

SURVIVAL ANALYSIS

(Lecture Notes)

by Qiqing Yu

Version 1/18/2018

This course will cover parametric, non-parametric and semi-parametric maximum likelihood estimation, in the Cox regression model and the linear regression model, with complete data and various types of censored data. The right censorship model, double censorship model, the mixed case interval censorship model, the mixture censorship model and the left-truncation model will be used to formulate the censored or truncated-data. Multivariate censorship models and the two-sample problem may also be introduced.

Math 557 MWF 8-9:30 am

Office: WH 132

Office hours: M. T. 3:30-4:30pm

Homework due: Wednesday;

Exam: March 9 (F)

Final: May 15 (Tu) 12:50pm -2:50pm (WH G002) closed book exam.

Grading: 50%hw + 20%exam + 30%final;

Reference Textbooks:

1. Analysis of survival data, by Cox and Oakes.
2. Survival Analysis, by Rupert G. Miller, JR.
3. The Statistical Analysis of Interval-censored Failure Time Data, by J. Sun.

Chapter 1. Introduction.

§1. Two main characters of survival analysis.

Suppose X is a random variable,
with the cumulative distribution function (cdf),

$$F(x) = P(X \leq x).$$

Assume X_1, \dots, X_n are independently and identically distributed (i.i.d.) from F .

Statistical inferences: $\begin{cases} \text{estimation} \\ \text{testing hypotheses,} \end{cases}$

on $S = 1 - F$ or the parameters in a certain model.

In particular, there are three different types:

parametric, *i.e.* $F(x) = F_o(x; \theta)$, where $\theta \in \Theta$, a parameter space;

nonparametric, *i.e.*, F is only known to be a c.d.f.

semi-parametric, *i.e.*, in between the above two.

In survival analysis, X is often

time to death of a patient after a treatment,

time to failure of a part of a system, etc.

One main character is $X \geq 0$, referred as *survival time or failure time*.

By S , it is much intuitive for doctors to compare different treatments or different systems,

the larger the survival probability S , the better.

Another main character of survival analysis:

Incomplete data.

Definition: An observation on X_i is called

$\begin{cases} \text{complete (exact, or uncensored)} & \text{if the exact value of } X_i \text{ is observed;} \\ \text{incomplete} & \text{o.w.;} \end{cases}$

A data set is called $\begin{cases} \text{complete} & \text{if all observations are exact;} \\ \text{incomplete} & \text{otherwise.} \end{cases}$

An incomplete observation on X_i is called *interval censored (IC)* if $X_i \in I_i$, an interval with endpoints L_i and R_i being observed.

An IC observation can thus be represented by an **extended** random vector (L_i, R_i) .

(L_i, R_i) is called $\begin{cases} \textit{right censored (RC)} & \text{if } R_i = \infty \\ \textit{left censored (LC)} & \text{if } L_i = -\infty \\ \textit{strictly interval censored (SIC)} & \text{if } 0 < L_i < R_i < \infty. \end{cases}$

§1.2. **Right censoring.**

§1.2.1. Representations of an RC observation:

(L_i, R_i) — a vector,
 $I_i = (L_i, R_i)$ — an interval,
 $(L_i, 0)$ — a vector,
 L_i+ ,

Representations of an exact observation:

$(L_i, R_i) = (X_i, X_i)$
 $I_i = [X_i, X_i]$ — an interval,
 $(L_i, 1)$ — a vector,
 L_i .

Definition: A RC data set — observations are either exact or RC, but there exist exact observations. Representations of RC data:

$(L_i, R_i), i = 1, \dots, n,$ — random vectors,
 where $(L_i, R_i) = \begin{cases} (X_i, X_i) & \text{if the observation is exact,} \\ (L_i, \infty) & \text{if the observation is RC,} \end{cases}$
 $I_i, i = 1, \dots, n,$ — random intervals,
 where $I_i = \begin{cases} [X_i, X_i] & \text{if the observation is exact,} \\ (L_i, \infty) & \text{if the observation is RC.} \end{cases}$
 $(L_i, \delta_i), i = 1, \dots, n,$ — random vectors,
 where $\delta_i = \mathbf{1}_{(\text{the } i\text{-th observation is exact})}$, and
 $\mathbf{1}_A = \begin{cases} 1 & \text{if } A \text{ happens} \\ 0 & \text{o.w.} \end{cases}$ is the indicator function of the event A .
 L_i+ or $L_i, i = 1, \dots, n$.

Example of RC data:

1. Mortality data (population census, for computing the life expectancy of the population).

Let X_i be the age at which the i -th person died.
 Then at a census, we either knew X_i if the person died
 or knew L_i+ if he/she was alive, where L_i was his/her age then.

2. Type I censoring.

Each individual was followed by a fixed time c .
 Each X was recorded unless $X > c$.

$$X_{(1)}, \dots, X_{(i)}, c+, \dots, c+$$

where $X_{(1)} \leq \dots \leq X_{(i)}$ are order statistics of observed exact values of X_1, \dots, X_n .

3. Type II censoring.

Observation ceases after a predetermined number d of failures.

$$X_1, \dots, X_d, c, \dots, c, c+, \dots, c+$$

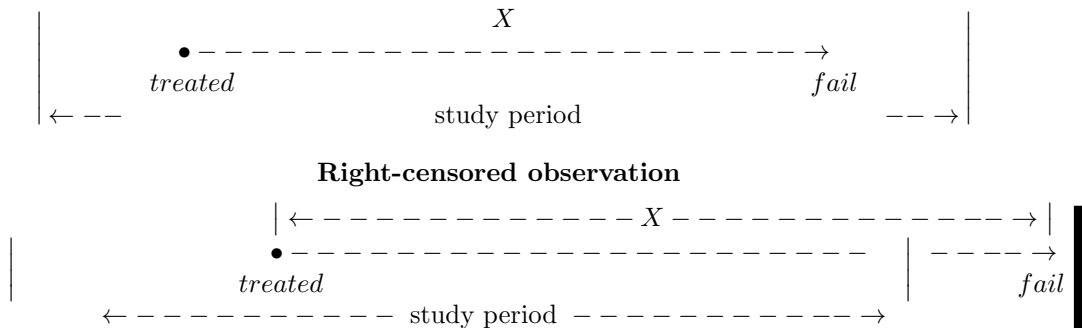
where $X_d = c$.

4. Random censoring.

In a medical follow-up study of 5 years, n cancer patients are enrolled (not necessary from the beginning). X is the time to death of a patient since a certain treatment (after the enrollment). We either know X or know $X > 5 - B$, where B is the beginning time of the treatment for the individual since the start of the study.

Graphical illustration for X and Y

Exact observation



Leukaemia data

Gehan, 1965 recorded times of remission of leukaemia patients.
 Some were treated with drug 6-mercaptopurine (6-MP),
 the others were serving as a control.

Table 1.1 (Cox and Oakes (1984) (pages 7,8)). Time of remission (weeks).

