

Software for the algorithms used in “Nonparametric Estimation under Shape Constraints”

Piet Groeneboom

Abstract: We discuss software, written in C and used in [Groeneboom and Jongbloed \(2014\)](#), which we make available on this site, using Rcpp, see, e.g., [Eddelbuettel \(2013\)](#). Parallel to this development, we also create GUI (Graphical User Interface) applications for Mac OS 10.10 (Yosemite). The software can be used to reproduce the results and pictures in [Groeneboom and Jongbloed \(2014\)](#).

AMS 2000 subject classifications: Primary 62G05, 62N01; secondary 62-04.

Keywords and phrases: MLE, SMLE, iterative convex minorant algorithm, Rcpp.

1. Introduction

In [Groeneboom and Jongbloed \(2014\)](#) frequent use is made of results of C programs, which implement algorithms for computing estimators iteratively, or for computing test statistics using bootstrap procedures. Unfortunately, these computer programs were written using standard routines from Numerical Recipes in C (see [Press et al. \(1992\)](#)), and the authors of this book prohibit publication of code which even partly contains routines from their book.

For this reason the computer programs were never made public and the only way we can make them public now seems to remove all Numerical Recipes routines and use equivalent (preferably better) routines from other sources or just to cook up these routines ourselves. This is indeed what we intend to do, and a first attempt is the treatment of the MLE, SMLE and hazard for the Competing Risk model under Current Status censoring, which has received a lot of attention from researchers, and is discussed at different spots in the book. We will add other routines later.

We will try to link the C/C++ programs to R, using Rcpp. In this way the programs are immediately available for Mac, Windows and Linux/Unix platforms and we also can connect to the graphical tools available in R. Parallel to this development, we also create GUI (Graphical User Interface) applications for Mac OS 10.10 (Yosemite), built in Xcode, allowing the user to open input files via the standard open file dialogue and run bootstrap simulations for computing confidence intervals.

For using R, using Rcpp, Windows users may have to install Rtools:

<http://cran.r-project.org/bin/windows/Rtools/index.html>

Mac users may have to use Xcode, or at least the corresponding command line tools, see (with advices for Windows, Mac and Linux users):

<https://support.rstudio.com/hc/en-us/articles/200486498-Package-Development-Prerequisites>

Mac users who use the GUI application, given in `compriskGUI.tar.gz` on:

<http://dutiosc.twi.tudelft.nl/~pietg/software.html>

do not have to download extra tools.

2. The competing risk model under current status censoring

2.1. The model

The competing risk model under current status censoring is described on p. 121 of [Groeneboom and Jongbloed \(2014\)](#) and we take the formulation from there. Consider a situation where for a certain object there are several (say $K \in \mathbb{N}$) possible causes of failure and that at a random point T in time the object is inspected. It is observed whether or not this object broke down before time T or not (its ‘current status’). In case the object has broken down, it is also observed which of the K possible causes (competing risks) lead to the breakdown. Write X for the time of breakdown and $Y \in \{1, 2, \dots, K\}$ for the corresponding cause. Together with inspection time T , the indicator vector

$$\Delta = (\Delta_1, \dots, \Delta_K) \text{ with } \Delta_k = 1_{[X \leq T, Y=k]}$$

is observed. Note that if all indicators are zero, this means that $X > T$. If $X \leq T$ also the breakdown cause is observed, so then exactly one of the K indicators equals 1. Assuming (X, Y) to have joint distribution function F and T to be independent of (X, Y) , this model is known as the competing risk model with current status observations. Data of this type arise naturally in cross-sectional studies with several failure causes.

Note that, given T , the vector

$$(\Delta_1, \Delta_2, \dots, \Delta_K, \Delta_{K+1}) \text{ with } \Delta_{K+1} = 1 - \sum_{k=1}^K \Delta_k \quad (2.1)$$

has a multinomial distribution with parameters 1 and $(F_1(T), \dots, F_K(T), 1 - F_+(T))$, where

$$F_k(t) = P(X \leq t, Y = k), \quad t \geq 0, k = 1, 2, \dots, K,$$

is the sub-distribution function of X for risk level k and $F_+ = \sum_k F_k$ is the marginal distribution function of X . Denoting by e_k the k th unit vector in \mathbb{R}^K , by $\#$ counting measure on $D = \{e_k : k = 1, \dots, K+1\}$ and G the distribution of T , we can define the measure $\mu = G \times \#$ on $\mathbb{R} \times D$. With respect to this (dominating) measure, the density of a single observation (T, Δ) is given by

$$p_F(t, \delta) = \prod_{k=1}^K F_k(t)^{\delta_k} (1 - F_+(t))^{1-\delta_+}, \quad (2.2)$$

where $\delta_+ = \sum_{k=1}^K \delta_k$.

Now consider an independent sample of size n , distributed as $(T, \Delta_1, \dots, \Delta_K)$,

$$(T_i, \Delta^i) = (T_i, \Delta_1^i, \dots, \Delta_K^i), \quad i = 1, \dots, n,$$

where, for $1 \leq i \leq n$

$$\Delta^i = (\Delta_1^i, \dots, \Delta_K^i) \text{ with } \Delta_k^i = 1_{\{X_i \leq T_i, Y=k\}}, \quad k = 1, \dots, K.$$

Also define

$$\Delta_{K+1}^i = 1 - \sum_{k=1}^K \Delta_k^i = 1_{\{X_i > T_i\}}.$$

Using (2.2) and independence of the observations, the log likelihood (divided by n) is given by

$$\begin{aligned} \ell(F) &= \int \log p_F(t, \delta) d\mathbb{P}_n(t, \delta) \\ &= \int \left\{ \sum_{k=1}^K \delta_k \log F_k(t) + (1 - \delta_+) \log(1 - F_+(t)) \right\} d\mathbb{P}_n(t, \delta). \end{aligned} \quad (2.3)$$

where \mathbb{P}_n is the empirical distribution of (T_i, Δ^i) , $i = 1, \dots, n$. An MLE $\hat{F} = (\hat{F}_1, \dots, \hat{F}_K)$ can then be defined by the property

$$\ell(\hat{F}) = \max_{F \in \mathcal{F}_K} \ell(F) \quad (2.4)$$

where

$$\begin{aligned} \mathcal{F}_K &= \{F = (F_1, \dots, F_K) : F_1, \dots, F_K \text{ are sub-distribution functions,} \\ &\text{such that for all } x \geq 0 : \sum_{k=1}^K F_k(x) \leq 1\}. \end{aligned} \quad (2.5)$$

