



## Maximum Likelihood for Interval Censored Data: Consistency and Computation

Robert Gentleman; Charles J. Geyer

*Biometrika*, Vol. 81, No. 3. (Aug., 1994), pp. 618-623.

Stable URL:

<http://links.jstor.org/sici?sici=0006-3444%28199408%2981%3A3%3C618%3AMLFICD%3E2.0.CO%3B2-T>

*Biometrika* is currently published by Biometrika Trust.

---

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/bio.html>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

---

The JSTOR Archive is a trusted digital repository providing for long-term preservation and access to leading academic journals and scholarly literature from around the world. The Archive is supported by libraries, scholarly societies, publishers, and foundations. It is an initiative of JSTOR, a not-for-profit organization with a mission to help the scholarly community take advantage of advances in technology. For more information regarding JSTOR, please contact [support@jstor.org](mailto:support@jstor.org).

# Maximum likelihood for interval censored data: Consistency and computation

BY ROBERT GENTLEMAN

*Department of Statistics, University of Auckland, Private Bag 92019, Auckland,  
New Zealand*

AND CHARLES J. GEYER

*School of Statistics, University of Minnesota, 206 Church Street S.E., Minneapolis,  
Minnesota, 55455, U.S.A.*

## SUMMARY

Standard convex optimization techniques are applied to the analysis of interval censored data. These methods provide easily verifiable conditions for the self-consistent estimator proposed by Turnbull (1976) to be a maximum likelihood estimator and for checking whether the maximum likelihood estimate is unique. A sufficient condition is given for the almost sure convergence of the maximum likelihood estimator to the true underlying distribution function.

*Some key words:* Convergence; Self-consistency algorithm; Uniqueness.

## 1. INTRODUCTION

Three data-collection schemes have been referred to as interval censored. Following Peto (1973) we use this term only to refer to the following situation. For each individual  $i$  there is a sequence of inspection times  $t_{i,1}, t_{i,2}, \dots$ . The exact failure time  $x_i$  of the individual is not observed. All that is known is which inspection times immediately preceded and followed the failure, that is the  $j$  such that  $t_{i,j-1} < x_i < t_{i,j}$ . Such data have been considered by Peto (1973), Turnbull (1976) and Finkelstein (1986), among others. A generalization of this situation has been considered by De Gruttola & Lagakos (1989), but they refer to it as doubly-censored data. Interval censored data, as we have defined it, differs substantially from grouped data (Heitjan, 1989) and the doubly-censored data of Chang & Yang (1987).

## 2. ESTIMATION

### 2.1. *The likelihood*

Suppose that survival times,  $X$ , arise from a distribution  $F_0$ , that each individual has a possibly infinite sequence of inspection times arising from some stochastic process  $Q$ , and that the inspection times and failure times are independent. This ensures that the censoring is noninformative. Also suppose that no time point occurs with positive probability under the inspection time process. This assumption is made to ensure that failures cannot coincide with inspection times. The observed data consist of the last inspection time prior to failure and the first inspection time after failure for each individual; i.e. the data are  $\{I_i\}_{i=1}^n$ , where  $I_i = (L_i, R_i)$  is the open interval known to contain the unobserved failure time.

These assumptions ensure that the probabilities of inspection times do not involve any of the parameters of interest and hence we may consider the likelihood conditional upon the

observed intervals,

$$L = \prod_{i=1}^n \{F_0(R_i-) - F_0(L_i)\}.$$

Let  $\{s_j\}_{j=0}^m$  denote the unique ordered elements of  $\{0, \{L_i\}_{i=1}^n, \{R_i\}_{i=1}^n\}$ . Then, as noted by Peto (1973), Turnbull (1976) and Finkelstein (1986), the likelihood depends on  $F_0$  only through the values  $\{F_0(s_j)\}_{j=1}^m$  and not on how  $F$  changes between  $s_j$ . Let  $\alpha_{ij}$  be the indicator of the event  $(s_{j-1}, s_j) \subseteq I_i$  and  $p_j = F_0(s_j-) - F_0(s_{j-1})$ ; then the likelihood can be written

$$L = \prod_{i=1}^n \left( \sum_{j=1}^m \alpha_{ij} p_j \right)$$

and the log likelihood as

$$l(p) = \sum_{i=1}^n \log \left( \sum_{j=1}^m \alpha_{ij} p_j \right).$$

Also let

$$d_k = \frac{\partial l}{\partial p_k} = \sum_{i=1}^n \frac{\alpha_{ik}}{\eta_i},$$

where  $\eta_i = \sum \alpha_{ij} p_j$ , with the summation over  $j = 1, \dots, m$ . The terms  $\eta_i$  correspond to the sum of probabilities associated with the  $i$ th individual and hence  $d_k$  is the sum of  $1/\eta_i$  for all individuals whose intervals,  $I_i$ , intersect the interval  $(s_{k-1}, s_k)$ .

