_i|s_j),$$

subject to the constraints $\sum_{j=1}^m w_j = 1$ and $w_j \geq 0$ for all j .

The authors propose an algorithm for this maximization problem which is similar to Turnbull's density estimation procedure described in Chapter 1.1. It consists of the following steps: First, the authors fix a value for θ and chose start values for \mathbf{w} . With these, they calculate the probability ν_{ij} that the covariate of the i -th individual is equal to s_j . This quantity is then used to determine the expected number τ_j of individuals with $Z_i = s_j$. Finally,

τ_j is taken to be an improved estimate of the covariate density \mathbf{w} , and can later be used to recalculate ν_{ij} and τ_j . This procedure is repeated until the improved estimate and the old estimate are sufficiently close. The following scheme illustrates this estimation procedure:

Step 1a: Fix the value for θ using $\theta^0 = (\alpha^0, \beta^0, \sigma_0^2)$, where

$$\begin{aligned}\alpha^0 &= \bar{y} - \frac{\beta^0}{n} \sum_{i=1}^n \hat{e}_i, \\ \beta^0 &= \frac{\sum_{i=1}^n (y_i - \bar{y}) \hat{e}_i}{\sum_{i=1}^n (\hat{v}_i^2 - \hat{e}_i^2) - (1/n)(\sum_{i=1}^n \hat{e}_i)^2}, \\ n\sigma_0^2 &= \sum_{i=1}^n (y_i - \alpha^0)^2 - (\beta^0)^2 \sum_{i=1}^n (\hat{v}_i^2 + \hat{e}_i^2),\end{aligned}$$

and

$$\hat{e}_i = (x_{L_i} + x_{R_i})/2, \quad \text{and} \quad \hat{v}_i^2 = ((x_{L_i} - \hat{e}_i)^2 + (x_{R_i} - \hat{e}_i)^2)/2.$$

Step 1b: Chose initial estimates for the w_j^0 , for instance take $w_j^0 = \frac{1}{m}$.

Step 1c: Evaluate $\nu_{ij}(\theta, \mathbf{w}^0)$ defined as

$$\nu_{ij} := P(X = s_j | y_i, [x_{L_i}, x_{R_i}]) = \frac{\alpha_{ij} f(y_i | s_j; \theta) w_j}{\sum_{k=1}^m \alpha_{ik} f(y_i | s_k; \theta) w_k},$$

replacing w_j by w_j^0 , and calculate

$$\tau_j(\theta, \mathbf{w}^0) = \frac{1}{n} \sum_{i=1}^n \nu_{ij}(\theta, \mathbf{w}^0).$$

Step 1d: Obtain the improved estimate w_j^1 by setting $w_j^1 = \tau_j(\theta, \mathbf{w}^0)$.

Step 1e: Return to step 1c replacing \mathbf{w}^0 by \mathbf{w}^1 .

Step 1f: Repeat steps 1c to 1e until the value of \mathbf{w}^1 does not change anymore. Denote it by $\hat{\mathbf{w}}^1$.

1.3.2 Maximum likelihood estimation of θ when \mathbf{w} is known

When the covariate density \mathbf{w} is known, the maximization of the likelihood

$$L^{**}(\theta) = \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} w_j f(y_i|\theta)$$

with respect to θ can be achieved via the usual maximum likelihood approach: The logarithm of L^{**} is derived with respect to α , β and σ^2 , and these derivations are set to zero and solved for the parameters. The authors show that the solution of the maximum likelihood equations $\frac{\partial}{\partial \theta} \log L^{**}(\theta) = 0$ is

$$\hat{\alpha} = \bar{y} - \frac{\beta}{n} \sum_{i=1}^n e_i(\theta, \mathbf{w}), \quad (1)$$

$$\hat{\beta} = \frac{\sum_{i=1}^n (y_i - \alpha) e_i(\theta, \mathbf{w})}{\sum_{i=1}^n (v_i(\theta, \mathbf{w}) + e_i^2(\theta, \mathbf{w}))}, \quad (2)$$

$$n\hat{\sigma}^2 = \sum_{i=1}^n (y_i - \alpha)^2 - \beta^2 \sum_{i=1}^n (v_i(\theta, \mathbf{w}) + e_i^2(\theta, \mathbf{w})), \quad (3)$$

where

$$e_i(\theta, \mathbf{w}) = \frac{\sum_{k=1}^m \alpha_{ik} s_k w_k \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \alpha - \beta s_k)^2\right\}}{\sum_{k=1}^m \alpha_{ik} w_k \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \alpha - \beta s_k)^2\right\}}, \quad (4)$$

and

$$v_i(\theta, \mathbf{w}) = \frac{\sum_{k=1}^m \alpha_{ik} (s_k - \exp_i(\theta, \mathbf{w}))^2 w_k \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \alpha - \beta s_k)^2\right\}}{\sum_{k=1}^m \alpha_{ik} w_k \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \alpha - \beta s_k)^2\right\}}. \quad (5)$$

The algorithm proposed by the authors maximizes L^{**} by first choosing initial values for equations (4) and (5), which are then used to calculate the estimates given in (1) to (3). Afterwards, (4) and (5) are determined again using the newly calculated estimates and the covariate density vector that resulted from the algorithm of the previous section. This procedure is repeated until the values for $\hat{\alpha}$, $\hat{\beta}$ and $\hat{\sigma}^2$ stabilize. The following scheme illustrates the estimation process.

Step 2a: Calculate θ^0 from formulas (1) to (3) by choosing the initial values for $e_i(\theta, \mathbf{w})$ and $v_i(\theta, \mathbf{w})$ to be

$$e_i^0(\theta, \mathbf{w}) = \frac{z_{L_i} + z_{R_i}}{2} \quad \text{and}$$

$$v_i^0(\theta, \mathbf{w}) = \frac{(z_{L_i} - e_i^0)^2 + (z_{R_i} - e_i^0)^2}{2}.$$

Step 2b: Evaluate $e_i(\theta^0, \hat{\mathbf{w}}^1)$ and $v_i(\theta^0, \hat{\mathbf{w}}^1)$ using equations (4) and (5) employing θ^0 and $\hat{\mathbf{w}}^1$ from step 1f above.

Step 2c: Obtain the improved estimate θ^1 from formulas (1) to (3), replacing $e_i(\theta, \mathbf{w})$ and $v_i(\theta, \mathbf{w})$ by $e_i(\theta^0, \hat{\mathbf{w}}^1)$ and $v_i(\theta^0, \hat{\mathbf{w}}^1)$.

Step 2d: Return to step 2a replacing θ^0 by θ^1 .

Step 2e: Repeat steps 2a to 2c until the difference between θ^0 and θ^1 is sufficiently small. Denote the final estimate by $\hat{\theta}^1$.

In total, the two-step algorithm for calculating simultaneously the density \mathbf{w} of the interval censored covariate and the estimator for the parameter vector θ , results in the combination of the two algorithms given above and is summarized in the following scheme:

Step I: Execute Step 1a up to Step 1f.

Step II: Execute Step 2a up to Step 2e.

Step III: Return to Step 1c replacing θ^0 by $\hat{\theta}^1$ and \mathbf{w}^0 by $\hat{\mathbf{w}}^1$.

Step IV: Repeat steps I to III until convergence of θ and \mathbf{w} .

Chapter 2

Linear regression with interval censored response and covariate

This chapter presents a new estimation theory for linear regression models when both covariate and response are interval censored. It is an extension of the method of Gómez et al. (2002) introduced previously. The model to be considered here is

$$Y_i = \alpha + \beta Z_i + \varepsilon_i, \quad i = 1, \dots, n \quad \text{model 2}$$

where the response Y_i is continuous and censored into the interval $[Y_{L_i}, Y_{R_i}]$, and the covariate Z_i is discrete and censored into the interval $[Z_{L_i}, Z_{R_i}]$. The model errors ε are assumed to have a normal distribution with mean zero and variance σ^2 .

Let s_j be the possible values for Z with corresponding weights w_j , $j = 1, \dots, m$, and denote the covariate density and distribution function as \mathbf{w} and \mathbf{W} , respectively. From the errors' normal distribution follows that the distribution of Y given s_j as a value of Z is also normal with mean $\alpha + \beta s_j$ and variance σ^2 :

$$f(y|s_j, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \alpha - \beta s_j)^2}{2\sigma^2}\right).$$

Here, $\theta = (\alpha, \beta, \sigma^2)$ is the vector of the model parameters which we want to estimate.

It will be assumed that the interval censoring for the covariate and the response occurs noninformatively. If a variable X is subject to noninformative censoring, this means that for any given values x_0, x_1, x_2 , the conditional density of this variable is the same as the density of the uncensored variable truncated into the observed censoring interval:

$$P(X = x_0 | X_L = x_1, X_R = x_2) = \begin{cases} \frac{P(X=x_0)}{P(X \in [x_1, x_2])} & : x_0 \in [x_1, x_2] \\ 0 & : \text{otherwise} \end{cases}.$$

Gómez et al. (2001b) show that the contribution to the likelihood of an unique individual with observed censoring interval $[x_L, x_R]$ which includes the true value of interest x , is proportional to $\int_{x_L}^{x_R} dW(x)$ where $W = P(X \leq x)$. With this fact, the likelihood for the observed data of *model 2* can be constructed as given in the next section.

2.1 Estimation procedure

The observed data for *model 2* consists of n independent and identically distributed realizations of Y and Z . Since these two variables are interval censored, one observes the intervals $([y_{L_i}, y_{R_i}], [z_{L_i}, z_{R_i}])$, $i = 1, \dots, n$. In order to obtain the estimates for the model parameters α , β and σ^2 , a maximum likelihood approach will be proposed as described in the following.