2.2. The Bangkok Metropolitan Administration cohort study

The Bangkok Metropolitan Administration injecting drug users cohort study ([Kitayaporn et al. \(1998\)](#) and [Vanichseni et al. \(2001\)](#)) was started in 1995 to assess (among other things) the feasibility of conducting a phase III HIV vaccine efficacy trial for injecting drug users in Bangkok. The data on a subset of 1365 injecting drug users who were below 35 years of age in this study were analyzed by [Maathuis and Hudgens \(2011\)](#) and [Li and Fine \(2013\)](#). In this group, 392 were HIV positive, with 114 infected with subtype B, 237 infected with subtype E, 5 infected by another mixed subtype and 36 infected with missing subtype. The subjects with other, mixed or missing subtypes were grouped in a single category.

In [Maathuis and Hudgens \(2011\)](#), the maximum likelihood estimator (MLE) for these data is computed and also a so-called naive estimator, based on analyzing one category such as the type B subjects, ignoring the data on the other types. In [Li and Fine \(2013\)](#) both the regular MLE and a smoothed version of the MLE (called the SMLE) are computed and theory developed in [Groeneboom, Jongbloed and Witte \(2010\)](#) is used for constructing confidence intervals. They also estimate the hazard and construct confidence intervals for the hazard, again using [Groeneboom, Jongbloed and Witte \(2010\)](#). The regular MLE cannot directly be used for this purpose because it corresponds to a discrete distribution, so that some kind of smoothing is needed to estimate the hazard and to construct the confidence intervals.

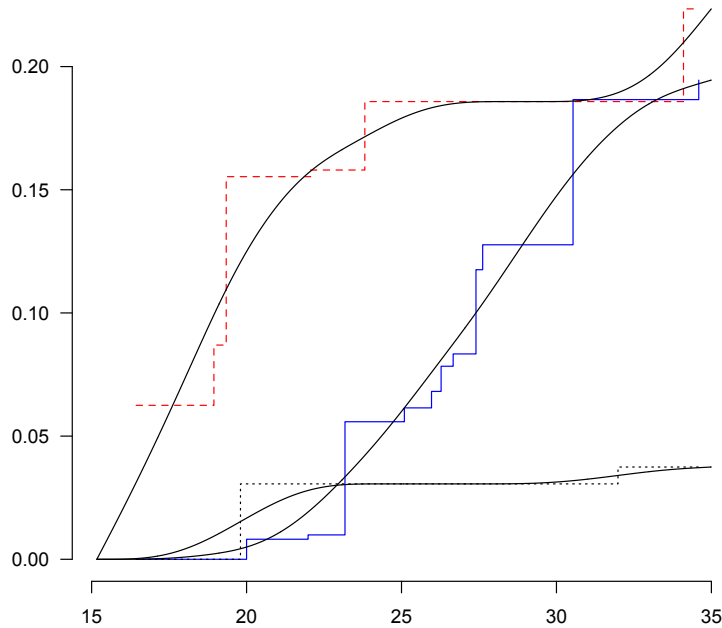


Fig 1: The MLE for the three categories in the Bangkok cohort study. The piecewise constant curves give the subdistribution functions, based on the MLE, for the different categories; dotted: type E, solid: type B, dashed; other types. The smooth solid curves give the corresponding estimates, based on the SMLE.

2.3. The iterative convex minorant algorithm

The iterative convex minorant algorithm is derived from Corollary 2.10 in [Groeneboom, Maathuis and Wellner, 2008](#), which is given below.

Lemma 2.1 (Corollary 2.10, [Groeneboom, Maathuis and Wellner, 2008](#)). *Let λ be given by*

$$\lambda = 1 - \int \frac{\delta_{K+1}}{1 - \hat{F}_+(u)} d\mathbb{P}_n(u, \delta) \geq 0. \quad (2.6)$$

Then $\hat{F} = (\hat{F}_1, \dots, \hat{F}_K)$ is an MLE if, for all $k = 1, \dots, K$ and each point of jump τ_{ki} of \hat{F}_k :

$$\int_{u \in [\tau_{ki}, s)} \left\{ \frac{\delta_k}{\hat{F}_k(u)} - \frac{\delta_{K+1}}{1 - \hat{F}_+(u)} \right\} d\mathbb{P}_n(u, \delta) \geq \lambda 1_{[\tau_{ki}, s)}(T_{(p)}), \quad s \in \mathbb{R}, \quad (2.7)$$

where equality holds if $s > \tau_{ki}$ is a point of increase of \hat{F}_k or if $s > T_{(p)}$, where $T_{(p)}$ is the largest of the strictly ordered order statistics.

To force (2.7) to hold, we can set up an iterative convex minorant algorithm. This algorithm is implemented in the C++ file `icm.cpp` and in particular in the routine `ICM_iteration`. A so-called cusum diagram is formed, with points $(0, 0)$ and points

$$\left(\sum_{i=1}^j w_{ki}, \sum_{i=1}^j (w_{ki} y_{ki} + v_{ki}) - \lambda 1_{\{j=n_k\}} \right), \quad j = 1, \dots, n_k, \quad (2.8)$$

where λ is as in (2.6), with \hat{F}_+ replaced by the F_+ at the present iteration step,

$$w_{ki} = \int_{u \in [T_i, T_j)} \left\{ \frac{\delta_k}{F_k(u)^2} + \frac{\delta_{K+1}}{\{1 - F_+(u)\}^2} \right\} d\mathbb{P}_n(u, \delta),$$

and T_i and T_j are successive points where $\delta_k = 1$,

$$v_{ki} = \int_{u \in [T_i, T_j)} \left\{ \frac{\delta_k}{F_k(u)} - \frac{\delta_{K+1}}{1 - F_+(u)} \right\} d\mathbb{P}_n(u, \delta),$$

and where y_{ki} is the value of the subdistribution function F_k at the point T_i at the present iteration; n_k is the number of points where $\delta_k = 1$. As explained in [Groeneboom and Jongbloed \(2014\)](#), the w_{ik} have an interpretation via the diagonal of the Hessian matrix of the maximization problem.