### 2.2. The Kuhn–Tucker conditions

To find the maximum likelihood estimate of the vector  $p$  we maximize  $l(p)$  with respect to  $p$  subject to the constraints

$$1 - \sum_{j=1}^m p_j = 0, \tag{1}$$

$$p_j \geq 0 \quad (j = 1, \dots, m). \tag{2}$$

For a concave programming problem with linear constraints, the Kuhn–Tucker conditions are necessary and sufficient for optimality (Rockafellar, 1970, Theorem 28.1, Corollary 28.2.2). A point  $\hat{p}$  is a maximum likelihood estimate if and only if there exist Lagrange multipliers  $\mu_j$  ( $j = 0, \dots, m$ ) such that the Kuhn–Tucker conditions (1)–(5) hold, with

$$\mu_j p_j = 0 \quad (j = 1, \dots, m), \tag{3}$$

$$\mu_j \geq 0 \quad (j = 1, \dots, m), \tag{4}$$

$$\frac{\partial}{\partial p_j} \left\{ l(p) + \sum_{j=1}^m p_j (\mu_j - \mu_0) \right\} = d_j + \mu_j - \mu_0 = 0 \quad (j = 1, \dots, m). \tag{5}$$

Multiplying (5) by  $p_j$  and summing yields

$$\mu_0 = \sum_j d_j p_j = \sum_{i,j} \frac{\alpha_{ij} p_j}{\eta_i} = \sum_i \frac{\eta_i}{\eta_i} = n$$

since  $\mu_j p_j = 0$  by (3). If  $p_j > 0$  then (3) implies that  $\mu_j = 0$ , and (5) then implies that  $d_j = \mu_0 = n$ . Conversely, if  $p_j = 0$  then (5) implies that  $\mu_j \geq 0$  so  $d_j = \mu_0 - \mu_j$  implies  $d_j \leq n$ . At a solution all of the  $\eta_i$  are strictly positive, since otherwise the  $d_k$  would not be finite.

For any  $p$  that satisfies the constraints (1) and (2), if  $p_j = 0$  set  $\mu_j = n - d_j$ , and if  $p_j > 0$  set  $\mu_j = 0$ . Then condition (3) is always satisfied. We call the  $\mu_j$  Lagrange multipliers whether or not

they satisfy (4). The left-hand side of (5),  $d_j + \mu_j - \mu_0$ , is called the reduced gradient, because it is the gradient with respect to the free variables. The Kuhn–Tucker conditions are satisfied if the Lagrange multipliers are nonnegative and the reduced gradient is zero.

Peto (1973) and Turnbull (1976) point out that  $p_j$  can be nonzero only if  $s_{j-1}$  is a left endpoint  $L_i$  for some individual  $i$  and  $s_j$  is a right endpoint  $R_k$  for some possibly different individual  $k$ . However, some of the  $p_j$  satisfying this criterion may also be zero.

### 2.3. Uniqueness of the maximum likelihood estimate

The maximum likelihood estimate need not be unique. Turnbull (1976) gives the example with  $\alpha_{ij} = \alpha_{ik}$  for all  $i$ . The maximum likelihood estimate will be unique if the log likelihood is strictly concave, that is the Hessian  $H$  is strictly negative definite. Let  $A$  denote the  $n \times m$  matrix with elements  $\alpha_{ij}$ , then  $H = A'DA$ , where  $D$  is the diagonal matrix with elements  $-1/\eta_i^2$ . Hence,  $H$  will be of full rank and the maximum likelihood estimate will be unique if  $\text{rank}(A) = m$ .