The likelihood for the observed data can be constructed by noting the following facts: The contribution of an arbitrary individual i to the likelihood consists of the contribution of this individual with respect to both the covariate and the response. Since the covariate Z is interval censored, its density must be estimated with a method similar to the one given in Turnbull (1976), yielding as a result the weights w_j (for more details on the method of Turnbull see Chapter 1.1). Thus, the contribution of individual i with respect to Z is $\sum_{j=1}^m \alpha_{ij} w_j$, where the indicator variable α_{ij} specifies whether or not the covariate value s_j is contained in the observed covariate interval $[z_{L_i}, z_{R_i}]$. On the other hand, the contribution of this individual with respect to the response Y given a fixed value of Z , is determined by the conditional density $f(y|s_j, \theta)$. Since the value of Y is not exactly observed but only its censoring interval $[y_{L_i}, y_{R_i}]$, the conditional density must be integrated over the range of this censoring interval in order to obtain the respective contribution to the

likelihood. The total contribution of individual i to the likelihood is then the combination of these two single likelihood contributions, and the complete likelihood accounting for all individuals is therefore given by

$$\begin{aligned} L(\theta, w_j) &= \prod_{i=1}^n P(Y_i \in [Y_{L_i}, Y_{R_i}], Z_i \in [Z_{L_i}, Z_{R_i}]) \\ &= \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j; \theta) dy, \end{aligned} \quad (2.1)$$

where $\alpha_{ij} = \begin{cases} 1 & : s_j \in [z_{L_i}, z_{R_i}] \\ 0 & : \text{otherwise} \end{cases}$,

and $w_j = P(Z = s_j)$ is the weight the covariate assigns to the point s_j .

The estimation of the parameter vector θ will be achieved through maximizing L . Similar as in the context of the regression theory of Gómez et al. (2002), this maximization can not be carried out directly because of the unknown covariate density function $\mathbf{w} = (w_1, \dots, w_m)$. Thus, L is maximized through an algorithm that iterates between maximizing L with respect to \mathbf{w} while holding θ fixed, and maximizing L with respect to θ while holding \mathbf{w} fixed. These two steps are described in detail below.

Nonparametric estimation of \mathbf{w} when θ is known

For a fixed value of θ , the maximum likelihood estimate of the vector \mathbf{w} , given the constraints $\sum_{j=1}^m w_j = 1$ and $w_j \geq 0$ for all j , is determined by using a procedure based on the equivalence between the maximum likelihood and the self-consistent estimators as explained in Turnbull (1976): First, initial values for the covariate density weights w_j , are chosen. With these, the conditional probabilities ν_{ij} that the covariate Z_i equals a given value s_j are calculated. Summing these probabilities over all individuals i leads to the expected number τ_j of individuals with a covariate value equal to s_j . This expected number is then taken to be an improved estimate of the covariate density \mathbf{w} , and can be used to recalculate ν_{ij} and τ_j . The whole procedure is repeated until the difference of the values of the improved and the old estimate is sufficiently small. The following scheme gives a summary:

Step A1 Take initial estimates for the w_j^0 , for example $w_j^0 = \frac{1}{m}$ for $j = 1, \dots, m$. Denote $\mathbf{w}^0 = (w_1^0, \dots, w_m^0)$.

Step A2 Evaluate $\nu_{ij}(\mathbf{w}^0, \theta)$ and $\tau_j(\mathbf{w}^0, \theta)$ defined as

$$\begin{aligned} \nu_{ij}(\mathbf{w}^0, \theta) &= P(Z_i = s_j | [z_{L_i}, z_{R_i}], [y_{L_i}, y_{R_i}]) \\ &= \frac{\alpha_{ij} w_j^0 \int_{y_{L_i}}^{y_{R_i}} f(y | s_j; \theta)}{\sum_{j=1}^m \alpha_{ij} w_j^0 \int_{y_{L_i}}^{y_{R_i}} f(y | s_j; \theta)}, \end{aligned}$$

$$\tau_j(\mathbf{w}^0, \theta) = \frac{1}{n} \sum_{i=1}^n \nu_{ij}(\mathbf{w}^0, \theta).$$

Step A3 Obtain the improved estimate \mathbf{w}^1 setting

$$\mathbf{w}^1 = \tau_j(\mathbf{w}^0, \theta).$$

Step A4 Go to step A2 replacing \mathbf{w}^0 by \mathbf{w}^1 and repeat the whole procedure until their values are sufficiently close.

Maximum likelihood estimation of θ when \mathbf{w} is known

When the covariate density is known, the maximization of the likelihood L with respect to θ can be achieved by solving the score equation $\frac{\partial}{\partial \theta} \log L = 0$. The resulting estimates for α , β and σ^2 are derived in Appendix A. They are calculated to

$$\hat{\beta} = \frac{\bar{d} - \bar{a}\bar{b}}{\bar{c} - \bar{b}^2}, \quad (2.2)$$

$$\hat{\alpha} = \bar{a} - \hat{\beta}\bar{b}, \quad (2.3)$$

$$\hat{\sigma}^2 = \bar{e} - 2\hat{\alpha}\bar{a} + \hat{\alpha}^2 - \hat{\beta}^2\bar{c}, \quad (2.4)$$

where \bar{a} , \bar{b} , \bar{c} , \bar{d} and \bar{e} is the average of a_i , b_i , c_i , d_i and e_i , $i = 1, \dots, n$, respectively, defined as

$$a_i = E(Y_i | [Z_{L_i}, Z_{R_i}], [Y_{L_i}, Y_{R_i}]),$$

$$b_i = E(Z_i | [Z_{L_i}, Z_{R_i}], [Y_{L_i}, Y_{R_i}]),$$

$$\begin{aligned}
c_i &= E(Z_i^2 | [Z_{Li}, Z_{Ri}], [Y_{Li}, Y_{Ri}]), \\
d_i &= E(Z_i Y_i | [Z_{Li}, Z_{Ri}], [Y_{Li}, Y_{Ri}]), \\
e_i &= E(Y_i^2 | [Z_{Li}, Z_{Ri}], [Y_{Li}, Y_{Ri}]).
\end{aligned}$$

The following propositions show that the estimates $\hat{\alpha}$, $\hat{\beta}$ and $\hat{\sigma}^2$ are similar to the maximum likelihood estimators in a simple linear model with exactly observed response and covariate.

Proposition 1

It holds that $\hat{\beta}$ as defined above converges in probability to the value $\frac{Cov(Z, Y)}{Var(Z)}$.

Proof

Applying the law of large numbers, it holds that

$$\bar{a} = \frac{1}{n} \sum_{i=1}^n a_i \xrightarrow{n \rightarrow \infty} E(a_i) = E(E(Y_i | [Z_{Li}, Z_{Ri}], [Y_{Li}, Y_{Ri}])) = E(Y_i),$$

$$\bar{b} = \frac{1}{n} \sum_{i=1}^n b_i \xrightarrow{n \rightarrow \infty} E(b_i) = E(E(Z_i | [Z_{Li}, Z_{Ri}], [Y_{Li}, Y_{Ri}])) = E(Z_i),$$

$$\bar{c} = \frac{1}{n} \sum_{i=1}^n c_i \xrightarrow{n \rightarrow \infty} E(c_i) = E(E(Z_i^2 | [Z_{Li}, Z_{Ri}], [Y_{Li}, Y_{Ri}])) = E(Z_i^2),$$

and

$$\bar{d} = \frac{1}{n} \sum_{i=1}^n d_i \xrightarrow{n \rightarrow \infty} E(d_i) = E(E(Z_i Y_i | [Z_{Li}, Z_{Ri}], [Y_{Li}, Y_{Ri}])) = E(Z_i Y_i).$$

Thus, it holds for the numerator of $\hat{\beta}$ that

$$\bar{d} - \bar{a}\bar{b} \xrightarrow{n \rightarrow \infty} E(ZY) - E(Y)E(Z) = Cov(Z, Y),$$

and for the denominator that

$$\bar{c} - \bar{b}^2 \xrightarrow{n \rightarrow \infty} E(Z^2) - E(Z)^2 = Var(Z, Y).$$

In total, this means that

$$\hat{\beta} = \frac{\bar{d} - \bar{a}\bar{b}}{\bar{c} - \bar{b}^2} \xrightarrow{n \rightarrow \infty} \frac{Cov(Z, Y)}{Var(Y)}.$$

□

Proposition 2

It holds that $\hat{\alpha}$ as defined above converges in probability to the value $E(Y) - \beta E(Z)$.

Proof

Applying the law of large numbers it holds that

$$\bar{a} = \frac{1}{n} \sum_{i=1}^n a_i \xrightarrow{n \rightarrow \infty} E(a_i) = E(E(Y_i | [Z_{L_i}, Z_{R_i}], [Y_{L_i}, Y_{R_i}])) = E(Y_i)$$

and

$$\bar{b} = \frac{1}{n} \sum_{i=1}^n b_i \xrightarrow{n \rightarrow \infty} E(b_i) = E(E(Z_i | [Z_{L_i}, Z_{R_i}], [Y_{L_i}, Y_{R_i}])) = E(Z_i).$$

Thus, together with Proposition 1, this means that

$$\hat{\alpha} = \bar{a} - \hat{\beta} \bar{b} \xrightarrow{n \rightarrow \infty} E(Y) - \beta E(Z). \quad \square$$

Proposition 3

It holds that $\hat{\sigma}^2$ as defined above converges in probability to the value $Var(Y) - \beta Var(Z)$.