We next determine the greatest convex minorant of the cusum diagram (2.8) and use its left derivative for forming the new values of F_k . The Lagrange multiplier λ is computed after each iteration via (2.6), using the values of the F_k at that iteration. We only have to determine the values of F_k at the points where $\delta_k = 1$, since these are the only points where the subdistribution function can have mass. It can be seen that at a stationary point of these iterations the conditions (2.7) are satisfied and hence an MLE is found. To go from one iteration to the next we determine the step size by a golden section search algorithm.

In our experience, this algorithm is at present the fastest algorithm for finding the MLE. It can easily be extended to interval censoring with more observations, like case 2 interval censoring, or so-called mixed case interval censoring, but we concentrate for the moment on the current status model.

2.4. The SMLE

The SMLE (smoothed maximum likelihood estimator) is computed using (9.75) in [Groeneboom and Jongbloed \(2014\)](#), where we use boundary correction on the left and right. The SMLE is defined on p. 367 of [Groeneboom and Jongbloed \(2014\)](#) for the data of the Bangkok cohort study, and given by

$$\tilde{F}_{nk,h}(t) = \int \left\{ \mathbb{K} \left(\frac{t-x}{h} \right) + \mathbb{K} \left(\frac{t+x-2a}{h} \right) - \mathbb{K} \left(\frac{2b-t-x}{h} \right) \right\} d\hat{F}_{nk}(x),$$

where $a = 15$, $b = 35$, h is the bandwidth (taken to be $(b-a)n^{-1/5}$) and

$$\mathbb{K}(u) = \int_{-\infty}^u K(y) dy, \quad (2.9)$$

choosing for K for the triweight kernel, given by

$$K(u) = \frac{35}{32} (1 - u^2)^3 1_{[-1,1]}(u).$$

The SMLE is computed in the last routine `bdf` of `icm.cpp`.

As explained in [Groeneboom and Jongbloed \(2014\)](#), the SMLE has a faster asymptotic convergence rate ($n^{-2/5}$, the usual convergence rate of density estimation), if one is willing to assume smoothness, than the MLE (which has convergence rate $n^{-1/3}$ under the usual conditions). Also, the SMLE has asymptotically a normal limit distribution (under the appropriate conditions, given in [Groeneboom and Jongbloed \(2014\)](#)), in contrast with the MLE, for which the asymptotical distribution still has no analytic characterization.

2.5. Rcpp

The Rcpp part of the algorithm is implemented in the routine `ComputeMLE(DataFrame input)` in `icm.cpp`. The input to the routine is a data frame which comes from a file with two columns: the first column contains the observation times and the second the causes of failure $1, \dots, K$ or 0 if there is no failure at that observation time. The R script `MLEThai.R` essentially only exists of the lines:

```
library(Rcpp)
icm<-sourceCpp("icm.cpp")
A<-read.table("dataThai.txt")
output <- ComputeMLE(A)
```

and has as output a list of three things: the MLE, the SMLE (computed on a equidistant grid of 1000 points) and the value of the log likelihood. The first line activates the Rcpp library (one may have to load the package Rcpp first), the second line compiles the C++ file `icm.cpp`, the third line transforms the data into a form suitable for input to the function `ComputeMLE`, and the fourth line computes the MLE and SMLE from this function (and also produces the value of the log likelihood). The present data file `dataThai.txt` can of course be replaced by any data file of the same structure; the file should not have headers, and the number of observations and number of risks, which should have the labels $1, \dots, K$, is counted by the C++ program `icm.cpp` itself.

The program also performs a reduction to take ties into account. For example, the data file `dataThai.txt` contains 1365 observations for three risks. After reduction for ties, the number of remaining (unique) observations is 1211, and at each unique observation the frequencies of the occurrences of $\delta_k = 1$ are given, for $k = 1, \dots, K$, where $K = 3$ in this case.

One can also use an input file where the ties are already taken into account, and where the data consist of the frequencies of the observations for the different risks are given. In that one case the input file contains $K + 2$ columns: the first column contains the observation times, the next column the frequencies for the observations where no failure is observed and next one gets the columns with the frequencies for the observed failures of type k , $k = 1, \dots, K$. The script demonstrating this approach is given in `MLEThai2.R`, and the input file is `dataThai2.txt`. The C++ file for this approach is given by `icm2.cpp`.

For illustration purposes the drawing of Figure 1.7 on p. 11 of [Groeneboom and Jongbloed \(2014\)](#), where the MLE and SMLE are shown, is added to the scripts. This figure is also shown above in Figure 1.

2.6. Comparison with MLEcens

In first instance, the R package `MLEcens` ([Maathuis \(2013\)](#)) was developed by Marloes Maathuis for analyzing bivariate interval censored data. For this type of data, the non-uniqueness of the MLE is more of an issue than for the present case of the MLE for competing risk data under current status censoring. The method proceeds by first computing rectangles which are candidates for containing mass and next computing the masses of the MLE on these rectangles. The really time consuming step of the algorithm is the computation of the MLE, not the preliminary reduction step of the computation of the candidate rectangles.

This method can also be used for computing the MLE for the present model, in which case the rectangles are replaced by intervals. The output of `MLEcens` gives indeed the masses of the MLE on these intervals,

together with a list of the intervals. The computation of the MLE is based on an old C program, using support reduction, as discussed in [Groeneboom, Jongbloed and Wellner \(2008\)](#).

This algorithm is totally different from the iterative convex minorant algorithm. The support reduction algorithm successively builds up a distribution of mass from a few rectangles until further addition of new rectangles with mass will not increase the likelihood any more. This is done by two types of iteration: *inner iterations* which perform a least squares minimization, using weights given by the *outer iterations* which perform a change of the masses using Armijo's rule for determining a step size in a direction given by the inner iterations. Details of this procedure are given in [Groeneboom, Jongbloed and Wellner \(2008\)](#) for the so-called "Aspect experiment" in quantum statistics.