Group 0 (6-MP): 6+, 6, 6, 6, 7, 9+, 10+, 10, 11+, 13, 16, 17+, 19+, 20+, 22, 23, 25+, 32+, 32+, 34+, 35+ (m=21),

Group 1 (control): 1, 1, 2, 2, 3, 4, 4, 5, 5, 8, 8, 8, 8, 11, 11, 12, 12, 15, 17, 22, 23 (n=21).

§1.2.2. **The right censorship model (RC model)** (Kaplan and Meier, 1958, JASA).

Assume:

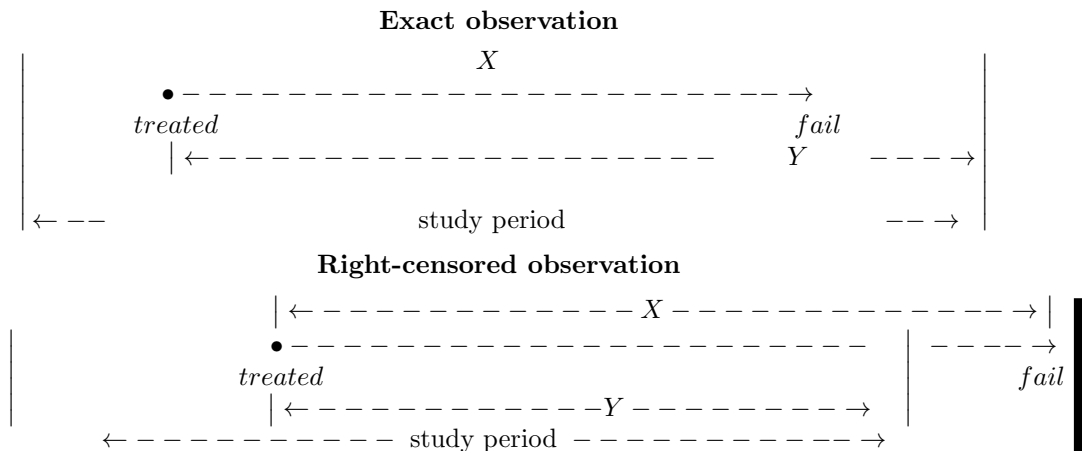
Y — random censoring variable.

X and Y are independent ($X \perp Y$) i.e., $P(X \leq x, Y \leq y) = P(X \leq x)P(Y \leq y)$.

Observable random vector $(L, R) = \begin{cases} (X, X) & \text{if } X \leq Y \\ (Y, \infty) & \text{if } X > Y. \end{cases}$

Equivalent observations: $(Z, \delta) = (\min(X, Y), \mathbf{1}_{(X \leq Y)})$.

Graphical illustration for Y and X



Note: If $Y = c$ w.p.1, it becomes type I censoring. **How to express in the graph ?**

§1.2.3. Two incorrect approaches for RC data (before 1958):

Method 1. Discard all RC observations;

Method 2. Treat RC observations as exact observations.

Question: What is wrong?

Reasoning:

For complete data, X_1, \dots, X_n (i.i.d. from X),
 if $\mu = E(X)$ exists, an estimator of μ is \bar{X} .

Its properties:

$E(\bar{X}) = \mu$ (unbiased);

$\bar{X} \rightarrow \mu$ w.p.1. (strongly consistent);

A nonparametric estimator of F and S are

$$\hat{F}(x) = \sum_{i=1}^n \mathbf{1}_{(X_i \leq x)} / n (= \overline{W}) \quad \mathbf{W}=? \quad \text{and} \quad \hat{S}(x) = \sum_{i=1}^n \mathbf{1}_{(X_i > x)} / n (= \overline{Z})$$

—Empirical distribution function (edf) and Empirical survival function, respectively.

Their properties:

$$E(\hat{F}(x)) = F(x) \text{ (unbiased) } \mathbf{Why?};$$

$$\hat{F}(x) \rightarrow F(x) \text{ w.p.1. (strongly consistent) } \mathbf{Why?};$$

Using method 1, we have

$$m = \sum_{i=1}^n \mathbf{1}_{(X_i \leq Y_i)} \text{ exact observations.}$$

So the “sample mean” is

$$\tilde{\mu} = \frac{\sum_{i=1}^n X_i \mathbf{1}_{(X_i \leq Y_i)}}{\sum_{i=1}^n \mathbf{1}_{(X_i \leq Y_i)}}$$

Question:

$$E(\tilde{\mu}) = \mu? \tag{*}$$

Or

$$E(\tilde{\mu}) = \frac{E(\sum_{i=1}^n X_i \mathbf{1}_{(X_i \leq Y_i)})}{E(\sum_{i=1}^n \mathbf{1}_{(X_i \leq Y_i)})} = \frac{E(X_i \mathbf{1}_{(X_i \leq Y_i)})}{E(\mathbf{1}_{(X_i \leq Y_i)})} ? \tag{**}$$

Counter-example: Suppose $n = 2$, X_i, Y_i are i.i.d. from $\text{bin}(1, 1/2) + 1$.

$$\mu = E(X) = ?$$

$$\vdash: E(\tilde{\mu}) \neq \mu.$$

$$E(\tilde{\mu}) = \begin{cases} \int t f_{\tilde{\mu}}(t) dt & ? \\ \sum_t t f_{\tilde{\mu}}(t) & ? \end{cases}$$

case	X_1	Y_1	X_2	Y_2	$\tilde{\mu}$	p	$p \times \tilde{\mu}$
1	1	1	1	1	1	1/16	1/16
2	2	1	2	1	∞	1/16	∞
3	1	2	1	1		1/16	
4	2	2	2	1		1/16	
.
<i>sum</i>							∞

Thus $\tilde{\mu}$ is biased, and both (*) and (**) fail.

Remark. Here we define $\frac{0}{0} = \infty$. If we define $\frac{0}{0} = 1$, we can also show $\tilde{\mu}$ is biased.

Question: $\tilde{\mu} \xrightarrow{a.s.} \mu$?

$$\begin{aligned} \lim_{n \rightarrow \infty} \tilde{\mu} &= \frac{\lim_{n \rightarrow \infty} \sum_{i=1}^n X_i \mathbf{1}_{(X_i \leq Y_i)} / n}{\lim_{n \rightarrow \infty} \sum_{i=1}^n \mathbf{1}_{(X_i \leq Y_i)} / n} \\ &= \frac{\lim_{n \rightarrow \infty} \overline{U}}{\lim_{n \rightarrow \infty} \overline{V}} \\ &= E(U) / E(V) \text{ w.p.1.} \\ &= \frac{E(X_1 \mathbf{1}_{(X_1 \leq Y_1)})}{E(\mathbf{1}_{(X_1 \leq Y_1)})} \\ &= ? \\ &= ? \end{aligned}$$

The second counterexample. If $Y \equiv 1$ and X has a df $f(x) = \begin{cases} 1/3 & \text{if } x \in (0, 1) \\ 2/3 & \text{if } x \in (1, 2). \end{cases}$, then

(1) $\mu = 7/6$, but w.p.1 $\lim_{n \rightarrow \infty} \tilde{\mu} = \frac{1/6}{1/3} = 1/2$.