Proof

Again, with the law of large numbers and Proposition 1, one obtains

$$\begin{aligned} \hat{\sigma}^2 &\xrightarrow{n \rightarrow \infty} E(Y^2) - 2\hat{\alpha}E(Y) + \hat{\alpha}^2 - \beta^2 E(Z^2) \\ &= E(Y^2) - 2(E(Y) - \beta E(Z))E(Y) \\ &\quad + (E(Y) - \beta E(Z))^2 - \beta^2 E(Z^2) \\ &= E(Y^2) - E(Y)^2 - \beta^2 (E(Z^2) - E(Z)^2) \\ &= Var(Y) - \beta^2 Var(Z). \end{aligned} \quad \square$$

For the determination of the parameter estimates of *model 2*, a procedure is proposed that uses start values for \bar{a} to \bar{e} . It iterates between calculating $\hat{\alpha}$, $\hat{\beta}$ and $\hat{\sigma}^2$ and re-determining the values for \bar{a} to \bar{e} as explained in the scheme given below.

Step B1 Take initial estimates for a_i , b_i , c_i , d_i and e_i , for example

$$\begin{aligned} a_i^0 &= \frac{y_{L_i} + y_{R_i}}{2}, \\ b_i^0 &= \frac{z_{L_i} + z_{R_i}}{2}, \\ c_i^0 &= \frac{(z_{L_i} - b_i^0)^2 + (z_{R_i} - b_i^0)^2}{2}, \\ d_i^0 &= \frac{(z_{L_i} - b_i)(y_{L_i} - a_i) + (z_{R_i} - b_i)(y_{R_i} - a_i)}{2}, \\ e_i^0 &= \frac{(y_{L_i} - a_i^0)^2 + (y_{R_i} - a_i^0)^2}{2}. \end{aligned}$$

Step B2 Use these values in (2.2) to (2.4) to compute the initial estimate $\theta^0 = (\hat{\alpha}^0, \hat{\beta}^0, \hat{\sigma}_0^2)$.

Step B3 Re-evaluate a_i up to e_i with their theoretical formulas given in Appendix A by employing θ_0 .

Step B4 Obtain the improved estimate θ^1 by solving equations (2.2) to (2.4).

Step B5 Go to step B3 substituting θ^0 by θ^1 .

Step B6 Cycle steps B3 to B5 until the difference between the values of θ^0 and θ^1 is sufficiently small.

The complete algorithm to obtain the joint maximum likelihood estimate for \mathbf{w} and θ follows from the combination of the two conditional algorithms given above. It has been implemented in the program *semipara.cpp* and can be found on the floppy disc. The criteria for convergence of the estimates was chosen to be the relative norm differences of the estimates at iteration stage l :

$$\frac{\|\hat{w}^{l-1} - \hat{w}^l\|}{\|\hat{w}^{l-1}\|} \quad \text{and} \quad \frac{\|\hat{\theta}^{l-1} - \hat{\theta}^l\|}{\|\hat{\theta}^{l-1}\|}.$$

The estimates were defined to converge if the respective relative norm difference was less than 0.001. A flow-chart of the structure of this program is given in Chapter 3.

2.2 Confidence intervals for the model parameters

The MAPLE program given in Appendix B can be used to construct approximate confidence intervals for the parameter estimates resulting from the newly proposed estimation procedure. It uses the observed information matrix and quantiles of the normal distribution, and the different steps in the calculation process of the program are explained in the following:

Consider a given data set which consists of values y_{L_i} and y_{R_i} for the observed response intervals, values s_j for the discrete covariate with respective density weights w_j , and the estimated regression parameters $\hat{\alpha}$, $\hat{\beta}$ and $\hat{\sigma}^2$. The first part of the program reads this data into variables. With these, the log-likelihood as defined in equation (2.1) is constructed and its first and second derivatives with respect to the regression parameters are calculated. Then, the Hessian matrix is formed from all second derivatives and the observed information matrix is calculated by multiplying the Hessian with minus one. Eventually, the inversion of the observed information matrix provides an estimate for the variances of $\hat{\alpha}$, $\hat{\beta}$ and $\hat{\sigma}^2$. These estimated variances are then employed in the construction of the approximate confidence intervals.

2.3 Multiple regression

This section extends the proposed regression theory to the case that *model 2* additionally incorporates an exactly observed covariate vector. This means, the model now under consideration is

$$Y = \alpha + \vec{\beta}'_1 \vec{X} + \beta_2 Z + \varepsilon,$$

where $\vec{X} = (X_1, \dots, X_p)$ is a vector of exactly observed covariates, $\vec{\beta}'_1$ is the corresponding p -dimensional parameter vector, Y is the interval censored response, Z is an interval censored covariate, and ε is a continuous $N(0, \sigma^2)$ random variable independent of \vec{X} and Z .

The observed data for individual i is then $\vec{x}_i = (x_{1_i}, \dots, x_{p_i})'$, $[z_{L_i}, z_{R_i}]$ and $[y_{L_i}, y_{R_i}]$. By defining $\theta = (\alpha, \vec{\beta}'_1, \beta_2, \sigma^2)$ and using the notation and

assumptions of *model 2*, the likelihood function in the new context is given as

$$L_n^*(\mathbf{w}, \theta) = \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y | (\vec{x}_i, s_j); \theta), \quad (2.5)$$

where $\mathbf{w} = (w_1, \dots, w_m)$, $w_j = P(Z = s_j)$, $\alpha_{ij} = I\{s_j \in [z_{L_i}, z_{R_i}]\}$ and $f(y | (\vec{x}_i, s_j); \theta)$ is the conditional density of Y given $(\vec{X} = \vec{x}_i, Z = s_j)$:

$$f(y | (\vec{x}_i, s_j); \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \alpha - \vec{\beta}'_1 \vec{x}_i - \beta_2 s_j)^2}{2\sigma^2}\right).$$

The idea of the estimation procedure for the model parameters α , $\vec{\beta}'_1$, β_2 and σ^2 is the same as for *model 2*, only that the likelihood function is now given by (2.5). This means, L^* is maximized simultaneously for \mathbf{w} and θ by cycling between steps A and B of the earlier proposed algorithm. In the present context, Step A now consists of the same self-consistent equations as given earlier but using the new expression for $\nu_{ij}(\mathbf{w}, \theta)$, which is

$$\nu_{ij}(\mathbf{w}, \theta) = P(Z_i = s_j | [y_{L_i}, y_{R_i}], [z_{L_i}, z_{R_i}], \vec{x}_i) = \frac{\alpha_{ij} w_j \int_{y_{L_i}}^{y_{R_i}} f(y | (\vec{x}_i, s_j); \theta)}{\sum_{j=1}^m \alpha_{ij} w_j \int_{y_{L_i}}^{y_{R_i}} f(y | (\vec{x}_i, s_j); \theta)}.$$

Step B is modified in so far that it now incorporates the maximum likelihood estimators resulting from the new context of the multiple regression. These are obtained from maximizing the logarithm of likelihood (2.5) for fixed \mathbf{w} and are derived in Appendix C.

2.4 Model errors coming from the exponential family or Weibull distribution

In the previous sections, the regression parameters were estimated assuming the model errors to be normally distributed with mean zero and variance σ^2 . The normal distribution is known to be a member of the so-called exponential family of distributions, which is defined in the following way:

Definition

Let X be a random variable with density function f determined by the parameter vector η . One says that f belongs to the exponential family of distributions if it can be expressed as

$$f(x; \eta) = h(x)c(\eta)\exp[Q(\eta)t(x)],$$

where $Q(\eta)$ and $t(x)$ are vectors of common dimension k such that $Q(\eta)t(x) = \sum_{i=1}^k Q_i(\eta)t_i(x)$.

For example, the $N(0, \sigma^2)$ -distribution is obtained when taking $h(x) = 1$, $c(\eta) = (2\pi\sigma^2)^{-1/2}$, $Q(\eta) = (0, \frac{1}{2\sigma^2})$ and $t(x) = (x, -x^2)$. Other members of the exponential family are the gamma, binomial and Poisson distribution.

In what follows it will be shown that the proposed regression theory still holds when the model errors come from any distribution which is a member of the exponential family. This means that the likelihood to be considered now is

$$L^{**}(\theta, w_j) = \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) dy,$$

where

$$f(y_i | s_j; \theta) = h(y_i - \alpha - \beta s_j) c(\eta) \exp[Q(\eta)t(y_i - \alpha - \beta s_j)]$$

and $\theta = (\alpha, \beta, \eta)$.

The proceeding for obtaining the maximum likelihood estimate for θ in the new context is the same as in the original setting, namely maximizing the logarithm of L^{**} with respect to the parameters α , β and η . The resulting partial derivatives are given in Appendix D. It can be shown that the solutions (F1) – (F3) of Appendix D include equations (E1) – (E3) for normally distributed ε :

Corollary 1

Equation (F1) reduces to equation (E1) when the model errors are normal.