We compare the output of the two algorithms for the Bangkok data below. The iterative convex minorant algorithm gives on my computer, after 20 iterations, the following output for the MLE:

	[,1]	[,2]	[,3]	[,4]	[,5]
[1,]	16.43532	0.000000000	0.06250000	0.00000000	0.06250000
[2,]	18.94593	0.000000000	0.08695652	0.00000000	0.08695652
[3,]	19.34292	0.000000000	0.15531609	0.00000000	0.15531609
[4,]	19.80287	0.000000000	0.15531609	0.03060505	0.18592114
[5,]	20.00000	0.008140789	0.15531609	0.03060505	0.19406192
[6,]	21.98494	0.009897753	0.15531609	0.03060505	0.19581889
[7,]	22.06708	0.009897753	0.15801138	0.03060505	0.19851418
[8,]	23.17591	0.055833685	0.15801138	0.03060505	0.24445011
[9,]	23.81656	0.055833685	0.18579624	0.03060505	0.27223497
[10,]	25.09240	0.061458723	0.18579624	0.03060505	0.27786000
[11,]	25.96851	0.068139019	0.18579624	0.03060505	0.28454030
[12,]	26.27789	0.078359872	0.18579624	0.03060505	0.29476115
[13,]	26.66667	0.083361566	0.18579624	0.03060505	0.29976285
[14,]	27.40315	0.117539808	0.18579624	0.03060505	0.33394109
[15,]	27.61944	0.127664960	0.18579624	0.03060505	0.34406624
[16,]	30.53525	0.186595483	0.18579624	0.03060505	0.40299676
[17,]	31.98631	0.186595483	0.18579624	0.03745355	0.40984527
[18,]	34.09993	0.186595483	0.22342172	0.03745355	0.44747075
[19,]	34.59274	0.194506509	0.22342172	0.03745355	0.45538178

The first column gives observations (ages), the second to fourth column the values of the three subdistribution functions of which (at least) one has a jump at the corresponding observation time, and the 5th column gives the values of the sum function F_+ at these points. There are of course many more observations, but the MLE does not have jumps at these points.

After 20 iterations, the values of the estimates satisfy the two (Fenchel duality) conditions (7.56) and (7.57) of [Groeneboom and Jongbloed \(2014\)](#) at the level 10^{-10} (absolute value of inner product of nabla vector and solution smaller than 10^{-10} and minimum of partial sums of the nabla vector bigger than 10^{-10} , respectively). The same convergence criterion is used in `MLEcens`.

If I apply `MLEcens` to the same data, I get the following output on my computer. For the masses on the rectangles I get the vector of values (denoted by $\$p$ in the output):

```
0.062500000 0.024456522 0.068359569 0.030605045 0.008140789 0.001756965
0.002695293 0.045935932 0.027784852 0.005625038 0.006680296 0.010220853
0.005001694 0.034178242 0.010125152 0.058930524 0.006848503 0.037625482
0.007911025 0.544618225
```

If we take the cumulative sums of these numbers, using the R function `cumsum`, we get:

```
0.06250000 0.08695652 0.15531609 0.18592114 0.19406192 0.19581889
0.19851418 0.24445011 0.27223497 0.27786000 0.28454030 0.29476115
0.29976285 0.33394109 0.34406624 0.40299676 0.40984527 0.44747075
0.45538178 1.00000000
```

in which we recognize the last column of the output of the iterative convex minorant algorithm, corresponding to the sum function F_+ , except for the last number (which just gives the jump to the total mass).

As next output, I get the table `$rects` in the output:

	[,1]	[,2]	[,3]	[,4]
[1,]	15.16769	16.43532	1.75	2.25
[2,]	18.94045	18.94593	1.75	2.25
[3,]	19.34018	19.34292	1.75	2.25
[4,]	19.79192	19.80287	2.75	3.25
[5,]	19.99179	20.00000	0.75	1.25
[6,]	21.97399	21.98494	0.75	1.25
[7,]	22.05065	22.06708	1.75	2.25
[8,]	23.15400	23.17591	0.75	1.25
[9,]	23.80287	23.81656	1.75	2.25
[10,]	25.04860	25.09240	0.75	1.25
[11,]	25.94114	25.96851	0.75	1.25
[12,]	26.26968	26.27789	0.75	1.25
[13,]	26.64750	26.66667	0.75	1.25
[14,]	27.40041	27.40315	0.75	1.25
[15,]	27.61670	27.61944	0.75	1.25
[16,]	30.52977	30.53525	0.75	1.25
[17,]	31.97262	31.98631	2.75	3.25
[18,]	34.09719	34.09993	1.75	2.25
[19,]	34.54894	34.59274	0.75	1.25
[20,]	34.99521	100.0000	0.75	3.25

Here one recognizes in the second column the numbers in the first column of the table, produced by the iterative convex minorant. The numbers 1.75, 2.25, etc. in the last columns may seem somewhat peculiar, but have the following meaning. The interval $[0.75, 1.25]$ corresponds to risk 1, the interval $[1.75, 2.25]$ to risk 2, the interval $[2.75, 3.25]$ to risk 3, and the interval $[0.75, 3.25]$ to an observation where there is no failure. This way of coding shows the descentance of the algorithm from the algorithm for bivariate interval censoring.

So, if one places the mass of the intervals at the right endpoint, one gets the same results as with the iterative convex minorant algorithm. It is also seen that the numbers in the first two columns are pretty close, except for the irrelevant last row (as an aside, one could also argue that this last interval should start at 34.59274 rather than at 34.99521), so there is only a slight space of freedom for defining different MLEs.

To further compare the two algorithms for a more challenging data set, we compared the performance for the example of 25000 observations on data, generated by the exponential-type subdistribution functions $F_k(t) = (k/3)\{1 - e^{-kt}\}$, $t \geq 0$, $k = 1, 2$, and analyzed on pages 190 and 191 of [Groeneboom and Jongbloed \(2014\)](#). Using the script `MLE25000.R`, the computation for and plotting of Figure 2 below, where the SMLE also was computed, took about 10 seconds on my computer. On the other hand, just to compute the MLE took 100 seconds using the package `MLEcens`. If the number of observations is growing, the algorithm in `MLEcens` is considerably slowed down by the matrix inversions (or growing number of linear equations to be solved), whereas the iterative convex minorant algorithm does not have to solve equations of this type (only using the diagonal of the Hessian matrix).