(2) The “edf” is

$$\tilde{F}(x) = \frac{\sum_{i=1}^n \mathbf{1}_{(X_i \leq Y_i, X_i \leq x)}}{\sum_{i=1}^n \mathbf{1}_{(X_i \leq Y_i)}}$$

W.p.1. its limit is

$$\frac{E(\mathbf{1}_{(X \leq Y, X \leq x)})}{E(\mathbf{1}_{(X \leq Y)})} (= \frac{P(X \leq 1, X \leq x)}{P(X \leq 1)}).$$

When $x = 1$, $\lim_{n \rightarrow \infty} \tilde{F}(1) = 1 \neq 1/3 = F(1)$.

Thus both estimators are inconsistent.

Homework : Verify (1) and (2).

Remark. In the derivation of Examples 1 and 2, we are making use of the right censorship model. These examples illustrate the importance of a correct approach in dealing with censored data.

§1.2.4. **Homework:**

1. Using the second method, are the modified edf and sample mean consistent? Justify your statement.

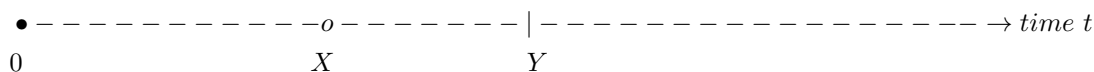
§1.3. **Case 1 interval-censoring.**

Definition. If a data set only contains RC observations and LC observations, it is called *case 1 IC data* (C1 data) or *current status data*.

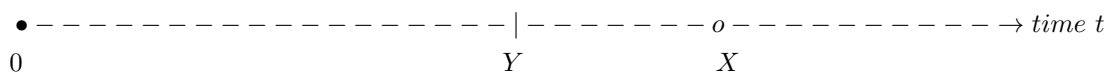
Example. Consider an animal sacrifice study in which a laboratory animal has to be dissected to check whether a tumor has developed. In this case, X is the onset of tumor and Y is the time of the dissection, and we only can infer at the time of dissection whether the tumor is present or has not yet developed.

Other examples are mentioned in Ayer *et al.* (1955), Keiding (1991) and Wang and Gardiner (1996).

LC :



RC :



The case 1 interval censorship model:

Assume:

Y is a random inspection time;

X and Y are independent;

The observable random vector is $(L, R) = \begin{cases} (-\infty, Y) & \text{if } X \leq Y \\ (Y, \infty) & \text{if } X > Y. \end{cases}$

Equivalent forms:

Interval: $I = \begin{cases} (-\infty, Y] & \text{if LC} \\ (Y, \infty) & \text{if RC.} \end{cases}$

vector: (Y, δ) , where $\delta = \mathbf{1}_{(X \leq Y)}$.

Given a sample from the C1 model, is it possible all are right censored ?

§1.4. **Double censoring**

Definition. If a data set contains RC, LC and exact observations, but not SIC observations, it is called a *doubly-censored data* (DC data).

Example. Leiderman *et al.* (1973) presented a study on the time needed for an infant to learn to perform a particular task (crawling) during the first year. The sampled infants were all born within 6 months of the start of the study. At the time of the start of the study, some children had already known how to perform the task; so their observed times were left-censored. Some children learned the task during the time-span of the study, and their ages were recorded. At the end of the study, some of the children had not yet learned the task, and hence their observed times were right-censored.

The double censorship model:

Assume

(Z, Y) is a random censoring vector;

$Z \leq Y$ w.p.1.;

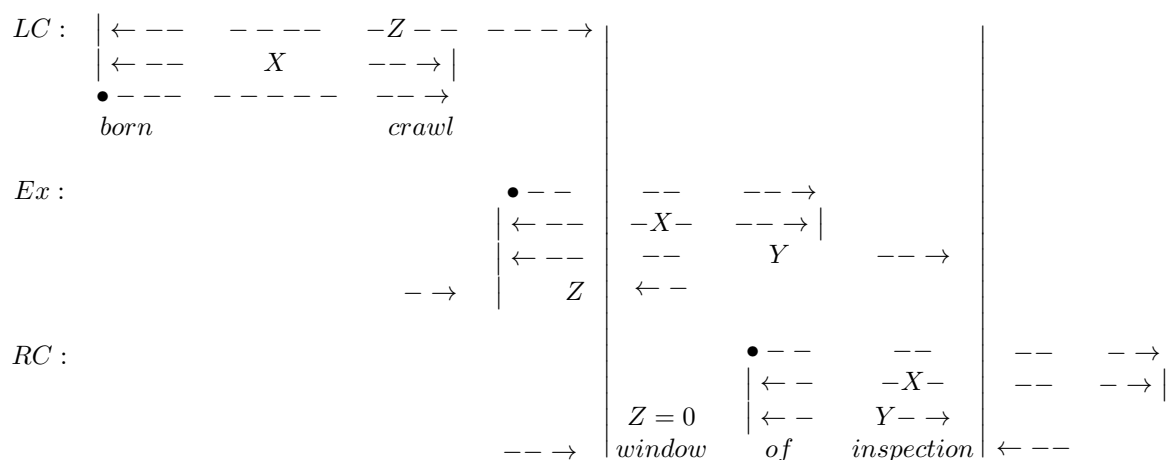
X and (Z, Y) are independent;

The observable random vector is $(L, R) = \begin{cases} (-\infty, Z) & \text{if } X \leq Z \\ (X, X) & \text{if } Z < X \leq Y \\ (Y, \infty) & \text{if } X > Y. \end{cases}$

Equivalent forms:

Interval: $I = \begin{cases} (-\infty, Z] & \text{if } X \leq Z \\ [X, X] & \text{if } Z < X \leq Y \\ (Y, \infty) & \text{if } X > Y. \end{cases}$

Vector (U, δ) , where $U = \max\{Z, \min\{X, Y\}\}$, $\delta = \begin{cases} 1 & \text{if exact} \\ 2 & \text{if RC} \\ 3 & \text{if LC.} \end{cases}$



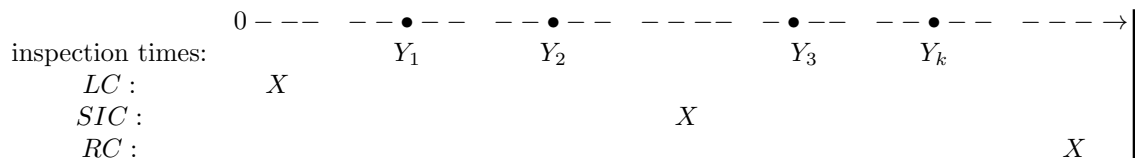
Remark. We could also set $Z = \text{time at birth} - \text{time to starting study}$, which may take negative values.

§1.5. Case 2 interval censoring.

Definition If a data set contains SIC observations and/or RC or LC observations, but not exact observations, it is called a *case 2 IC data* (C2 data).