Proof

When the ε are normally distributed, it holds that $h(\varepsilon_i) = 1$, $c(\eta) = \frac{1}{\sqrt{2\pi\sigma^2}}$,

$Q(\eta) = -\frac{1}{2\sigma^2}$ and $t(\varepsilon_i) = (y_i - \alpha - \beta s_j)^2$. With that, equation (F1) results to

$$\begin{aligned}
(F1) &= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} -\frac{h'(\varepsilon_i)}{h(\varepsilon_i)} f(y_i | s_j; \theta) dy}{C_i(\theta)} \right. \\
&\quad \left. - \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) [Q(\eta) t'(\varepsilon_i)] dy}{C_i(\theta)} \right) \\
&= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} -\frac{0}{1} f(y_i | s_j; \theta) dy}{C_i(\theta)} \right. \\
&\quad \left. - \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) [-\frac{1}{2\sigma^2} (-2)(y_i - \alpha - \beta s_j)] dy}{C_i(\theta)} \right) \\
&= \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) [\frac{1}{\sigma^2} (y_i - \alpha - \beta s_j)] dy}{C_i(\theta)} = (E1) \quad \square
\end{aligned}$$

Corollary 2

Equation (F2) reduces to equation (E2) when the model errors are normal.

Proof

When the ε are normally distributed, it holds that $h(\varepsilon_i) = 1$, $c(\eta) = \frac{1}{\sqrt{2\pi\sigma^2}}$, $Q(\eta) = -\frac{1}{2\sigma^2}$ and $t(\varepsilon_i) = (y_i - \alpha - \beta s_j)^2$. With that, equation (F2) results to

$$\begin{aligned}
(F2) &= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} -\frac{s_j h'(\varepsilon_i)}{h(\varepsilon_i)} f(y_i | s_j; \theta) dy}{C_i(\theta)} \right. \\
&\quad \left. - \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) [Q(\eta) s_j t'(\varepsilon_i)] dy}{C_i(\theta)} \right) \\
&= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} -\frac{s_j 0}{1} f(y_i | s_j; \theta) dy}{C_i(\theta)} \right)
\end{aligned}$$

$$\begin{aligned}
& - \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) \left[-\frac{1}{2\sigma^2} (-2)(y_i - \alpha - \beta s_j) s_j \right] dy}{C_i(\theta)} \\
& = \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) \left[\frac{1}{\sigma^2} (y_i - \alpha - \beta s_j) s_j \right] dy}{C_i(\theta)} = (E2) \quad \square
\end{aligned}$$

Corollary 3

Equation (F3) reduces to equation (E3) when the model errors are normal.

Proof

When the ε are normally distributed, it holds that $h(\varepsilon_i) = 1$, $c(\eta) = \frac{1}{\sqrt{2\pi\sigma^2}}$, $Q(\eta) = -\frac{1}{2\sigma^2}$ and $t(\varepsilon_i) = (y_i - \alpha - \beta s_j)^2$. With that, equation (F3) results to

$$\begin{aligned}
(F3) &= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} -\frac{c'(\eta)}{c(\eta)} f(y_i | s_j; \theta) dy}{C_i(\theta)} \right. \\
&\quad \left. + \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) [Q'(\eta)t(\varepsilon_i)] dy}{C_i(\theta)} \right) \\
&= \sum_{i=1}^n \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \left(-\frac{c'(\eta)}{c(\eta)} + [Q'(\eta)t(\varepsilon_i)] \right) f(y_i | s_j; \theta) dy}{C_i(\theta)} \\
&= \sum_{i=1}^n \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \left(-\frac{1}{\sigma^2} + \left[\frac{1}{\sigma^4} (y_i - \alpha - \beta s_j)^2 \right] \right) f(y_i | s_j; \theta) dy}{C_i(\theta)} = (E3) \square
\end{aligned}$$

2.4.1 Weibull distribution

The proposed regression theory can also be applied when the model errors come from the Weibull distribution, as will be shown in the following. The likelihood of the data in this context is

$$L^{***}(\theta, w_j) = \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) dy,$$

where

$$f(y_i|s_j; \theta) = \alpha \beta s_j^{\beta-1} \exp(-\alpha s_j^\beta)$$

and $\theta = (\alpha, \beta)$.

Setting the partial derivatives of $l^{***} = \log L^{***}$ to zero and solving for α and β yields the maximum likelihood estimates given in Appendix E.

Chapter 3

Simulations

Since theoretical results for the goodness of the proposed estimates are difficult to obtain, their performance is checked through a simulation study. It involves different data scenarios for *model 2* with the aim to assess to what extent the proposed parameter estimates are able to reflect these data situations. Table 3.1 shows the simulation scenarios used in the study.

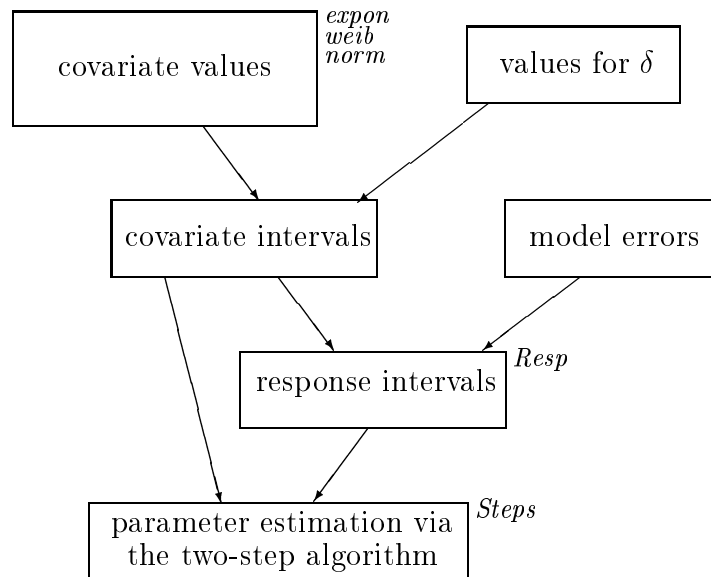
Table 3.1: Scenarios for the simulation study

number of observations	200 and 500
covariate distributions	$Exp(\frac{1}{8}), Weib(\frac{1}{6}, \frac{3}{2}), N(6, 4)$
percentage of censoring	0.3 and 0.7
value for α	4
values for β	2 and 5
value for σ^2	1

The simulations are carried out by the program *semipara.cpp* on the floppy disc, and a short summary of how this program works is given now: The model errors ε are generated from a $N(0, 1)$ -distribution, and the values for the covariate Z are simulated from the exponential, Weibull or normal distribution. These values are used to construct the covariate intervals $[Z_L, Z_R]$ after the following scheme: Depending on the covariate distribution, there is a certain number of values $j, j = 1, \dots, k$, which the covariate can take on. An indicator variable δ_{ij} determines with a given probability p , whether or not the covariate for individual i is observed at value j . Then, one looks at each value z_i and goes back to the nearest observed value j and takes

it as the value for z_{L_i} . Similarly, z_{R_i} is that observed value j coming first after z_i . The corresponding response intervals $[y_{L_i}, y_{R_i}]$ result from the formulas $y_{L_i} = \alpha + \beta z_{L_i} + \varepsilon_i$ and $y_{R_i} = \alpha + \beta z_{R_i} + \varepsilon_i$. Eventually, the two-step algorithm described in Chapter 2 is applied to the generated response and covariate intervals for the estimation of the model parameters α , β and σ^2 .

The following flow-chart illustrates the simulation process of the program *semipara.cpp*. The steps of the program are written inside the boxes and the arrows indicate which step enters in the calculation of another step. As most calculations are executed by procedures within the program, their names are written outside the corresponding box which will make it easier to find one's way when looking at the code of the program.



Other procedures used in this program are listed below together with a short description of their usage:

FileOpen: opens all files needed for reading and writing.

Spalloc: allocates memory for the vectors and matrices.

ran2: generates random uniform variates.

Simpson: integrates an user-defined function applying Simpson's method.

The last mentioned procedure *Simpson* is used for the calculation of different integrals over the conditional density $f(y|s_j; \theta)$ which is needed among others in the calculation of the conditional means a_i to e_i given in Chapter 2. As these integrals cannot be calculated analytically in C, a numerical approximation applying Simpson's method is used. The idea of Simpson's method is to approximate the area under a given graph by a sequence of quadratics. That is, the range of the upper and lower interval limit is divided into an even number of subintervals and their width is calculated. Then, the function value at the left endpoints of the first three subintervals is calculated as well as the area of the parabola through these three points. This process is repeated moving two subintervals to the right. Simpson's method is said to be the most exact among those existing for numerical integration. Though, it is obviously not as exact as the analytical form. This must be taken into consideration when assessing the simulation results of the estimates.

The performance of the program *semipara.cpp* with respect to speed and convergence is highly satisfying. Running it on a 400 megahertz Pentium II processor with 128 MB RAM main memory using the SUSE LINUX 7.1 operating system yielded convergence of the parameter estimates after 5 to 30 iterations depending on the number of observations and the level of censoring. The time needed for the calculations varied between 5 and 60 seconds.

3.1 Simulation theory

The simulation study involves the generation of data coming from different statistical distributions. The theory applied for the generation of these distributions is given now (for references see Box and Müller, 1958, or Morgan, 1984).

Uniform distribution

For the generation of a Uniform(0,1) random variable, a *Congruential Pseudo-Random Number Generator* is used. By applying the recursion formula $x_{n-1} = ax_n + b \text{ mod } m$ with seed x_0 and a, b, m given numbers, a sequence of integers will be obtained, each of which lies between 0 and $m - 1$. An approximation to Uniform(0,1) random variables u_i can then be achieved by setting $u_i = x_i/m$.

Exponential and Weibull distribution

As the Exponential and Weibull distributions are continuous, one can make use of the *Inversion Method* to generate their distribution functions. Suppose one wishes to simulate a continuous random variable X with distribution function $F(x) = P(X \leq x)$, and suppose further that the inverse function $F^{-1}(u)$ is well-defined for $u \in [0, 1]$. Then, it is well known that if U is a $(0, 1)$ -Uniform random variable, $X = F^{-1}(U)$ has the required distribution.