3. Bootstrap confidence intervals for the hazards in the competing risk model

Confidence intervals for the hazards in Bangkok data, discussed above, were given in [Li and Fine \(2013\)](#). The intervals, given there, are based on asymptotic theory and do not use the bootstrap. It is not at all trivial that the bootstrap can be used for the confidence intervals, if one wants to base this on the computation of the MLE, since the bootstrap for the MLE itself will no doubt fail, as one can infer from [Kosorok \(2008\)](#) and [Sen, Banerjee and Woodroffe \(2010\)](#). The reason that one still can get the confidence interval by bootstrapping seems to be a consequence of the fact that the estimate of the hazard, based on the MLE, has an asymptotic behavior described by (what is called in [Groeneboom and Jongbloed \(2014\)](#)) *local smooth functional theory*,

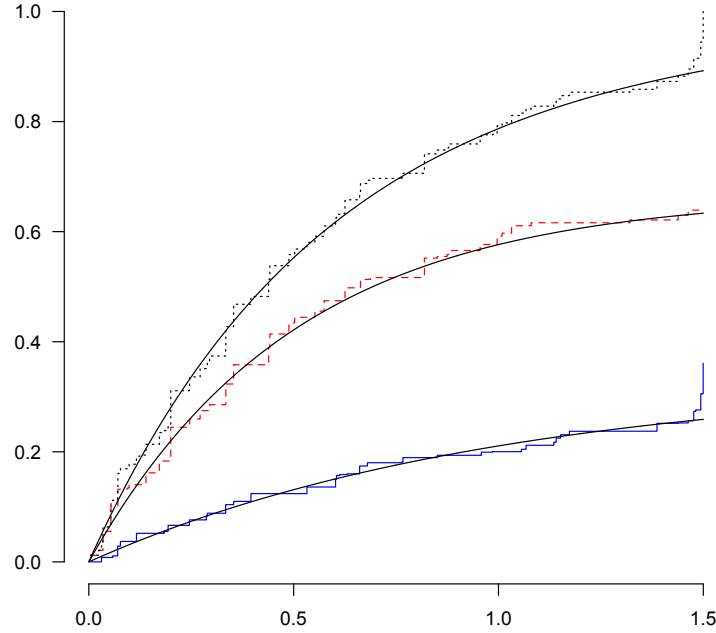


Fig 2: The two distribution functions $F_k(t) = (k/3)\{1 - e^{-kt}\}$ for the competing risk model with current status data and their MLE's for a sample of size $n = 25000$. The upper curves give the sum function F_+ and its estimator \hat{F}_{n+} .

which “washes away” the peculiar local limit behavior of the MLE. Nevertheless, the theory for this is far from complete and is described by integral equations one has to be able to solve first, as also argued in [Groeneboom \(2013\)](#). And once one has solved the integral equations, one has to deduce the relevant equivalence to a “toy estimator” from this, which usually is also a somewhat daunting task.

We will not further dwell on these matters here, but instead describe the R script. We take the data set `dataThai.txt` again and activate the C++ program via the R script `MLEThai.R`. This will produce three hazard estimates for the three groups and also 95% confidence intervals for the hazards, computed at 99 points of the interval $[15, 35]$. Two of the pictures of the confidence intervals can be seen in Figure 12.2 on p. 369 of [Groeneboom and Jongbloed \(2014\)](#). We give the the three pictures below. In [Groeneboom and Jongbloed \(2014\)](#) both studentized and non-studentized confidence intervals are discussed; we give the results for the non-studentized ones, since there was not much difference anyway. The C++ code `icm.cpp` can be consulted to see what is done: 1000 bootstrap samples are drawn (with replacement) from the original data, consisting of the pairs (T_i, Δ_i) , and for each sample the MLE is computed iteratively, using the iterative convex minorant algorithm. Once the MLE is computed, one can compute the SMLE and the density estimate, and on the basis of these the hazard estimate is computed for the bootstrap sample. It is interesting to note that they look somewhat different from the corresponding intervals for the same data in [Li and Fine \(2013\)](#), which might illustrate the distance between asymptotic theory and bootstrap approximation.

4. A GUI application and confidence intervals for the SMLE

As mentioned in the Introduction, we also started writing GUI (Graphical User Interface) applications, parallel to the development of the `Rcpp` approach. These applications were developed for Mac OS X, using