Example. In medical research when each patient had several follow-ups and the event of interest was only known to take place either before the first follow-up, or between two consecutive follow-ups, or after the last one.

Graphical illustration



Examples of C2 data can be found in breast cancer research (Finkelstein and Wolfe, 1985) and AIDS studies (Becker and Belbye, 1991).

Possible models:

1.5.1. A simple model:

Groeneboom and Wellner (1992) proposed the following case 2 interval-censorship model.

Assume:

1. U and V are random inspection times such that $U < V$ w.p.1;
2. X and (U, V) are independent;

3. The observable random vector is

$$(L, R) = \begin{cases} (-\infty, U) & \text{if } X \leq U \\ (U, V) & \text{if } U < X \leq V \\ (V, \infty) & \text{if } X > V. \end{cases}$$

Remark. In a follow-up study, each patient has \mathcal{N} visits, where $\mathcal{N} \geq 1$ is a random integer. The inspection times are $Y_1 < \dots < Y_{\mathcal{N}}$. It is reasonable to assume that X and $(\mathcal{N}, \{Y_i : i \geq 1\})$ are independent. Then, on the event $\{\mathcal{N} = k\}$, define

$$(U, V) = \begin{cases} (Y_1, Y_2) & \text{if } X \leq Y_1 \\ (Y_{k-1}, Y_k) & \text{if } X > Y_k \\ (Y_{i-1}, Y_i) & \text{if } Y_{i-1} < X \leq Y_i, i \in \{2, \dots, k\} \end{cases}$$

where $Y_0 = 0$. Then (U, V) and X are not independent. In other words, the case 2 model is simple, but its assumption is not realistic.

1.5.2. An alternative model: Wellner realized the drawback and proposed the Case k IC model, in which each patient has exactly k visits. Case 1 and Case 2 are special case of the case k models. However, it is not realistic except the case 1 model.

1.5.3. Another model: Petroni and Wolfe (1994) assume that inspection times Y_j 's can only be taken at y_i 's, where

$$y_1 < y_2 < \dots < y_k,$$

which are predetermined (can be viewed as the reservation time), and $q(y_i) = P(\text{a patient keeps the appointment at time } y_i) \in (0, 1]$.

This results in inspection times $Y_1 < \dots < Y_N$, together with their distribution, where N is again a random integer. Define $Y_0 = -\infty$ and $Y_{N+1} = \infty$.

Assume that X and (Y_1, Y_2, \dots) are independent.

Then $(L, R) = (Y_{j-1}, Y_j)$ if $X \in (Y_{j-1}, Y_j]$ for some j .

This model assumes that Y_j 's are discrete, taking only k values.

The real data are actually continuous.

1.5.4. A realistic model (mixed case IC model, Schick and Yu (2000)):

Assume:

N is a random positive integer;

$Y_1 < Y_2 < \dots < Y_k < \dots$ are inspection times;

Conditional on $N = k$, X and $\{Y_1, \dots, Y_k\}$ are independent

(or for simplicity, X , N , and $(Y_i, i \geq 1)$ are independent);

The observable random vector is

$$(L, R) = \begin{cases} (-\infty, Y_1) & \text{if } X \leq Y_1 \\ (Y_i, Y_{i+1}) & \text{if } Y_i < X \leq Y_{i+1}, i = 1, \dots, N-1 \\ (Y_N, \infty) & \text{if } X > Y_N. \end{cases}$$

Remark .

1. When $N = k$ w.p.1, it is called a case k model, where $k = 1, 2, \dots$

2. The mixed case IC model can be viewed as a mixture of various case k models.

3. When Y_i 's are discrete, then the mixed case IC model becomes the model in §1.5.3.

Example of generating 100 observations under the mixed case IC model through simulation.

Main idea:

generate $X \sim f_X$;

generate $N \sim f_N$;

generate $Y_1, \dots, Y_N \sim f_{\mathbf{Y}}$;

find j such that $Y_{j-1} < X \leq Y_j$ to obtain (L, R) .

repeat 100 times.

Example of generating Y_i 's:

```

(1) Generate  $N$   $Z_i \sim \text{exp}(\lambda)$ ,  $Y_1 = Z_1, Y_2 = Y_1 + Z_2, \dots, Y_N = Y_{N-1} + Z_N$ .
(2) generate  $Y_i \sim U(2i, 2i + 2)$ ,  $i = 1, \dots, N$ .
Assume  $N \sim \text{Poisson}(5)+1$ ,  $X \sim \text{exp}(3)$  (mean=3),  $Y_i \sim U(2i, 2i + 2)$ .
L=rep(0,100) # initialize L
R=rep(0,100) # initialize R
for (i in 1 : 100) { # loop for 100 data
N=rpois(1,5) + 1 # generate 1 random variable from Poisson(5) +1
X=rexp(1,1/3) # generate 1 random variable from exp(3)
#J = 1
#if (N>1)
#for (j in 2:N) J = c(J,j) # J is a vector of (1,2,...,N)
J=1:N # J is a vector of (1,2,...,N)
Y=runif(N,0,2)+2*J # generate N rv from U(2j,2j+2), j=1,...,N
if (X <= Y[1]) {
L[i] = 0
R[i] = Y[1]
}
else {
if (X > Y[N]) {
L[i] = Y[N]
R[i] = 1000
}
else {
j=length(Y[Y<X])+1 # j=2, while (X > Y[j] & j <= N) j=j+1
L[i] = Y[j-1]
R[i] = Y[j]
} } }
U=c(L,R)
dim(U)=c(100,2) # matrix of dimension 100 x 2
U # print the matrix

```

Question: What is wrong with the following scheme ?

```

X=rexp(100)
L=X-1
R=X+1

```

§1.6. Mixed IC censoring.

Definition. If a data set contains both exact and SIC observations, and/or RC or LC observations, it is called a *mixed IC data* (MIC data). It is also called partly IC data.

Example. (the National Longitudinal Survey of Youth 1979-98 (NLSY)). The 1979-98 cross-sectional and supplemental samples consist of 11,774 respondents, who were between the ages of 14 and 22 in 1979. Interviews were conducted yearly from 1979 through 1994; since then data were recorded bi-annually. One entry is the age at first marriage. There are SIC, exact, RC and LC observations in the data.

Possible models:

A simple model (MIC model (1), Yu *et al.* (1995)):

Assume:

1. (U, V) is an extended random censoring vector such that $U < V$ w.p.1;
2. X and (U, V) are independent;
3. The observable random vector is

$$(L, R) = \begin{cases} (X, X) & \text{if } X \notin (U, V] \\ (U, V) & \text{if } X \in (U, V]. \end{cases}$$

Remark. In reality, U and V are

$-\infty$ and the left censoring variable, respectively, if left censoring occurs;
the right censoring variable and ∞ , respectively, if right censoring occurs,
the two consecutive inspection times if SIC occurs.

Then assumption 2 in the model is false according to the interpretation. However, like the case 2 model for case 2 data, the model is very simple and easy to interpret for their variables.