There may be situations in which the likelihood is concave, but not strictly concave, and the maximum likelihood estimate is unique nevertheless. Theorem 9.3.2 of Fletcher (1987), specialized to our problem, states the following. Let  $\hat{p}$  be a solution to the Kuhn–Tucker equations with suitable Lagrange multipliers  $\mu$ . Define

$$W = \{w \in \mathbb{R}^m : w_j = 0 \text{ if } \mu_j > 0; w_j \geq 0, \text{ if } \hat{p}_j = 0; \sum_j w_j = 0\}.$$

Then the maximum likelihood estimate  $\hat{p}$  is unique if, whenever  $w \in W$  and  $w \neq 0$ ,

$$w'Hw < 0. \quad (6)$$

We can get a condition much easier to verify if we drop some of the constraints and verify the condition (6) with the set  $W$  replaced by

$$W' = \{w \in \mathbb{R}^m : w_j = 0 \text{ if } \mu_j > 0\}.$$

Since we check a larger set, this condition implies the other and is sufficient.

This can be further simplified by letting  $A = (A_1 A_2)$  be a partition of  $A$ , where  $A_1$  consists of those columns  $j$  such that  $\mu_j > 0$ . Also partition vectors  $w = (w_1 w_2)$  in the same way. The sufficient condition involving  $W'$  can then be stated as

$$w_2'A_2'DA_2w_2 \neq 0, \quad w_2 \neq 0,$$

where one direction of the inequality comes from (6) and the other from concavity. Since  $D$  is negative definite, this occurs if and only if  $A_2w_2 \neq 0$ , which proves the following.

**THEOREM 1.** *A sufficient condition for uniqueness of the maximum likelihood estimate is that the matrix  $A_2$  consisting of the columns of  $A$  corresponding to  $j$  such that  $\mu_j = 0$  has rank equal to its number of columns.*

## 3. CONSISTENCY

Maximum likelihood estimation for interval censored data is strongly consistent. The maximum likelihood estimator converges almost surely to the truth in a topology to be described presently. For simplicity we assume that  $F_2(0) = 0$ , and that all of the inspection times are greater than zero. We also assume that with probability one there are only a finite number of inspection times in any finite interval so that each realization of the inspection time process can be written  $t = (t_0, t_1, \dots, t_{m(t)})$ , where

$$0 = t_0 < t_1 < \dots < t_{m(t)} = +\infty$$

and  $m(t)$  is either finite or  $\infty$ . The assumption that all times are positive serves merely to avoid doubly infinite sequences.

The log-likelihood for our problem is then

$$l(F) = \sum_{i=1}^n \sum_{j=1}^{m(t_i)} 1(t_{i,j} > x_i > t_{i,j-1}) \log \{F(t_{i,j}-) - F(t_{i,j-1})\}.$$

Note that in the inner summation exactly one of the indicators is nonzero so this summation is simply another notation for  $\log \{F(R-) - F(L)\}$ . A proof of consistency requires a suitable compactification of the parameter space, which we take to be the set  $\bar{\Theta}$  of all subdistribution functions with the topology of vague convergence which is compact by Helley's selection theorem. The expectation of the log likelihood ratio,  $\lambda(F) = E\{l(F) - l(F_0)\}$ , is an upper semicontinuous, nonnegative concave function by Fatou's lemma, Jensen's inequality, and the assumption that no inspection time occurs with positive probability. So the set  $C = \{F: \lambda(F) = 0\}$  is a closed subset of  $\bar{\Theta}$ . The distribution functions in  $C$  cannot be distinguished by maximum likelihood. Hence, following Redner (1981), we identify all of the points in  $C$  with  $F_0$ .

Then we have the following theorem.

**THEOREM 2.** *Under the assumptions stated above, the maximum likelihood estimate  $\{\hat{F}_n\}$  converges strongly to the equivalence class  $C$  of the true distribution in the topology of vague convergence.*

This says that the sequence  $\{\hat{F}_n\}$  is eventually in every neighbourhood of  $C$ . The proof of the theorem follows Wang (1975), and is available from the authors.

The equivalence class  $C$  is the set of all distribution functions  $F$  such that  $F(t_j) = F_0(t_j)$  for  $j = 0, \dots, m(t)$  and almost all inspection time sequences  $t = (t_0, t_1, \dots, t_{m(t)})$ . If the inspection time process densely samples  $[0, \infty)$ , the equivalence class  $C$  will contain only  $F_0$ .