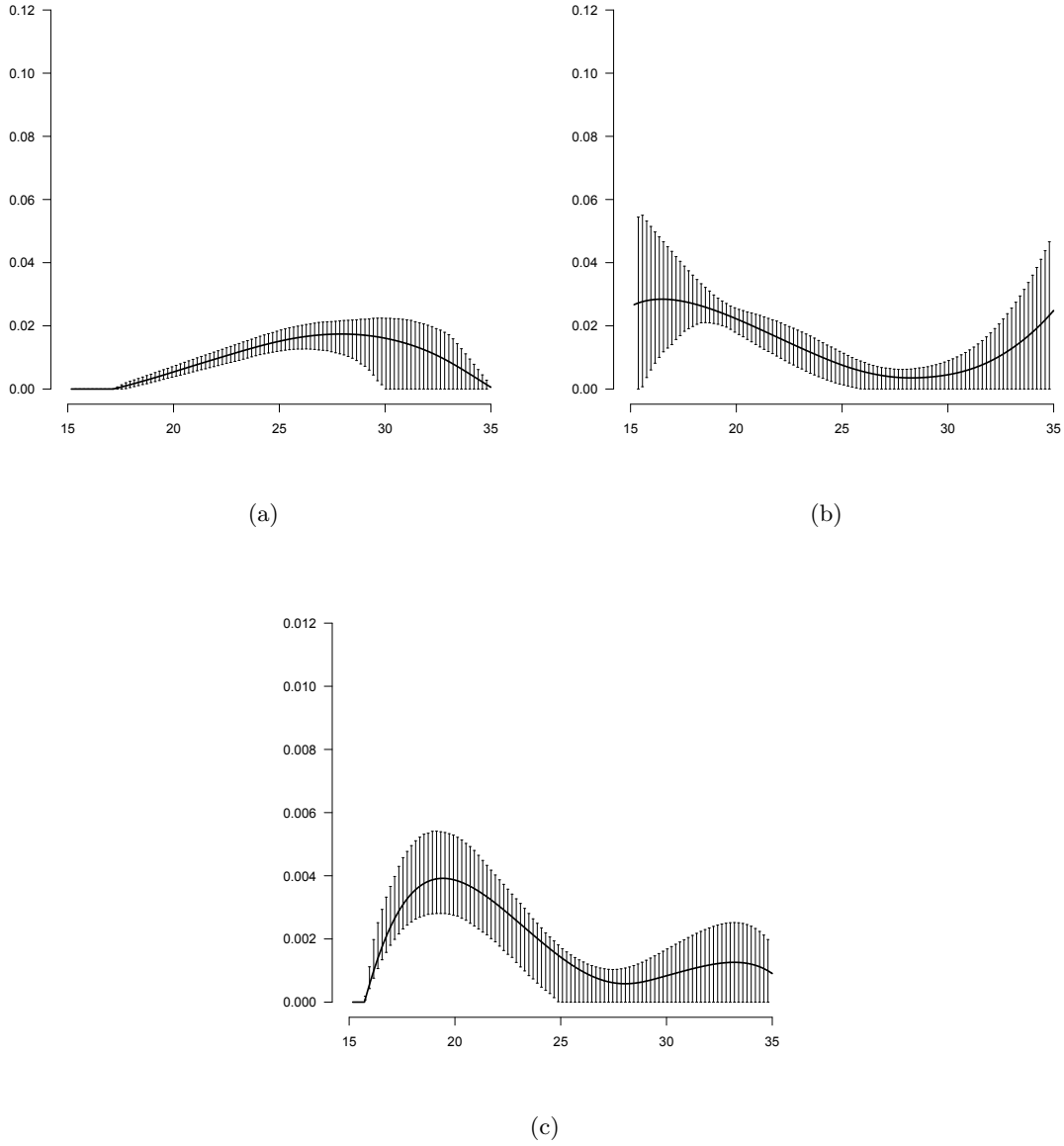


Fig 3: 95% confidence intervals for the hazards for the group infected with type B (Figure (a)), with type E (Figure (b)) and remaining group (Figure (c)) in the Bangkok cohort data.

Xcode and a combination of `Objective-C` and `C++`, but it should also be possible to implement these for Windows, using Microsoft Visual `C++`.

The application allows the user to specify the number of bootstrap samples to be taken for the computation of the confidence intervals. It has a menu, where one can choose the input file by the standard “File-open ...” submenu. One can also choose via the submenu “Statistic” between confidence intervals for the SMLE or confidence intervals for the hazard. One can both use the ungrouped and grouped file formats. In the latter case the observations are unique and at each such unique point the information about the frequencies of failures for the competing risks is given, and the frequency of observations of no failure at that point. The screenshot Figure 4 shows a moment in the analysis of the Bangkok data, where the number of observations was 1365, and where after the correction for ties 1211 unique observations remained. To have an idea of

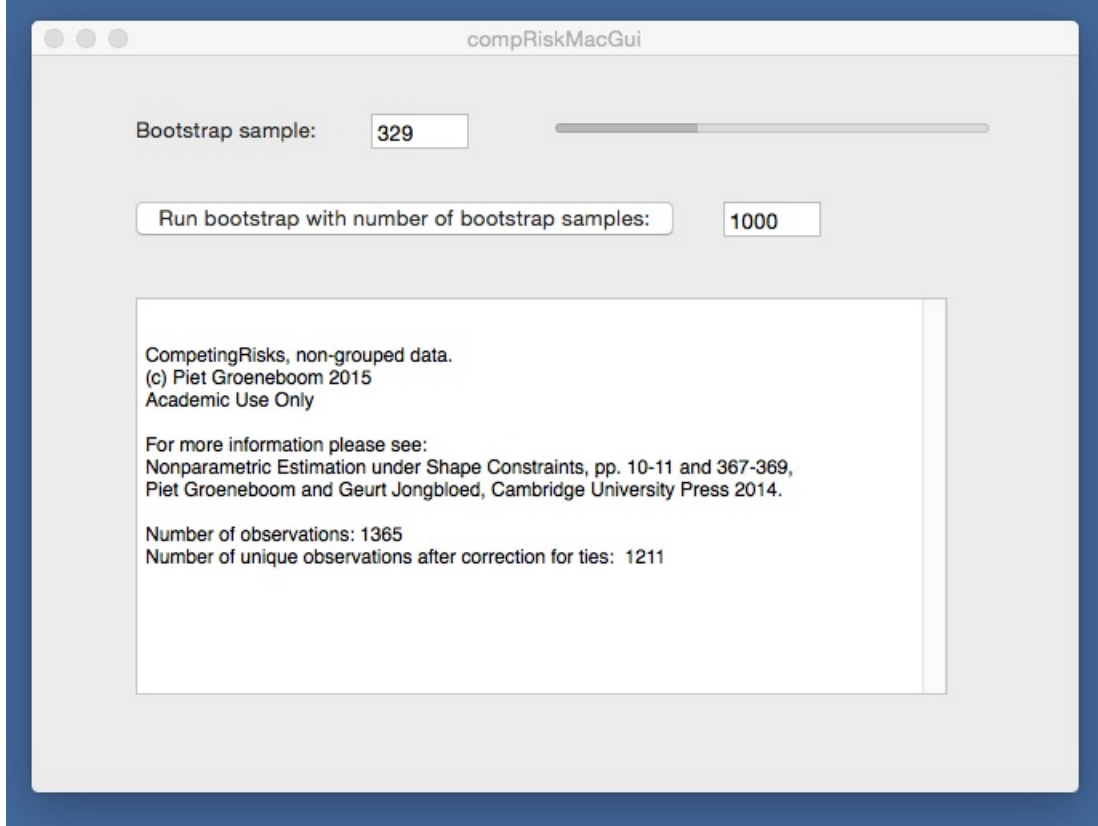


Fig 4: A screenshot of the GUI application

the progress of the bootstrap samples, there is a (blue) progress bar, showing the proportion of bootstrap samples that have been treated. The number of treated bootstrap samples is also shown in the (so-called) textfield left if it. See the discussion on Section 2.5 for the difference between the input file with two columns and 1365 observation points for the ungrouped data and the input file with $K + 2$ (in this case: 5) columns and 1211 (unique) observation points for the grouped data.