A realistic model (MIC model (2), Yu *et al.* (2001)):

Assume:

N is a random integer;
 $T, Y_1 < Y_2 < \dots < Y_k < \dots$ are inspection times, $Y_0 = -\infty$;
 X and $(N, T, Y_1, \dots, Y_k, \dots)$ are independent;
 $P(N = 0) > 0$ and $P(N > 1) > 0$;
The observable random vector is

$$(L, R) = \begin{cases} (X, X) & \text{if } X \leq T \text{ and } N = 0 \\ (T, \infty) & \text{if } X > T \text{ and } N = 0 \\ (-\infty, Y_1) & \text{if } X \leq Y_1 \text{ and } N \geq 1 \\ (Y_i, Y_{i+1}) & \text{if } Y_i < X \leq Y_{i+1}, i = 1, \dots, N - 1 \text{ and } N \geq 1 \\ (Y_N, \infty) & \text{if } X > Y_N \text{ and } N \geq 1 \end{cases}$$

The model can be viewed as a mixture of a RC model and a mixed case interval censorship model. That is,

$$F_{L,R}(l, r) = \sum_{k \geq 0} F_{L,R|N}(l, r|k) f_N(k).$$

Other models

Petroni and Wolfe (1994) and Huang (1999) construct two different models for the mixed IC data. Huang's model is basically a mixture of an uncensored model and a case k model, and thus is a special case of our MIC model (2) with $P(N = i) = 0$ for $i \neq 0$ or k and with $T \equiv \infty$. The formulation of Petroni and Wolfe's model is basically the model described in §1.5.3 with the additional assumption that X is discrete as well. Thus it limits its extension to the continuous cases. Huang's model requires that X may be observed in the whole range of X , which is often not the case in reality.

§1.7. Left censoring.

Easy.

§1.8.

Table 1. Classification of IC data

(<i>observations :</i>	<i>LC</i>	<i>SIC</i>	<i>RC</i>	<i>exact</i>)
	<i>RC data</i>			+	+	
	<i>LC data</i>	+			+	
	<i>C1 data</i>	+		+		
	<i>DC data</i>	+		+	+	
	<i>C2 data</i>	+	+	+		
	<i>MIC data</i>	+	+	+	+	

§1.9. Homework: Generate a set of C2 data under the mixed case interval censorship model with a size of 100 and $P(N = i) > 0, i = 1, \dots, 8$. What will you do if you want to estimate the $F(x)$? (For example, one may consider the following treatment: Let

$X_i^* = \begin{cases} \frac{L_i + R_i}{2} & \text{if SIC,} \\ L_i & \text{if RC,} \\ R_i & \text{if LC,} \end{cases}$ then pretend that X_i is observed and its value is X_i^* . Finally,

estimate $F(t)$ by

$$\tilde{F}(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(X_i^* \leq t)}$$

What do you expect in terms of the asymptotic property (consistency) of your or the above estimator? Use the simulated data to compute the above estimate and compare to $F(t)$

(repeat 10 times) and the limiting value of the estimator (you can select only one specific t).

BIBLIOGRAPHY

- * Ayer, M., Brunk, H.D., Ewing, G.M., Reid, W.T. and Silverman, E. (1955). An empirical distribution function for sampling incomplete information. *Ann. Math. Statist.* 26, 641-647.
- * Becker, N. and Melbye, M. (1991). Use of a log-linear model to compute the empirical survival curves from interval-censored data, with application to data on tests for HIV positivity. *Austral. J. Statistics*, 33, 125-133.
- * Cox, D.R. and Oakes, D. (1984). *Analysis of Survival Data*. Chapman & Hall NY, 70-71.
- * Finkelstein, D.M. and Wolfe, R.A. (1985). A semi-parametric model for regression analysis of interval-censored failure time data. *Biometrics*, 41, 933-945.
- * Gehan, E. A. (1965). A generalized Wilcoxon test for comparing arbitrarily single-censored samples. *Biometrika*, 52, 203-23.
- * Groeneboom, P. and Wellner, J. A. (1992). *Information bounds and nonparametric maximum likelihood estimation*. Birkhäuser Verlag, Basel.
- * Huang, J. (1999). Asymptotic properties of nonparametric estimation based on partly interval-censored data. *Statistica Sinica*, 9, 501-520.
- * Huang, J.L., Lee, C.S. and Yu, Q.Q. (2007). A generalized log-rank test for interval-censored failure time data via multiple imputation. *Statistics in Medicine*, (accepted).
- * Kaplan, E. L. and Meier, P. (1958). Nonparametric estimation from incomplete observations. *J. Amer. Stat. Assoc.*, 53, 457-481.
- * Keiding, N. (1991) Age-specific incidence and prevalence: A statistical perspective (with discussion) *J. Roy. Statist. Soc. Ser. A*, 154, 371-412.
- * Leiderman, P.H., Babu, D., Kagia, J., Kraemer, H.C. and Leiderman, G.F. (1973). African infant precocity and some social influences during the first year. *Nature*, 242, 247-249.
- * Miller Jr., R. G. (1981). *Survival analysis*. Wiley NY. Odell, P.M., Anderson, K.M. and D'Agostino, R.B. (1992). Maximum likelihood estimation for interval-censored data using a Weibull-based accelerated failure time model. *Biometrics*, 48, 951-959.
- * Petroni, G. R. and Wolfe, R. A. (1994). A two-sample test for stochastic ordering with interval-censored data. *Biometrics*, 50, 77-87.
- * Schick, A. and Yu, Q.Q. (2000). Consistency of the GMLE with mixed case interval-censored data. *Scand. J. Statist.*, 27, 45-55.
- * Sun, Jianguo (2006). *The statistical analysis of interval-censored failure time data* Springer.
- * Wang, Z. and Gardiner, J. C. (1996). A class of estimators of the survival function from interval-censored data. *Ann. Statist.*, 24, 647-658.
- * Yu, Q.Q., Li, L.X. and Wong, G.Y.C. (2000). On consistency of the self-consistent estimator of survival functions with interval censored data. *Scan. J. of Statist.* Vol 27 35-44.
- * Yu, Q. Q., Wong, G. Y. C. and Li, L. X. (2001). Asymptotic properties of self-consistent estimators with mixed interval-censored data. *Ann. Inst. Statist. Math.* 53 469-486.

Chapter 2. Distribution of failure time

§2.1. Hazard.

Suppose that X is r.v. with cdf F and density function f .

Definition. $S(t) = P(X > t)$ is often called the *survival function* of a r.v. X .

$S(t-) = P(X \geq t)$ is sometimes called the *survival function* of X (Cox and Oakes (1984)).

If X is continuous, $S(t-) = S(t)$.