Normal distribution

For the simulation of the Normal distribution, the *Polar Marsagliar Method* is applied: If U is a Uniform(0,1) random variable, then $V = 2U - 1$ is a Uniform(-1,1) random variable. By selecting two independent Uniform(-1,1) random variables V_1 and V_2 , a random point in the square $[-1, 1] \times [-1, 1]$ can be specified which has polar coordinates (\tilde{R}, Θ) given by $\tilde{R}^2 = V_1^2 + V_2^2$ and $\tan(\Theta) = V_2/V_1$. The repeated selection of such points provides a random scatter of points inside this square, and rejection of points outside the unit-circle produces a uniform random scatter of points within this circle. For any of these points, the polar coordinates \tilde{R} and Θ are independent random variables, Θ is a Uniform(0,2 π) random variable and \tilde{R}^2 is a Uniform(0,1) random variable. One can write

$$\sin(\Theta) = \frac{V_2}{\tilde{R}} = \frac{V_2}{\sqrt{V_1^2 + V_2^2}}, \quad \cos(\Theta) = \frac{V_1}{\sqrt{V_1^2 + V_2^2}}.$$

Eventually, a pair of independent $N(0, 1)$ -variables is obtained by defining M_1 and M_2 as

$$M_1 = \sqrt{-2\log(\tilde{R}^2)} \frac{V_2}{\sqrt{V_1^2 + V_2^2}}, \quad M_2 = \sqrt{-2\log(\tilde{R}^2)} \frac{V_1}{\sqrt{V_1^2 + V_2^2}}.$$

3.2 Results of the simulations

Table 3.2 and 3.3 show the results of the simulation study for *model 2* under the different scenarios given in Table 3.1 above. Each column gives the median and mean value [standard deviation] calculated using 500 replicates for the estimated model parameters.

Table 3.2: Estimated regression parameters when $\alpha = 4$, $\beta = 2$ and $\sigma^2 = 1$

	Median	Mean [Std]	Median	Mean [Std]	Median	Mean [Std]
	for α		for β		for σ^2	
Exponential($\frac{1}{8}$)						
n=200,p=0.3	3.801	3.799 [0.228]	2.011	2.011 [0.032]	1.193	1.192 [0.130]
n=500,p=0.3	3.827	3.823 [0.141]	2.007	2.007 [0.032]	1.199	1.201 [0.084]
n=200,p=0.7	3.971	3.972 [0.159]	1.997	1.997 [0.021]	0.994	1.100 [0.111]
n=500,p=0.7	3.974	3.977 [0.099]	1.997	1.997 [0.013]	1.005	1.010 [0.068]
Weibull($\frac{1}{6}, \frac{3}{2}$)						
n=200,p=0.3	4.030	4.028 [0.246]	1.973	1.972 [0.069]	1.309	1.317 [0.134]
n=500,p=0.3	4.043	4.033 [0.163]	1.971	1.973 [0.044]	1.330	1.327 [0.091]
n=200,p=0.7	3.977	3.979 [0.183]	1.999	1.999 [0.049]	0.958	0.961 [0.101]
n=500,p=0.7	3.981	3.981 [0.117]	2.001	2.000 [0.032]	0.978	0.980 [0.071]
Normal(6,4)						
n=200,p=0.3	4.219	4.215 [0.497]	1.937	1.940 [0.085]	0.950	0.945 [0.118]
n=500,p=0.3	4.213	4.223 [0.303]	1.939	1.938 [0.052]	0.948	0.952 [0.069]
n=200,p=0.7	4.055	4.033 [0.358]	1.983	1.984 [0.059]	0.930	0.933 [0.105]
n=500,p=0.7	4.059	4.058 [0.222]	1.980	1.981 [0.037]	0.931	0.938 [0.069]

Table 3.3: Estimated regression parameters when $\alpha = 4$, $\beta = 5$ and $\sigma^2 = 1$

	Median	Mean [Std]	Median	Mean [Std]	Median	Mean [Std]
	for α		for β		for σ^2	
Exponential($\frac{1}{8}$)						
n=200,p=0.3	3.531	3.510 [0.288]	5.062	5.064 [0.045]	1.882	1.879 [0.235]
n=500,p=0.3	3.559	3.549 [0.180]	5.056	5.058 [0.028]	1.866	1.885 [0.138]
n=200,p=0.7	3.944	3.939 [0.168]	5.004	5.003 [0.022]	1.093	1.093 [0.123]
n=500,p=0.7	3.952	3.951 [0.100]	5.003	5.002 [0.012]	1.106	1.105 [0.072]
Weibull($\frac{1}{6}, \frac{3}{2}$)						
n=200,p=0.3	3.817	3.817 [0.306]	5.040	5.043 [0.090]	2.233	2.253 [0.303]
n=500,p=0.3	3.836	3.833 [0.193]	5.038	5.039 [0.056]	2.267	2.258 [0.203]
n=200,p=0.7	3.970	3.974 [0.198]	5.005	5.004 [0.053]	1.042	1.042 [0.110]
n=500,p=0.7	3.972	3.974 [0.118]	5.004	5.005 [0.032]	1.078	1.076 [0.071]

Normal(6,4)						
n=200,p=0.3	4.423	4.409 [0.568]	4.921	4.920 [0.100]	1.283	1.294 [0.166]
n=500,p=0.3	4.436	4.433 [0.347]	4.918	4.917 [0.060]	1.277	1.282 [0.100]
n=200,p=0.7	4.107	4.125 [0.378]	4.971	4.969 [0.063]	0.960	0.970 [0.110]
n=500,p=0.7	4.124	4.124 [0.232]	4.970	4.970 [0.039]	0.984	0.984 [0.075]

Both tables show that the values of the median and the mean do not differ much within the simulation scenarios. For $\beta = 2$, the estimation results for the parameter α are best when the covariate distribution is Weibull. For an exponential covariate distribution, this parameter is slightly underestimated, and for a normal distribution it is slightly overestimated. It can be also noticed that the standard deviation is twofold when the covariate distribution is normal. The estimation of the parameter β is very accurate for all covariate distributions and the standard deviations are also smaller than those for the parameter α . The results for the estimation of the error variance σ^2 is most satisfying for an exponential and Weibull covariate distribution with a low level of censoring ($p = 0.7$). At a high censoring level, the value of the error variance is overestimated. The results for a normally distributed covariate are similar for both low and high censoring levels but generally underestimate the error variance.

For $\beta = 5$, the estimation results for the parameter α are most satisfying when the percentage of censored data is low, regardless of the covariate distribution. When the percentage of censoring is high, the value of α is underestimated in case of the exponential and Weibull distribution, and overestimated in case of the normal distribution. Among these three covariate distributions, the Weibull performs best. With respect to the model parameter β , the simulation results show that the estimation procedure performs well for all three covariate distributions and estimates close to the true parameter value are obtained. The error variance σ^2 is estimated most satisfactorily for a low censoring level, otherwise it is overestimated. The value of the slope β has obviously an effect in the estimation of the error variance because the overestimation was not that high for $\beta = 2$.

It can be also noticed that the number of observations affects the value of the standard deviation of the estimates in so far that it gets smaller if the

number of observations gets larger.

Table 3.4 gives a summary of those simulation scenarios for which the parameter estimates perform best.

Table 3.4: Summary of the simulation results

	best performance for $\beta = 2$	best performance for $\beta = 5$
$\hat{\alpha}$	Weib, exp/norm and p=0.7	exp/norm/Weib and p=0.7
$\hat{\beta}$	all scenarios	all scenarios
$\hat{\sigma}^2$	exp/Weib and p=0.7	exp/norm/Weib and p=0.7

Chapter 4

Discussion of other approaches

Two other approaches for the estimation problem of *model 2* were investigated in addition to the semiparametric approach described in Chapter 2. The first approach is an empirical one with the idea of adapting the well-known uncensored regression estimators to the context of interval censored data. The second approach imitates the least squares method of uncensored regression analysis and transfers it to the interval censored setting. The following sections summarize the problems encountered in the process of examining these approaches.

4.1 Empirical approach

Consider the linear model $Y = \alpha + \beta Z + \varepsilon$ where Y is the response variable and Z the covariate, both uncensored. It is known from regression theory that for this model the least squares estimates

$$\hat{\beta} = \frac{\widehat{cov}(Y, Z)}{\widehat{var}(Z)} \quad \text{and} \quad \hat{\alpha} = \hat{E}(Y) - \hat{\beta}\hat{E}(Z) \quad (*)$$

are unbiased and have minimum variance when the conditions of the Gauss-Markov theorem are met.

When Y and Z are interval censored, one could think in trying to estimate the involved covariance, variance and expected values through the common density function of Z and Y , which can be calculated with the method developed by Betensky and Finkelstein (1999) described in Chapter 1.2. From the

estimated common density \hat{h} , say, one could then calculate the marginal densities \hat{f} and \hat{g} , say, of Y and Z , respectively. From these three distribution functions one could finally estimate the covariance, variance and expected values from

$$\hat{E}(Z) = \int_{Z_{L_i}}^{Z_{R_i}} z g(z) dz, \quad \hat{E}(Y) = \int_{Y_{L_i}}^{Y_{R_i}} y f(y) dz,$$

$$v\hat{a}r(Z) = \int_{Z_{L_i}}^{Z_{R_i}} (z - \hat{E}(Z))^2 g(z) dz,$$

$$c\hat{o}v(Y, Z) = \int_{Z_{L_i}}^{Z_{R_i}} \int_{Y_{L_i}}^{Y_{R_i}} (z - \hat{E}(Z))(y - \hat{E}(Y)) h(y, z) dy dz,$$

and calculate the estimators $\hat{\alpha}$ and $\hat{\beta}$ with the formulas given in (*).