For performing the bootstrap, it seemed easiest to stick with the 1365 observation points, from which we resample uniformly 1365 bootstrap sample points of type (T_i, Δ^i) and compute the statistics for the bootstrap sample. This means that in the case of the input of the 1211 (grouped) data points, the data were first transformed to the other format, where we only have two columns with 1365 points (T'_i, γ_i) , where $\gamma_i \in \{0, 1, \dots, K\}$ denotes the risk for which there is a failure at time T'_i (or equals zero if it denotes an observation for which there is no failure).

Since the computations are again done by the C++ part of the program, the computations take about the same time as in the Rcpp implementation. For example, the bootstrap experiment with 1000 bootstrap samples of 1365 observations for the confidence intervals for the hazards takes about 20 seconds (the amount of time is computed in the program and given in the printed output). One can also compute the 95% confidence intervals for the SMLE. This takes somewhat longer on my computer, about 50 seconds, since in this case the book on p. 367 is (again) followed, in taking a kind of “studentized” confidence intervals, for which an estimate of the variance has to be computed for each sample, see (12.7) on p. 367. This was done in view of the “reproducibility” of Figure 12.1 on p. 368. It turns out, however, that, just as in the case of the confidence intervals for the hazards, the non-studentized confidence intervals for the SMLE, not using this scaling by the estimate of the standard deviation, are rather similar.

For the computation of the confidence intervals for the MLE for current status and interval censored data, a different kind of bootstrap is proposed in Sen and Xu (2015). They suggest bootstrapping from the SMLE, by resampling only the indicators, with probabilities given by the SMLE, and using the (fixed) observation

points of the original sample. For example, in the current status model the bootstrap samples contain the observations (T_i, Δ_i^*) , where the T_i are the observations of the original sample, and Δ_i^* is, conditionally on the original sample, a Bernoulli random variable, with success probability $\tilde{F}_{n,h}(T_i)$ and where $\tilde{F}_{n,h}$ is an SMLE, with a bandwidth h , tending to zero slower than $n^{-1/3} \log n$. They prove in their Theorem 3 in Section 2.3 (in the supplementary material to their paper) that this leads to consistent bootstrapping of the MLE. This could also be extended to confidence intervals for the SMLE in the competing risk model, by smoothing the bootstrap values of the MLE.

In contrast, we would, for the current status model, resample the (T_i, Δ_i) (so both the values T_i and Δ_i become random in our bootstrap sample) and compute the MLE and subsequently the SMLE directly on the basis of such a bootstrap sample. It is indeed surprising that this leads to inconsistent estimates of the distribution of the MLE, but consistent estimates of the corresponding distribution of the SMLE. The latter seems to be a consequence of the fact that the SMLE, although based on the MLE, seems to “by-pass” the aberrant limit behavior of the bootstrap version of the MLE.

To apply the method, suggested in Sen and Xu (2015) in the present situation, we generated the distribution of the $(\Delta^i)^*$ from the corresponding multinomial distribution, given by the parameters $\tilde{F}_{n,h}^{(1)}(T_i), \dots, \tilde{F}_{n,h}^{(K)}(T_i)$, where $\tilde{F}_{n,h}$ is an SMLE, computed with bandwidth h . Some preliminary experiments with this method suggest that the bias plays a much larger role here and that the method of computing the confidence intervals, given by our method, might be preferable. It is also not clear how one can obtain confidence intervals for the hazards, using the method of Sen and Xu (2015). More research is clearly needed here.

The bootstrap confidence of our method for the Thailand data are shown in Figure 5.

5. Acknowledgements

I want to thank Richard Gill for testing all **Rcpp** algorithms on Mac, Windows and Linux. Harold Hvistendahl tested the original C++ program for bootstrapping on Windows (where it seemed to work fine) and the original **Rcpp** version of this (which first crashed in all environments). He also added some text to the output, which seemed a good idea. Dimitri Tischenko spotted a bug in the first **Rcpp** version, using Valgrind. After repairing this bug, we did not experience crashes any more. Frank van der Meulen tested the GUI application on his MacBook Pro with system 10.10. I feel very grateful for their contributions.

References

- EDDELBUEITTEL, D. (2013). *Seamless R and C++ Integration with Rcpp*. Springer, New York.
- GROENEBOOM, P. (2013). Nonparametric (smoothed) likelihood and integral equations.
- GROENEBOOM, P., JONGBLOED, G. and WELLNER, J. A. (2008). The support reduction algorithm for computing non-parametric function estimates in mixture models. *Scand. J. Statist.* **35** 385–399. [MR2446726 \(2009m:62115\)](#)
- GROENEBOOM, P., JONGBLOED, G. and WITTE, B. I. (2010). Maximum smoothed likelihood estimation and smoothed maximum likelihood estimation in the current status model. *Ann. Statist.* **38** 352–387.
- GROENEBOOM, P. and JONGBLOED, G. (2014). *Nonparametric Estimation under Shape Constraints*. Cambridge Univ. Press, Cambridge.
- GROENEBOOM, P., MAATHUIS, M. H. and WELLNER, J. A. (2008). Current status data with competing risks: consistency and rates of convergence of the MLE. *Ann. Statist.* **36** 1031–1063. [MR2418648 \(2009h:62039\)](#)
- KITAYAPORN, D., VANICHSENI, S., MASTRO, T. D., RAKTHAM, S., VANIYAPONGS, T., DES JARLAIS, D. C., WASI, C., YOUNG, N. L., SUJARITA, S., HEYWARD, W. L. and ESPARZA, J. (1998). Infection with HIV-1 subtypes B and E in injecting drug users screened for enrollment into a prospective cohort in Bangkok, Thailand. *J. Acquir. Immune Defic. Syndr. Hum. Retrovirol.* **19** 289–295.
- KOSOROK, M. R. (2008). Bootstrapping the Grenander estimator. In *Beyond parametrics in interdisciplinary research: Festschrift in honor of Professor Pranab K. Sen. Inst. Math. Stat. Collect.* **1** 282–292. Inst. Math. Statist., Beachwood, OH.
- LI, C. and FINE, J. P. (2013). Smoothed nonparametric estimation for current status competing risks data. *Biometrika* **100** 173–187. [MR3034331](#)