#### 4. COMPUTATION

The method proposed by Turnbull (1976), a version of the EM algorithm, is easy to implement but is known to have slow convergence. Alternative methods are the constrained Newton-Raphson method of Peto (1973) and the similar active set methods of optimization theory (Fletcher, 1987, § 11.2). The latter are more difficult to implement but have quadratic convergence.

Another problem with Turnbull's algorithm is that there can be self-consistency points other than the maximum likelihood estimate. These are not stationary points of the log-likelihood. They are maxima relative to faces of the parameter space, but moving away from such points into the interior increases the likelihood. An example of this is where  $F(t)$  puts mass only on the interval  $(0, 3]$ . Suppose that the data are the intervals  $(0, 1]$ ,  $(1, 3]$ ,  $(1, 3]$ ,  $(0, 2]$ ,  $(0, 2]$ ,  $(2, 3]$ . Then it can be verified that  $p(0, 1] = \frac{1}{2}$ ,  $p(1, 2] = 0$ ,  $p(2, 3] = \frac{1}{2}$  is a self-consistent estimator while  $p(0, 1] = \frac{1}{3}$ ,  $p(1, 2] = \frac{1}{3}$ ,  $p(2, 3] = \frac{1}{3}$  is the maximum likelihood estimate. An examination of the Kuhn-Tucker conditions at  $(\frac{1}{2}, 0, \frac{1}{2})$  shows that they are violated at this point.

The occurrence of self-consistency points other than the maximum likelihood estimate is troubling for two reasons. First, continuity of the EM steps implies that the algorithm makes arbitrarily small steps near a self-consistency point so it is not possible to test for convergence by examining either the sequence of iterates or the likelihood along the sequence. Secondly, as will be illustrated in the next section, it is a reasonable procedure to restart the EM algorithm with very small parameter values set to zero to 'polish' the parameter values. This will produce incorrect results if the zeros are incorrectly determined, since the EM iteration never changes the zeros.

Both of these problems can be cured by the simple expedient of examining the Kuhn-Tucker conditions. If they are used as the convergence test, convergence to the maximum likelihood estimate is guaranteed. The computational effort required to check the Kuhn-Tucker conditions is minimal. All of the necessary quantities are calculated during the self-consistency iteration. Interestingly, Turnbull does derive a characterization of the maximum likelihood estimate equivalent to the Kuhn-Tucker conditions, but he does not recommend that it be used to test for convergence of the self-consistency algorithm.

## 5. EXAMPLE

The data in Table 1 come from Finkelstein & Wolfe (1985). It gives the interval in which cosmetic deterioration for early breast cancer patients treated with radiotherapy occurred in 46 individuals. The Kuhn–Tucker conditions indicate that there are only 14 intervals that need be considered; these intervals and the  $p_j$  associated with them are reported in the first three columns of Table 2. The matrix  $(\alpha_{ij})$  is of full rank; hence the maximum likelihood estimate is unique.

Table 1. *Intervals in which deterioration occurred*

(45, —]	(6, 10]	(0, 7]	(46, —]	(46, —]	(7, 16]	(17, —]	(7, 14]
(37, 44]	(0, 8]	(4, 11]	(15, —]	(11, 15]	(22, —]	(46, —]	(46, —]
(25, 37]	(46, —]	(26, 40]	(46, —]	(27, 34]	(36, 44]	(46, —]	(36, 48]
(37, —]	(40, —]	(17, 25]	(46, —]	(11, 18]	(38, —]	(5, 12]	(37, —]
(0, 5]	(18, —]	(24, —]	(36, —]	(5, 11]	(19, 35]	(17, 25]	(24, —]
(32, —]	(33, —]	(19, 26]	(37, —]	(34, —]	(36, —]		