Simulations using the same simulation scenarios as in the semiparametric approach showed that the estimates for α resulting from the empirical approach are not very accurate. Table 4.1 below gives the means [mean squared errors] of $\hat{\alpha}$ and $\hat{\beta}$, calculated from 1000 replications of each setting.

Table 4.1: Simulation results for the empirical approach where $\alpha = 4$ and $\beta = 2$

distribution	parameters	$\hat{\alpha}$ [MSE]	$\hat{\beta}$ [MSE]
Exponential ($\frac{1}{8}$)	n=100, p=0.3	4.390 [0.41]	2.086 [0.02]
	n=500, p=0.3	4.329 [0.41]	2.086 [0.02]
	n=100, p=0.7	4.547 [0.40]	2.022 [<0.01]
	n=500, p=0.7	4.478 [0.25]	2.024 [<0.01]
Weibull ($\frac{1}{6}, \frac{3}{2}$)	n=100, p=0.3	4.611 [0.52]	1.984 [0.03]
	n=500, p=0.3	4.662 [0.47]	1.939 [<0.01]
	n=100, p=0.7	4.593 [0.42]	1.995 [<0.01]
	n=500, p=0.7	4.559 [0.33]	1.976 [<0.01]
Normal (6,4)	n=100, p=0.3	4.145 [0.98]	2.112 [0.05]
	n=500, p=0.3	4.099 [0.20]	2.103 [0.02]
	n=100, p=0.7	4.168 [0.48]	2.111 [0.03]
	n=500, p=0.7	4.091 [0.10]	2.103 [0.01]

It can be seen that the value of α is strongly overestimated when the covariate distribution is exponential or Weibull. Only in case of a normally distributed covariate, this estimate is near the true value. The mean squared error is quite high for all three covariate distributions, so it must be concluded that the values of the estimator differ considerably within the 1000 replications. With respect to the parameter β , the simulation results show that the estimates are quite accurate and the mean squared errors are small.

One could conclude from Table 4.1 that the estimation results for a normally distributed covariate are not too bad, but this conclusion is not very appropriate due to the high mean squared errors for $\hat{\alpha}$. Furthermore, the estimation results are only stable when the number of observation is very high ($n = 500$), which does indicate a poor performance on small data sets. Also, the percentage of censoring effects the value of the mean squared error, but the influence seems not to be as high as that of the number of observations, especially in the case of a normally distributed covariate.

The main disadvantage, though, of the empirical approach is that it does not provide an estimate for the model error variance σ^2 . In the uncensored data setting, $\hat{\sigma}^2$ is calculated from the formula

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta}z_i)^2,$$

which has no proper equivalent in the interval censored data setting. The method of replacing the unobserved values y_i and z_i by the midpoints of their observed censoring intervals is generally known to lead to considerable biases in the estimators and is also not a methodologically correct approach.

4.2 Least squares approach

The least squares method in uncensored regression analysis achieves parameter estimation by minimizing the sum of squares

$$\sum_{i=1}^n (y_i - \alpha - \beta z_i)^2,$$

that is, the vertical distances between the observed data points and the fitted line. One could think in applying this method to the interval censored data

setting by minimizing the distances between the observed data rectangles and the fitted line. To avoid the definition of such a distance, one could directly try to minimize

$$\sum_{i=1}^n E \left((y_i - \alpha - \beta z_i)^2 \mid z_i \in [z_{L_i}, z_{R_i}], y_i \in [y_{L_i}, y_{R_i}] \right),$$

which is the expected sum of squares conditioned on the observed data rectangles $[z_{L_i}, z_{R_i}] \times [y_{L_i}, y_{R_i}]$. This would be equivalent to minimizing

$$\sum_{i=1}^n \int_{z_{L_i}}^{z_{R_i}} \int_{y_{L_i}}^{y_{R_i}} (y - \alpha - \beta z)^2 h_i(z, y) dy dz, \quad (**)$$

where $h_i(z, y)$ is the joint density of Z and Y truncated into the rectangle $[z_{L_i}, z_{R_i}] \times [y_{L_i}, y_{R_i}]$.

The solution of this equations would require the calculation of the density h_i , which can be achieved with the method of Betensky and Finkelstein (1999), as well as the mathematical minimization of the given sum with respect to the parameters α and β , which could be carried out by a mathematical software like MAPLE. For the purpose of running simulations in order to assess the performance of the estimators, the problem occurs how to connect these two steps so that they can be executed consecutively by the computer without interference from the outside. This problem could not be solved until now because of two facts: The MAPLE software is too inefficient to calculate the common density h_i , and the C language can not be used to solve minimization problems. Trying to calculate first h_i in C and then solving the minimization problem in MAPLE fails because it does not seem to exist a command that automatically starts a MAPLE program from the C interface. Theoretical calculations of the properties of the parameter estimates resulting from minimizing (**) are quite complex and difficult to interpret.

Chapter 5

Outlook

For the purpose of assessing the goodness of the estimated *model 2*, a residual theory should be developed in the future to complete the proposed regression theory. It is not sufficient to consider an ad-hoc approach like Gómez et al. (2002) did, because it could be seen from the results of the simulation study in Chapter 3 that these residuals perform quite unsatisfactorily in most of the considered data situations. It is rather desirable to extend the concept of the residual theory given in Part II of this thesis to the case that the response variable is interval censored as well.

Appendix A

Derivation of the ML equations when the errors are normally distributed

Consider the likelihood function

$$L = \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) w_j dy = \prod_{i=1}^n \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) dy,$$

where α_{ij} equals one if $s_j \in [z_{L_i}, z_{R_i}]$ and zero elsewhere. $\theta = (\alpha, \beta, \sigma^2)$ is the parameter vector to be estimated, and $f(y|s_j, \theta)$ is given by

$$f(y|s_j, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y - \alpha - \beta s_j)^2\right).$$

Define

$$L := \prod_{i=1}^n C_i(\theta),$$

where $C_i(\theta)$ is the contribution of the i -th individual to the likelihood L . Then,

$$\log L = \sum_{i=1}^n \log C_i(\theta).$$

In order to get the ML-estimators of θ , the ML equations are solved:

$$(E1) \quad \frac{\partial \log L}{\partial \alpha} = 0,$$

$$(E2) \quad \frac{\partial \log L}{\partial \beta} = 0,$$

$$(E3) \quad \frac{\partial \log L}{\partial \sigma^2} = 0.$$

Consider the quantities a_i , b_i , c_i , d_i and e_i defined as

$$a_i := E(Y|[Z_{L_i}, Z_{R_i}], [Y_{L_i}, Y_{R_i}]) = \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} y f(y|s_j, \theta) dy}{C_i(\theta)},$$

$$b_i := E(Z|[Z_{L_i}, Z_{R_i}], [Y_{L_i}, Y_{R_i}]) = \frac{\sum_{j=1}^m s_j \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) dy}{C_i(\theta)},$$

$$c_i := E(Z^2|[Z_{L_i}, Z_{R_i}], [Y_{L_i}, Y_{R_i}]) = \frac{\sum_{j=1}^m s_j^2 \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) dy}{C_i(\theta)},$$

$$d_i := E(ZY|[Z_{L_i}, Z_{R_i}], [Y_{L_i}, Y_{R_i}]) = \frac{\sum_{j=1}^m s_j \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} y f(y|s_j, \theta) dy}{C_i(\theta)},$$

$$e_i := E(Y^2|[Z_{L_i}, Z_{R_i}], [Y_{L_i}, Y_{R_i}]) = \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} y^2 f(y|s_j, \theta) dy}{C_i(\theta)}.$$

Then, solving equation (E1) leads to

$$\begin{aligned} (E1) &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \frac{y - \alpha - \beta s_j}{\sigma^2} f(y|s_j, \theta) dy = 0 \\ &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} y f(y|s_j, \theta) dy \\ &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} (\alpha + \beta s_j) f(y|s_j, \theta) dy \end{aligned}$$

$$\begin{aligned}
&\Leftrightarrow \sum_{i=1}^n a_i = \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} (\alpha + \beta s_j) f(y|s_j, \theta) dy \\
&\Leftrightarrow \sum_{i=1}^n a_i = \alpha \sum_{i=1}^n \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) dy}{C_i(\theta)} \\
&\quad + \beta \sum_{i=1}^n \frac{\sum_{j=1}^m s_j \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) dy}{C_i(\theta)} \\
&\Leftrightarrow \sum_{i=1}^n a_i = n\alpha + \beta \sum_{i=1}^n b_i \Leftrightarrow n\alpha = \sum_{i=1}^n a_i - \beta \sum_{i=1}^n b_i \\
&\Rightarrow \hat{\alpha} = \frac{1}{n} \sum_{i=1}^n a_i - \hat{\beta} \frac{1}{n} \sum_{i=1}^n b_i = \bar{a} - \hat{\beta} \bar{b}.
\end{aligned}$$