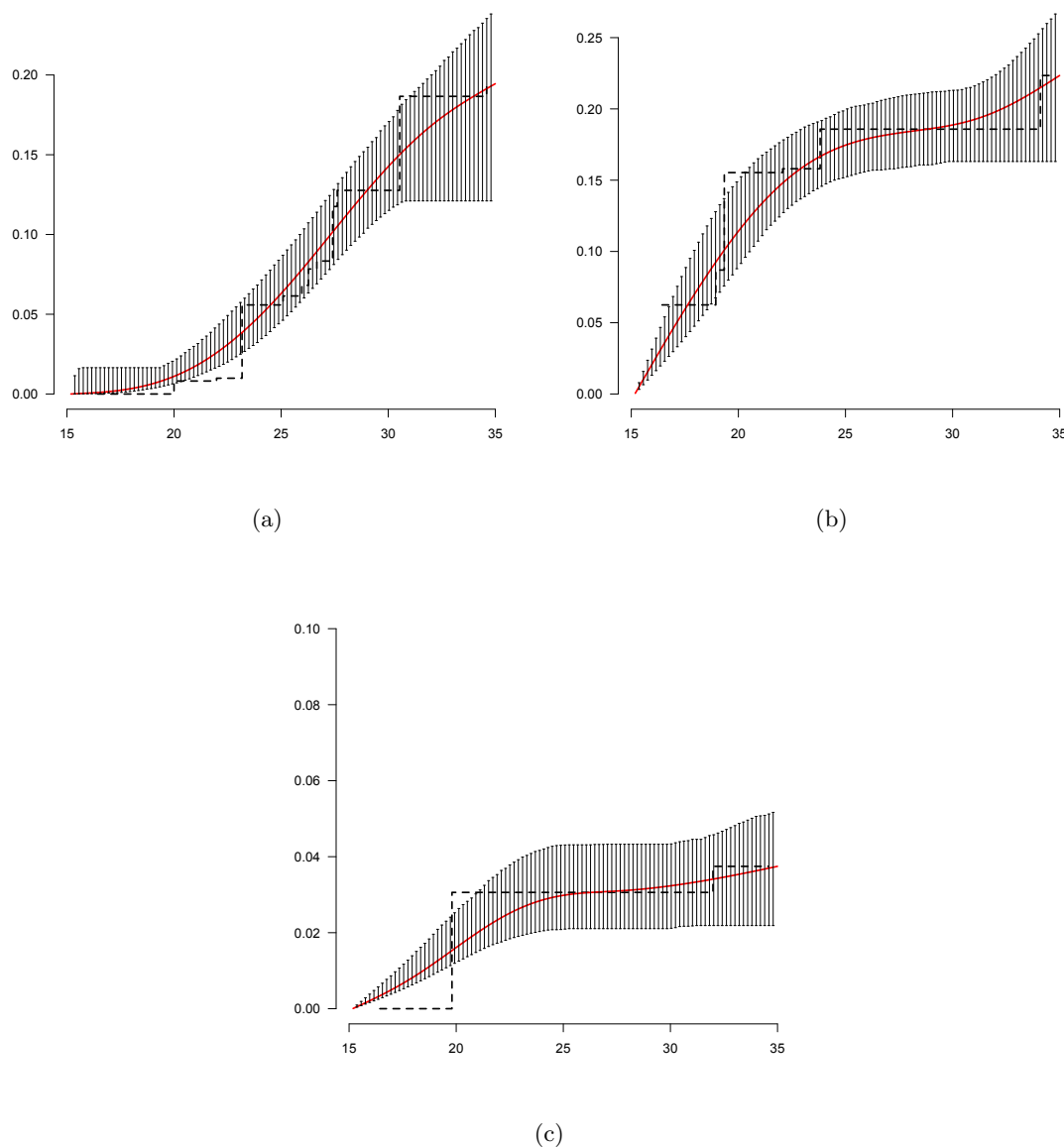


Fig 5: 95% confidence intervals for the subdistribution functions for the group infected with type B (Figure (a)), with type E (Figure (b)) and remaining group (Figure (c)) in the Bangkok cohort data.

MAATHUIS, M. H. (2013). MLEcens. R package. Version 0.1-4.

MAATHUIS, M. H. and HUDGENS, M. G. (2011). Nonparametric inference for competing risks current status data with continuous, discrete or grouped observation times. *Biometrika* **98** 325–340. [MR2806431 \(2012e:62108\)](#)

PRESS, W. H., TEUKOLSKY, S. A., VETTERLING, W. T. and FLANNERY, B. P. (1992). *Numerical recipes in C*, Second ed. Cambridge University Press, Cambridge The art of scientific computing. [MR1201159 \(93i:65001b\)](#)

SEN, B., BANERJEE, M. and WOODROOFE, M. B. (2010). Inconsistency of bootstrap: the Grenander estimator. *Ann. Statist.* **38** 1953–1977. [MR2676880 \(2011f:62046\)](#)

SEN, B. and XU, G. (2015). Model based bootstrap methods for interval censored data. *Comput. Statist.*

Data Anal. **81** 121–129. [MR3257405](#)
VANICHSENI, S., KITAYAPORN, D., MASTRO, T. D., MOCK, P. A., RAKTHAM, S., C., D. J. D., SUJARITA, S., SRISUWANVILAI, L. O., YOUNG, N. L., WASI, C., SUBBARAO, S., HEYWARD, W. L., ESPARZA, L. and CHOOPANYA, K. (2001). Continued high HIV-1 incidence in a vaccine trial preparatory cohort of injection drug users in Bangkok, Thailand. *AIDS* **15** 397–405.