In general, $S(t-) = \lim_{u \uparrow t} S(u)$.

$$S(x) = 1 - F(x).$$

The d.f. of X ,

$$f(t) = \begin{cases} S(t-) - S(t) & \text{if } X \text{ is discrete (why ?)} \\ -S'(t) & \text{if } X \text{ is continuous (why ?)} \end{cases}$$

$$S(t) = \begin{cases} \sum_{x>t} f(x) & \text{if } X \text{ is discrete} \\ \int_t^\infty f(x)dx & \text{if } X \text{ is continuous} \end{cases}$$

Definition. $h(t) = f(t)/S(t-)$ is called the *hazard function* of X .

$H(t) = -\log S(t)$ is called the *integrated (or cumulative) hazard* of X .

Interpretation of h and H :

$$h(x) = \begin{cases} P(X = x | X \geq x) & \text{if } X \text{ is discrete (why ?)} \\ \lim_{\Delta \rightarrow 0} \frac{P(X \in [x, x+\Delta) | X \geq x)}{\Delta} & \text{if } X \text{ is continuous (why ?)} \end{cases}$$

If X is continuous,

$$\begin{aligned} H'(t) &= (-\log S(t))' \\ &= f(t)/S(t) \\ &= h(t) \text{ if } f(t) \text{ exists.} \end{aligned}$$

That is

$$H(t) = \int_{-\infty}^t h(x)dx \quad \text{integrated hazard.}$$

If X is discrete,

$$\begin{aligned} S(t) &= \prod_{x_i \leq t, x_i \in D_f} \frac{S(x_i)}{S(x_i-)} \\ &= \prod_{x_i \leq t} (1 - h(x_i)) \end{aligned}$$

where $D_f = \{x : f(x) > 0\}$ and $\prod_{a \in \emptyset} a = 1$. Thus

$$H(t) = -\log \prod_{x_i \leq t} (1 - h(x_i)) \quad (\approx \sum_{x_i \leq t} h(x_i) \text{ if } h(x_i) \approx 0 \text{ for all } i) \text{ (cumulated hazards).}$$

§2.2. Some distributions

For the purpose of parametric analysis later on, as well as for possible simulation studies, we need to be familiar with certain distributions which are related to the survival analysis. Thus we study some typical distributions here. Note that they all correspond to some positive random variables. We thus use the following conventions.

$f(x)$ ($F(x)$, $h(x)$, or $H(x)$), $x \geq 0$ implies that $f(x) = 0$ etc. for $x < 0$;

$S(x)$, $x \geq 0$ implies that $S(x) = 1$ for $x < 0$.

§2.2.1. Exponential distribution

$f(x) = \rho e^{-\rho x}$, $x > 0$, $\rho > 0$;

$S(x) = e^{-\rho x}$, $x > 0$; (**easy way to remember**)

$h(x) = \frac{f(x)}{S(x)} = \rho$, $x > 0$ (constant hazard);

$H(x) = \int_0^x h(t)dt = \rho x$, $x > 0$ (or $= -\ln S(x)$);

$E(X) = 1/\rho$.

Note that $E(X) = \int_0^\infty S(t)dt$.

Reason: If X is a continuous survival time, then

$$E(X) = \int_0^\infty x f(x)dx = \int_0^\infty \int_0^x f(x)dt dx = \int_0^\infty \int_t^\infty f(x)dx dt = \int_0^\infty S(t)dt.$$

In general, if X is a survival time, then $E(X) =$

$$E\left(\int_0^X 1 dt\right) = E\left(\int_0^\infty \mathbf{1}_{(t \leq X)} dt\right) = \int_0^\infty E(\mathbf{1}_{(t \leq X)}) dt = \int_0^\infty P(X \geq t) dt = \int_0^\infty S(t-) dt$$

Are these two equations equivalent ?

If X is not a survival time, do we have $E(X) = \int_0^\infty S(t-)dt$?

$$\int_0^x 1dx = \int_0^\infty 1(t \leq x)dt ?$$

$$\begin{aligned} E(X) &= E(\int_0^X 1dt) = E(\int_0^X 1(X \geq 0)dt) + E(\int_0^X 1(X < 0)dt) \\ &= E(\int_0^X 1(X \geq 0)dt) - E(\int_X^0 1(X < 0)dt) \\ &= E(\int_0^\infty \mathbf{1}_{(t \leq X)} dt) - E(\int_{-\infty}^0 \mathbf{1}_{(t \geq X)} dt) \\ &= \int_0^\infty E(\mathbf{1}_{(t \leq X)}) dt - \int_{-\infty}^0 E(\mathbf{1}_{(t \geq X)}) dt \\ &= \int_0^\infty S(t-)dt - \int_{-\infty}^0 F(t)dt \end{aligned}$$

§2.2.2. Weibull distribution

$$\begin{aligned} f(x) &= \kappa \rho (\rho x)^{\kappa-1} e^{-(\rho x)^\kappa}, \quad x > 0, \rho, \kappa > 0; \\ S(x) &= e^{-(\rho x)^\kappa}, \quad x > 0 \text{ (easy way to remember);} \\ h(x) &= \kappa \rho (\rho x)^{\kappa-1}, \quad x > 0; \\ H(x) &= (\rho x)^\kappa, \quad x > 0; \end{aligned}$$

§2.2.3 Log normal distribution

Let $Z \sim N(\log \rho^{-1}, \tau^2)$, $\rho, \tau > 0$, and $X = e^Z$, then X has a log normal distribution.

$$f(x) = f_Z(\log x) \left| \frac{dz}{dx} \right| = \frac{1}{\sqrt{2\pi}\tau x} \exp\left(-\frac{(\log x - \log \rho^{-1})^2}{2\tau^2}\right) = \frac{1}{\sqrt{2\pi}\tau x} \exp\left(-\frac{(\log(x\rho))^2}{2\tau^2}\right),$$

$x > 0$. S , h and H cannot be simplified.

§2.2.4 Log-logistic distribution

If $X = e^T$ and T has logistic distribution, then X has the log logistic distribution.

A logistic distribution has a d.f

$$f_T(x) = \frac{\tau^{-1} \exp(\frac{x-\nu}{\tau})}{(1 + \exp(\frac{x-\nu}{\tau}))^2}, \quad \tau > 0.$$

$$f(x) = f_T(\log x) \left| \frac{dt}{dx} \right| = \frac{\tau^{-1} \exp(\frac{\log x - \nu}{\tau})}{(1 + \exp(\frac{\log x - \nu}{\tau}))^2 x}, \quad x \geq 0;$$

Let $\nu = -\log \rho$ and $\kappa = 1/\tau$ (reparameterization),

$$\begin{aligned} f(x) &= \frac{\kappa \exp(\log(\rho x)^\kappa)}{[1 + \exp(\log(\rho x)^\kappa)]^2 x} \\ &= \frac{\kappa (\rho x)^\kappa}{[1 + (\rho x)^\kappa]^2 x} \\ &= \frac{\kappa \rho^\kappa x^{\kappa-1}}{[1 + (\rho x)^\kappa]^2}, \quad x \geq 0. \end{aligned}$$

It yields S . **How ?**

Another smart way to derive S .