Equally, solving equation (E2) results in

$$\begin{aligned}
(E2) &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \frac{y - \alpha - \beta s_j}{\sigma^2} s_j f(y|s_j, \theta) dy = 0 \\
&\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j s_j \int_{Y_{L_i}}^{Y_{R_i}} (y - \alpha) f(y|s_j, \theta) dy \\
&= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \beta s_j^2 \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) dy \\
&\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j s_j \int_{Y_{L_i}}^{Y_{R_i}} (y - \alpha) f(y|s_j, \theta) dy = \beta \sum_{i=1}^n c_i \\
&\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \left(\sum_{j=1}^m \alpha_{ij} w_j s_j \int_{Y_{L_i}}^{Y_{R_i}} y f(y|s_j, \theta) dy - \alpha \sum_{j=1}^m \alpha_{ij} w_j s_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) dy \right) \\
&= \beta \sum_{i=1}^n c_i \\
&\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j s_j \int_{Y_{L_i}}^{Y_{R_i}} y f(y|s_j, \theta) dy - \alpha \sum_{i=1}^n b_i = \beta \sum_{i=1}^n c_i
\end{aligned}$$

$$\begin{aligned} &\Leftrightarrow \sum_{i=1}^n d_i - \alpha \sum_{i=1}^n b_i = \beta \sum_{i=1}^n c_i \\ &\Leftrightarrow \bar{d} - \alpha \bar{b} = \beta \bar{c} \Leftrightarrow \beta = \frac{\bar{d} - \alpha \bar{b}}{\bar{c}}, \end{aligned}$$

and replacing α by its estimate $\hat{\alpha}$ from (E1) results that

$$\hat{\beta} = \frac{\bar{d} - \hat{\alpha} \bar{b}}{\bar{c} - \hat{\alpha} \bar{b}^2}.$$

Finally, from equation (E3) one obtains

$$\begin{aligned} (E3) &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) \left(\frac{-1}{2\sigma^2} + \frac{(y - \alpha - \beta s_j)^2}{2\sigma^4} \right) dy = 0 \\ &\Leftrightarrow \frac{1}{\sigma^4} \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) (y - \alpha - \beta s_j)^2 dy \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) dy \\ &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) (y - \alpha - \beta s_j)^2 dy = n\sigma^2. \end{aligned}$$

Noting that $(y - \alpha - \beta s_j)^2 = (y - \alpha)^2 + \beta^2 s_j^2 - 2\beta s_j(y - \alpha)$, this is equal to

$$\begin{aligned} &\sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} (y - \alpha)^2 f(y|s_j, \theta) dy = n\sigma^2 - \beta^2 \sum_{i=1}^n c_i + 2\beta^2 \sum_{i=1}^n c_i \\ &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} (y - \alpha)^2 f(y|s_j, \theta) dy = n\sigma^2 + \beta^2 \sum_{i=1}^n c_i \\ &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} (y^2 - 2\alpha y + \alpha^2) f(y|s_j, \theta) dy = n\sigma^2 + \beta^2 \sum_{i=1}^n c_i \\ &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \left(\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} y^2 f(y|s_j, \theta) dy - 2\alpha \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} y f(y|s_j, \theta) dy \right) \end{aligned}$$

$$\begin{aligned}
& +\alpha^2 \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y|s_j, \theta) dy = n\sigma^2 + \beta^2 \sum_{i=1}^n c_i \\
\Leftrightarrow & \sum_{i=1}^n e_i - 2\alpha \sum_{i=1}^n a_i + n\alpha^2 = n\sigma^2 + \beta^2 \sum_{i=1}^n c_i \\
\Leftrightarrow & \sum_{i=1}^n e_i - 2\alpha \sum_{i=1}^n a_i + n\alpha^2 - \beta^2 \sum_{i=1}^n c_i = n\sigma^2 \\
\Rightarrow & \hat{\sigma}^2 = \bar{e} - 2\hat{\alpha}\bar{a} + \hat{\alpha}^2 - \hat{\beta}^2\bar{c}.
\end{aligned}$$

Appendix B

Maple program for the calculation of approximate confidence intervals

```
> with(LinearAlgebra):
```

Specifying the number of observations n and the number of examinations m

```
> n:=2; m:=6;
```

Reading the data

```
> data:=matrix(9,6,readdata('A:\\data.txt',9));
```

Assigning the variables needed in the loglikelihood

```
> ID:=matrix(n,m);  
> for i from 1 to n do for j from 1 to m do  
  ID[i,j]:=data[i,j] end do end do;  
> for j from 1 to m do w[j]:=data[3,j] end do;  
> for j from 1 to m do s[j]:=data[4,j] end do;  
> for i from 1 to n do yl[i]:=data[5,i] end do;  
> for i from 1 to n do yr[i]:=data[6,i] end do;  
> alphahat:=data[7,1];
```

```
> betahat:=data[8,1];
> sigma2hat:=data[9,1];
```

Definition of the log-likelihood

```
> i:='i'; j:='j';
> for j from 1 to m do
  f[j]:=(1/(sqrt(2*Pi*sigma^2)))*exp(-((y-alpha-beta*s[j])^2)/
    (2*sigma^2)) end do;
> loglike:=
  sum('log(sum('ID[i,j]*w[j]*int(f[j],y=y1[i]..yr[i])),
    'j'=1..6))','i'=1..n);
```

Calculation of the score function of loglike

```
> i:='i'; j:='j';
> der11:=diff(loglike,alpha);
> der12:=diff(loglike,beta);
> der13a:=algsubs(sigma^2=V,loglike);
> der13b:=subs(sigma=sqrt(V),der13a);
> der13c:=diff(der13b,V);
> der13:=subs(V=sigma^2,der13c);
```

Calculation of the second derivatives of loglike

```
> der111:=diff(der11,alpha);
> der112:=diff(der11,beta);
> der113a:=algsubs(sigma^2=V,der11);
> der113b:=subs(sigma=sqrt(V),der113a);
> der113c:=diff(der113b,V);
> der113:=subs(V=sigma^2,der113c);
>
> der122:=diff(der12,beta);
> der123a:=algsubs(sigma^2=V,der12);
> der123b:=subs(sigma=sqrt(V),der123a);
> der123c:=diff(der123b,V);
> der123:=subs(V=sigma^2,der123c);
>
> der133a:=algsubs(sigma^2=V,der13);
```

```
> der133b:=subs(sigma=sqrt(V),der133a);
> der133c:=diff(der133b,V);
> der133:=subs(V=sigma^2,der133c);
```

Construction of the Hessian matrix

```
> matt:=Matrix(1..3,1..3,[[der111,der112,der113],
      [der112,der122,der123],[der113,der123,der133]]);
>
```

Calculating the observed information matrix

```
> alpha:=alphahat;beta:=betahat;sigma:=sqrt(sigma2hat);
> evalf(matt);
> fish:=evalf(-1*matt);
```

Inverting the observed information matrix which is an estimate for the variance of $\hat{\alpha}$, $\hat{\beta}$ and $\hat{\sigma}^2$

```
> variance:=MatrixInverse(fish);
```

Constructing the confidence intervals for the regression parameters

```
> alpha:='alpha';beta:='beta';sigma:='sigma';
> CI(alpha):=[alphahat-1.96*sqrt(variance[1,1])/sqrt(n),
      alphahat+1.96*sqrt(variance[1,1])/sqrt(n)];
> CI(beta):=[betahat-1.96*sqrt(variance[2,2])/sqrt(n),
      betahat+1.96*sqrt(variance[2,2])/sqrt(n)];
> CI(sigma):=[sigma2hat-1.96*sqrt(variance[3,3])/sqrt(n),
      sigma2hat+1.96*sqrt(variance[3,3])/sqrt(n)];
```


Appendix C

Derivation of the MLE for the multiple regression setting

With the notations given in Appendix A, setting the partial derivations of the likelihood to zero and solving for the parameters, one yields the following solutions:

For the parameter α it holds that

$$\begin{aligned}
 \frac{\partial \log L^*}{\partial \alpha} &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \frac{y - \alpha - \vec{\beta}'_1 \vec{x}_i - \beta_2 s_j}{\sigma^2} f(y | (\vec{x}_i, s_j); \theta) dy \stackrel{!}{=} 0 \\
 &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} y f(y | (\vec{x}_i, s_j); \theta) dy \\
 &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} (\alpha + \vec{\beta}'_1 \vec{x}_i + \beta_2 s_j) f(y | (\vec{x}_i, s_j); \theta) dy \\
 &\Leftrightarrow \sum_{i=1}^n a_i = n\alpha + \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \vec{\beta}'_1 \vec{x}_i f(y | (\vec{x}_i, s_j); \theta) dy + \beta_2 \sum_{i=1}^n b_i \\
 &\Leftrightarrow \sum_{i=1}^n a_i - \sum_{l=1}^p \beta_{1l} \sum_{i=1}^n x_{li} - \beta_2 \sum_{i=1}^n b_i = n\alpha.
 \end{aligned}$$

For the parameter $\vec{\beta}_1$ it holds that

$$\begin{aligned} \frac{\partial \log L^*}{\partial \vec{\beta}_1} &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \frac{\vec{x}_i(y - \alpha - \vec{\beta}_1' \vec{x}_i - \beta_2 s_j)}{\sigma^2} f(y | (\vec{x}_i, s_j); \theta) dy \\ &\stackrel{!}{=} 0 \\ &\Leftrightarrow \sum_{i=1}^n x_{ki} a_i = \alpha \sum_{i=1}^n x_{ki} + \sum_{l=1}^p \beta_{1l} \sum_{i=1}^n x_{li} x_{ki} + \beta_2 \sum_{i=1}^n x_{ki} b_i, \\ &\Leftrightarrow \sum_{i=1}^n x_{ki} a_i - \alpha \sum_{i=1}^n x_{ki} - \beta_2 \sum_{i=1}^n x_{ki} b_i = \sum_{l=1}^p \beta_{1l} \sum_{i=1}^n x_{li} x_{ki}, \end{aligned}$$

for $k = 1, \dots, p$.