Table 2. *Restricted set of intervals and the associated probabilities*

Left	Right	Probability	Probability	Reduced gradient	Lagrange multiplier	Probability	Reduced gradient	Lagrange multiplier
4	5	0.0463	0.0463	0.0002	0	0.0583	0.0000	0
6	7	0.0335	0.0334	0.0009	0	0	0	-6.097
7	8	0.0886	0.0886	0.0003	0	0.1128	0.0000	0
11	12	0.0708	0.0708	-0.0001	0	0.0680	0.0000	0
15	16	$4.19 \times 10^{-19}$	0	0	24.3	0	0	24.45
17	18	$2.65 \times 10^{-6}$	0	0	7.65	0	0	7.091
24	25	0.0927	0.0926	0.0000	0	0.0927	0.0000	0
25	26	$1.58 \times 10^{-7}$	0	0	9.36	0	0	9.42
33	34	0.0817	0.0818	0.0000	0	0.0817	0.0000	0
34	35	$4.88 \times 10^{-8}$	0	0	10.5	0	0	10.52
36	37	0.0007	0	0	2.87	0	0	2.87
38	40	0.1174	0.1206	0.0000	0	0.1185	0.0001	0
40	44	0.0031	0	0	2.79	0	0	2.786
46	48	0.4653	0.4658	0.0000	0	0.4654	0.0000	0

The first two columns give the intervals on which  $F(t)$  may have positive mass. The third column contains the maximum likelihood estimate of the masses; six were constrained to zero and the resulting estimate reported in the fourth column. The corresponding reduced gradient and Lagrange multipliers are given in the next two columns. The probability mass for the second interval was also constrained to zero and the resulting estimate is given in the seventh column, with the corresponding gradient and Lagrange multipliers in the last two columns.

Inspection of the probabilities indicates that several of them are very small and hence may be zero at the maximum likelihood estimate. They were set to zero and the EM algorithm applied to the resulting renormalized probability vector. The new candidate optimal point is reported in the fourth column of Table 2, the reduced gradient, defined in § 2.2, at this point is reported in the fifth column and the associated Lagrange multipliers in the sixth column. Notice that the Kuhn–Tucker conditions are approximately satisfied at the point reported in the fifth column of Table 2; hence we have found the maximum likelihood estimate at a point where six of the  $p_j$  are zero.

In this problem  $p_2$  may be set to zero without any of the  $\eta_i$  becoming zero. Doing this and applying the EM algorithm yields a self-consistent estimator that is not the maximum likelihood estimator as was described previously. However, an examination of the reduced gradient and the Lagrange multipliers at this point, the last three columns of Table 2, indicates that the Lagrange multiplier associated with  $p_2$  is negative and hence the Kuhn–Tucker conditions are violated at this point. It cannot be a maximum likelihood estimate.

## ACKNOWLEDGEMENT

R. Gentleman was supported by a grant from the Natural Sciences and Engineering Research Council of Canada. G. J. Geyer was supported by a postdoctoral fellowship from the National Science Foundation.

## REFERENCES

- CHANGE, M. N. & YANG, G. L. (1987). Strong consistency of a nonparametric estimator of the survival function with doubly censored data. *Ann. Statist.* **15**, 1536–47.
- DE GRUTTOLA, V. & LAKAGOS, S. W. (1989). Analysis of doubly-censored survival data with application to AIDS. *Biometrics* **45**, 1–11.
- FINKELSTEIN, D. M. (1986). A proportional hazards model for interval-censored failure time data. *Biometrics* **42**, 845–54.
- FINKELSTEIN, D. M. & WOLFE, R. A. (1985). A semiparametric model for regression analysis of interval-censored failure time data. *Biometrics* **41**, 933–45.
- FLETCHER, R. (1987). *Practical Methods of Optimization*, 2nd ed. New York: Wiley.
- HEITJAN, D. F. (1989). Inference from grouped continuous data: A review. *Statist. Sci.* **4**, 164–83.
- PETO, R. (1973). Experimental survival curves for interval-censored data. *Appl. Statist.* **22**, 86–91.
- REDNER, R. (1981). Note on the consistency of the maximum likelihood estimate for nonidentifiable distributions. *Ann. Statist.* **9**, 225–8.
- ROCKAFELLAR, R. T. (1970). *Convex Analysis*. Princeton University Press.
- TURNBULL, B. W. (1976). The empirical distribution function with arbitrarily grouped, censored and truncated data. *J. R. Statist. Soc. B* **38**, 290–5.
- WANG, J.-L. (1985). Strong consistency of approximate maximum likelihood estimators with applications in nonparametrics. *Ann. Statist.* **13**, 932–46.

[Received January 1991. Revised October 1993]