A logistic survival function is

$$S_T(x) = \frac{1}{1 + \exp(\frac{x-\nu}{\tau})}.$$

If $X = e^T$, then X has the log logistic distribution, with

$$S(x) = S_X(x) = P\{X > x\} = P\{T > \ln x\} = \frac{1}{1 + \exp(\frac{\ln x - \nu}{\tau})}.$$

Let $\nu = -\log \rho$ and $\kappa = 1/\tau$ (reparameterization),

then

$$S(x) = \frac{1}{1 + (\rho x)^\kappa}, \quad x \geq 0 \text{ (easy way to remember).}$$

$$H(x) = \log(1 + (\rho x)^\kappa), \quad x \geq 0.$$

$$h(x) = \frac{\kappa \rho^\kappa x^{\kappa-1}}{1 + (\rho x)^\kappa}, \quad x \geq 0.$$

§2.2.5. Gompertz-Makeham distribution.

$$\begin{aligned} h(t) &= \rho_0 + \rho_1 e^{\rho_2 t}, \quad t > 0, \quad \rho_0, \rho_1 > 0 \text{ (easy way to remember)}, \\ H(t) &= \int_0^t h(x) dx = [\rho_0 t + \frac{\rho_1}{\rho_2} (e^{\rho_2 t} - 1)], \quad t > 0; \\ S(t) &= e^{-H(t)} = e^{-[\rho_0 t + \frac{\rho_1}{\rho_2} (e^{\rho_2 t} - 1)]}, \quad t > 0; \\ f(t) &= -S'(t) = h(t)S(t), \quad t > 0; \end{aligned}$$

Remark:

1. It is an exponential distribution if $\rho_1 = 0$.
2. It is called Gompertz distribution if $\rho_0 = 0$.

§2.2.6. Compound exponential distribution.

Suppose that each individual survival time is exponentially distributed but that the rate varies randomly between individuals.

To represent this let P be a random variable with density f_P and the conditional density of X given $P = p$ is

$$f_{X|P}(x|p) = p e^{-px}, \quad x > 0.$$

Then $f(x) = \int p e^{-px} f_P(p) dp$. If P has gamma distribution, say,

$$f_P(p) = \frac{p^{\alpha-1} e^{-p/\beta}}{\Gamma(\alpha) \beta^\alpha}, \quad p, \alpha, \beta > 0.$$

Letting $\rho = \alpha\beta$, (the mean), then

$$\begin{aligned} f(x) &= \int_0^\infty p e^{-px} \frac{p^{\alpha-1} e^{-p/\rho}}{\Gamma(\alpha) (\rho/\alpha)^\alpha} dp = \int_0^\infty \frac{p^\alpha e^{-p(x+\alpha/\rho)}}{\Gamma(\alpha) (\rho/\alpha)^\alpha} dp \\ &= \frac{\Gamma(\alpha+1)}{(x + \alpha/\rho)^{\alpha+1} \Gamma(\alpha) (\rho/\alpha)^\alpha} \\ &= \alpha (\alpha/\rho)^\alpha (x + \alpha/\rho)^{-\alpha-1}. \end{aligned}$$

The latter df is called the **Pareto distribution**, with

$$\begin{aligned} S(x) &= (\alpha/\rho)^\alpha (x + \alpha/\rho)^{-\alpha}. \\ H(x) &= \alpha (-\ln(\alpha/\rho) + \ln(x + \alpha/\rho)). \\ h(x) &= \alpha / (x + \alpha/\rho) \text{ (easy way to remember)}. \end{aligned}$$

§2.2.7. Discrete distributions.

The common discrete random variables do not have concise forms for h and H . Thus so far, we only consider continuous r.v.s. Now consider a binomial distribution. $X \sim \text{bin}(2, p)$.

$$\begin{aligned} f(x) &= \binom{2}{x} p^x q^{2-x}, \quad x \in \{0, 1, 2\}, \quad q = 1 - p. \\ S(x) &= \begin{cases} 1 & \text{if } x < 0 \\ 1 - q^2 & \text{if } x \in [0, 1), \\ p^2 & \text{if } x \in [1, 2), \\ 0 & \text{if } x \geq 2. \end{cases} \\ S(0-) &= 1, \quad S(1-) = 1 - (1 - p)^2, \quad S(2-) = p^2, \end{aligned}$$

$$h(x) = \begin{cases} q^2 & \text{if } x = 0, \\ \frac{2pq}{1-q^2} & \text{if } x = 1, \\ p^2/p^2 & \text{if } x = 2, \end{cases} = \begin{cases} q^2 & \text{if } x = 0, \\ \frac{2q}{2-p} & \text{if } x = 1, \\ 1 & \text{if } x = 2, \end{cases} \quad H(x) = \begin{cases} 0 & \text{if } x < 0 \\ -\ln(1 - q^2) & \text{if } x \in [0, 1), \\ -2\ln p & \text{if } x \in [1, 2), \\ \infty & \text{if } x \geq 2. \end{cases}$$

§2.2.8. Proportional hazards (PH) model.

An advantage of defining hazard functions is the introduction of the PH model. Define

$$\tau = \tau_T = \sup\{t : F_T(t) < 1\}$$

for a random variable T .

Definition. Let (\mathbf{X}, Y) be a random vector, where $\mathbf{X} \in \mathcal{R}^p$. Given $\mathbf{X} = \mathbf{x}$, the hazard of $Y|\mathbf{x}$ is

$$h(y|\mathbf{x}) = h_o(y)c(\mathbf{x}), \text{ for } y < \tau, \quad (1)$$

where $c(\mathbf{x}) \geq 0$, $c(\cdot)$ takes at least two distinct values, and h_o is a hazard function. Then we say (\mathbf{X}, Y) follows a proportional hazards (PH) model or Cox's regression model (Cox, 1972).

Remark 1. It is common to set $c(\mathbf{x}) = \exp(\beta\mathbf{x})$ so that $c(\mathbf{x}) \geq 0$, where $\beta\mathbf{x} = \beta'\mathbf{x}$.

Remark 2. If the random variable is discrete, then the choice of $c(\mathbf{x}) = e^{\beta\mathbf{x}}$ may cause problem. For instance, if $h_o(y) = 0.5$, $c(\mathbf{x}) = 3$, $P(Y = y|Y \geq y, \mathbf{X} = \mathbf{x}) = h(y|\mathbf{x}) = h_o(y)c(\mathbf{x}) > 1$.

Two alternatives:

- (1) either choose $c(\mathbf{x}) = \exp(-e^{\beta\mathbf{x}})$ to ensure that $c(\mathbf{x})$ is between 0 and 1,
- (2) or restrict the parameter space \mathcal{B} , the set that β belongs to.

For continuous random variable, we only need $c(\mathbf{x}) \geq 0$, as the hazard does not need to belong to $[0, 1]$, as long as it belongs to $[0, \infty]$.

Remark 3. In the original definition of the PH model (see Cox and Oakes (1984)), there is not a restriction $y < \tau$. We shall show in Example 1 that if h_o corresponds to a discrete random variable, statement (1) without the restriction does not define a hazard function.