For the parameter β_2 it holds that

$$\begin{aligned} \frac{\partial \log L^*}{\partial \beta_2} &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \frac{s_j(y - \alpha - \vec{\beta}_1' \vec{x}_i - \beta_2 s_j)}{\sigma^2} f(y | (\vec{x}_i, s_j); \theta) dy \\ &\stackrel{!}{=} 0 \\ &\Leftrightarrow \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \beta_2 s_j^2 f(y | (\vec{x}_i, s_j); \theta) dy \\ &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} s_j(y - \alpha - \vec{\beta}_1' \vec{x}_i) f(y | (\vec{x}_i, s_j); \theta) dy \\ &\Leftrightarrow \beta_2 \sum_{i=1}^n c_i = \sum_{i=1}^n d_i - \alpha \sum_{i=1}^n b_i - \sum_{l=1}^p \beta_{1l} \sum_{i=1}^n x_{li} b_i. \end{aligned}$$

For the parameter σ^2 it holds that

$$\begin{aligned} \frac{\partial \log L^*}{\partial \sigma^2} &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \left(-\frac{1}{2\sigma^2} + \frac{2(y - \alpha - \vec{\beta}_1' \vec{x}_i - \beta_2 s_j)^2}{4\sigma^4} \right) \\ &\quad f(y | (\vec{x}_i, s_j); \theta) dy \stackrel{!}{=} 0 \\ &\Leftrightarrow \frac{1}{2\sigma^2} \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y | (\vec{x}_i, s_j); \theta) dy \end{aligned}$$

$$= \frac{1}{2\sigma^4} \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} (y - \alpha - \vec{\beta}_1^t \vec{x}_i - \beta_2 s_j)^2 f(y | (\vec{x}_i, s_j); \theta) dy. \quad (*)$$

Noting that $(y - \alpha - \vec{\beta}_1^t \vec{x}_i - \beta_2 s_j)^2$ is equivalent to $(y - \alpha - \vec{\beta}_1^t \vec{x}_i)^2 + \beta_2^2 s_j^2 - 2\beta_2 s_j (y - \alpha - \vec{\beta}_1^t \vec{x}_i)$, it holds that

$$\begin{aligned} (*) &\Leftrightarrow n\sigma^2 - \beta^2 \sum_{i=1}^n c_i + 2\beta^2 \sum_{i=1}^n c_i \\ &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} (y - \alpha - \vec{\beta}_1^t \vec{x}_i)^2 f(y | (\vec{x}_i, s_j); \theta) dy \\ &\Leftrightarrow n\sigma^2 + \beta^2 \sum_{i=1}^n c_i = \sum_{i=1}^n e_i - 2\alpha \sum_{i=1}^n a_i - 2 \sum_{l=1}^p \beta_{1l} \sum_{i=1}^n x_{li} a_i \\ &\quad - 2\alpha \sum_{l=1}^p \beta_{1l} \sum_{i=1}^n x_{li} + n\alpha^2 + \sum_{l=1}^p \beta_{1l}^2 \sum_{i=1}^n x_{li}^2. \end{aligned}$$

Appendix D

Derivation of the MLE when the errors come from the exponential family

Setting $\log L^{**}$ to zero and solving for the parameters yields the following maximum likelihood estimates:

With respect to α one gets

$$\begin{aligned} \frac{\partial \log L^{**}}{\partial \alpha} &= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \frac{\partial h(y_i - \alpha - \beta s_j)}{\partial \alpha} c(\eta) \exp[Q(\eta)t(y_i - \alpha - \beta s_j)]}{C_i(\theta)} \right. \\ &\quad \left. + \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) \left[Q(\eta) \frac{\partial t(y_i - \alpha - \beta s_j)}{\partial \alpha} \right] dy}{C_i(\theta)} \right) \\ &= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} -\frac{h'(\varepsilon_i)}{h(\varepsilon_i)} f(y_i | s_j; \theta) dy}{C_i(\theta)} \right. \\ &\quad \left. - \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) [Q(\eta)t'(\varepsilon_i)] dy}{C_i(\theta)} \right) = (F1). \end{aligned}$$

With respect to β one gets

$$\begin{aligned}
\frac{\partial \log L^{**}}{\partial \beta} &= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \frac{\partial h(y_i - \alpha - \beta s_j)}{\partial \beta} c(\eta) \exp[Q(\eta)t(y_i - \alpha - \beta s_j)]}{C_i(\theta)} \right. \\
&\quad \left. + \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) \left[Q(\eta) \frac{\partial t(y_i - \alpha - \beta s_j)}{\partial \beta} \right] dy}{C_i(\theta)} \right) \\
&= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} -\frac{s_j h'(\varepsilon_i)}{h(\varepsilon_i)} f(y_i | s_j; \theta) dy}{C_i(\theta)} \right. \\
&\quad \left. - \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) [Q(\eta) s_j t'(\varepsilon_i)] dy}{C_i(\theta)} \right) = (F2).
\end{aligned}$$

With respect to η one gets

$$\begin{aligned}
\frac{\partial \log L^{**}}{\partial \eta} &= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} h(y_i - \alpha - \beta s_j) c'(\eta) \exp[Q(\eta)t(y_i - \alpha - \beta s_j)]}{C_i(\theta)} \right. \\
&\quad \left. + \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) [Q'(\eta)t(y_i - \alpha - \beta s_j)] dy}{C_i(\theta)} \right) \\
&= \sum_{i=1}^n \left(\frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} -\frac{c'(\eta)}{c(\eta)} f(y_i | s_j; \theta) dy}{C_i(\theta)} \right. \\
&\quad \left. + \frac{\sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) [Q'(\eta)t(\varepsilon_i)] dy}{C_i(\theta)} \right) = (F3).
\end{aligned}$$

Setting these equations to zero and solving for the parameters one obtains the maximum likelihood equations

$$\begin{aligned}\sum_{i=1}^n h_i &= -Q(\eta) \sum_{i=1}^n t'_i, \\ \sum_{i=1}^n z h_i &= -Q(\eta) \sum_{i=1}^n z t'_i, \\ \sum_{i=1}^n n \frac{c'(\eta)}{c(\eta)} &= -Q'(\eta) \sum_{i=1}^n t_i,\end{aligned}$$

where

$$\begin{aligned}t_i &= E(t(\varepsilon) | [y_{L_i}, y_{R_i}]; [x_{L_i}, x_{R_i}]), \\ t'_i &= E(t'(\varepsilon) | [y_{L_i}, y_{R_i}]; [x_{L_i}, x_{R_i}]), \\ z t'_i &= E(Z t'(\varepsilon) | [y_{L_i}, y_{R_i}]; [x_{L_i}, x_{R_i}]), \\ h_i &= E\left(\frac{h'(\varepsilon)}{h(\varepsilon)} \middle| [y_{L_i}, y_{R_i}]; [x_{L_i}, x_{R_i}]\right), \\ z h_i &= E\left(Z \frac{h'(\varepsilon)}{h(\varepsilon)} \middle| [y_{L_i}, y_{R_i}]; [x_{L_i}, x_{R_i}]\right).\end{aligned}$$

Appendix E

Derivation of the MLE when the errors come from the Weibull distribution

Setting l^{***} to zero and solving for the parameters yields the following solutions:

For the parameters $\hat{\alpha}$ it holds that

$$\begin{aligned} \frac{\partial l^{***}}{\partial \alpha} &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \beta s_j^{\beta-1} \exp(-\alpha s_j^\beta) - \alpha \beta s_j^{\beta-1} \exp(-\alpha s_j^\beta) s_j^\beta dy \\ &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) \left(\frac{1}{\alpha} - s_j^\beta \right) dy \stackrel{!}{=} 0 \\ &\Leftrightarrow \frac{n}{\alpha} = \sum_{i=1}^n f_i \quad \Rightarrow \quad \hat{\alpha} = n / \sum_{i=1}^n f_i. \end{aligned}$$

For the parameters $\hat{\beta}$ it holds that

$$\begin{aligned} \frac{\partial l^{***}}{\partial \beta} &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} \left(\alpha s_j^{\beta-1} + \alpha \beta (\beta - 1) s_j^{\beta-2} \right) \exp(-\alpha s_j^\beta) \\ &\quad + \alpha \beta s_j^{\beta-1} \exp(-\alpha s_j^\beta) \left(-\alpha \beta s_j^{\beta-1} \right) dy \\ &= \sum_{i=1}^n \frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} f(y_i | s_j; \theta) \left(\frac{1}{\beta} + (\beta - 1) s_j - \alpha \beta s_j^{\beta-1} \right) dy \stackrel{!}{=} 0 \end{aligned}$$

$$\Leftrightarrow \frac{n}{\beta} + (\beta - 1) \sum_{i=1}^n b_i = \alpha\beta \sum_{i=1}^n g_i$$

$$\Leftrightarrow \frac{n}{\beta} + \beta \sum_{i=1}^n b_i - \alpha\beta \sum_{i=1}^n g_i = \sum_{i=1}^n b_i.$$

In these expressions, the following conditional expected values are used:

$$\frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} s_j^\beta f(y_i | s_j; \theta) dy = E(Z^\beta | [y_{L_i}, y_{R_i}], [z_{L_i}, z_{R_i}]) = f_i$$

and

$$\frac{1}{C_i(\theta)} \sum_{j=1}^m \alpha_{ij} w_j \int_{Y_{L_i}}^{Y_{R_i}} s_j^{\beta-1} f(y_i | s_j; \theta) dy = E(Z^{\beta-1} | [y_{L_i}, y_{R_i}], [z_{L_i}, z_{R_i}]) = g_i.$$