Hereafter, take $p = 1$.

Example 1. (Counterexample of Eq. (1) without $y < \tau$). We shall consider an example of discrete random variables. If T is discrete and $P\{T = \tau\} = f_T(\tau) > 0$, then

$$h_T(\tau) = f_T(\tau)/S_T(\tau-) = f_T(\tau)/f_T(\tau) = 1$$

which is always true. It follows that for such a discrete random variable statement (1) does not hold at τ , as

$$h(\tau|x) = 1 \neq 1 \times c(x) = h_o(\tau)c(x) \text{ as } c(x) \neq 1 \text{ for some } x.$$

It does not matter for continuous random variables, as one can eliminate τ from the support.

Example 2. If $S_o(t)$ is a survival function of a continuous random variable, then

$$S(t|x) = (S_o(t))^{e^{\beta x}} \text{ satisfies the PH model.} \quad (2)$$

Reason: $f(t|x) = -S'(t|\mathbf{x}) = -e^{\beta x}(S_o(t))^{e^{\beta x}-1}S_o'(t) = e^{\beta x}S(t|x)\frac{f_o(t)}{S_o(t)}$,

$$h(t|x) = \frac{f(t|x)}{S(t|x)} = e^{\beta x}h_o(t).$$

Special cases:

- a. Weibull: $S_o(t) = e^{-t^\gamma}$, $t > 0$. $S(t|x) = \exp(-e^{\beta x}t^\gamma)$, $t > 0$.
 $h(t|x) = e^{\beta x}\gamma t^{\gamma-1}$, $h_o(t) = \gamma t^{\gamma-1}$, $t > 0$.
- b. Log-logistic: $S_o(t) = \frac{1}{1+t^\kappa}$. $S(t|x) = (\frac{1}{1+t^\kappa})^{e^{\beta x}}$, $t > 0$,
 $h(t|x) = \exp(\beta x)\frac{\kappa t^{\kappa-1}}{1+t^\kappa}$. $h_o(t) = \frac{\kappa t^{\kappa-1}}{1+t^\kappa}$, $t > 0$.
- c. Logistic: $S_o(t) = \frac{1}{1+e^t}$. $S(t|x) = (\frac{1}{1+e^t})^{e^{\beta x}}$,
 $h(t|x) = \exp(\beta x)\frac{1}{1+e^{-t}}$, $h_o(t) = \frac{1}{1+e^{-t}}$.

Remark.

The distribution generated from Case a is still a Weibull distribution for each (β, x) .

The distribution generated from Case b may not be a log-logistic distribution unless $\beta = 0$.

The distribution in Case c corresponds to a random variable Y with negative values in

its domain, though it is often in survival analysis that we only consider the non-negative domain. However, it still satisfies the PH model.

Definition. The family of the survival functions $S(t|\mathbf{x})$ satisfies statement (2) for all possible β is called a Lehmann family, or we can say that the distribution is from a proportional integrated hazards (PIH) model, as

$$H(t|x) = -\ln S(t|x) = e^{\beta x}(-\ln S_o(t)) = e^{\beta x} H_o(t).$$

Statement (2) does not hold for discrete random variable. When Cox proposes the PH model, he distinguishes the model from the Lehmann family or the PIH model. However, later in the literature, the PIH model and the PH model are mistaken to be the same (see Sun (2006) p.18). Example 3 shows that they are different.

Example 3. Suppose $Y_0 \sim \text{bin}(2, p)$. Then its hazard function is

$$h_o(t) = \begin{cases} (1-p)^2 & \text{if } t = 0, \\ \frac{2(1-p)}{2-p} & \text{if } t = 1, \\ 1 & \text{if } t = 2. \end{cases}$$

Suppose $h(y|x) = h_o(y)c(x)$ for $y = 0$ or 1 .

Then

$$\begin{aligned} h(0|x) &= (1-p)^2 c(x) \text{ yields } f(0|x) = (1-p)^2 c(x) \text{ as } S(0-|x) = 1. \\ h(1|x) &= \frac{2(1-p)}{2-p} c(x) \text{ yields } f(1|x) = \frac{2(1-p)}{2-p} c(x)(1 - (1-p)^2 c(x)). \\ & f(2|x) = 1 - f(0|x) - f(1|x). \end{aligned}$$

Verify that $f(\cdot|x)$ defines a discrete density function (for $c(x) \leq 1$. **Why add it ??**)

It follows that

$$S(0|x) = 1 - (1-p)^2 c(x).$$

However, if $p = 0.2$ and $c(x) = 0.3$,

$$(S_o(0))^{c(x)} = (1 - (1-p)^2)^{c(x)} \approx 0.7360219 \neq 0.808 \approx 1 - (1-p)^2 c(x) = S(0|x).$$

It indicates that if h_o is a hazard function of a discrete random variable, and $h(t|x) = h_o(t)c(x)$, its cdf may not be of the form $S(t|x) = (S_o(t))^{c(x)}$ or Equation (2).

§2.2.9. Accelerated lifetime (AL) model.

Definition. If S_o is a survival function and given a p dimensional vector $\mathbf{X} = \mathbf{x}$, $Y|\mathbf{x}$ has a survival function $S(y|\mathbf{x}) = S_o(y/\exp(\beta\mathbf{x}))$, $\beta \in \mathcal{R}^p$, then we say $Y|\mathbf{x}$ is from an accelerated lifetime model.

If $Y|\mathbf{x}$ is from an AL model, then $\ln Y = \beta\mathbf{x} + \epsilon$, where e^ϵ has the survival function S_o .

Examples.

Weibull: $S(y|\mathbf{x}) = \exp(-y^\kappa e^{\alpha\mathbf{x}})$, $y > 0$, $\alpha = -\beta\kappa$. ($S_o(y) = \exp(-(\rho y)^\gamma)$).

Log-logistic: $S(y|\mathbf{x}) = \frac{1}{1+(y \exp(\alpha\mathbf{x}))^\kappa}$. $\alpha = -\beta$. ($S_o(y) = \frac{1}{1+(\rho y)^\kappa}$).

§2.3. Homework:

- A.1. Let X_1, \dots, X_n be independent continuous nonnegative random variables with hazard functions $h_1(\cdot), \dots, h_n(\cdot)$. Prove that $X = \min\{X_1, \dots, X_n\}$ has hazard function $\sum_{j=1}^n h_j(t)$.
- A.2. In a compound exponential distribution, let the rate be represented by the random variable P . Prove that $E(X) = E(1/P)$ and $\text{Var}(X) = 2E(1/P^2) - [E(1/P)]^2$.
- A.3. Derive the integrated hazard function $H(t|x)$ under the PH model in Example 3 of §2.2.8.
- A.4. Show that the two examples in §2.2.9 are indeed from the AL model. Moreover, one is from the PH model and the other is not.

References.

- * Cox, D.R. (1972). Regression models and life tables (with discussion). *Journal of the Royal Statistical Society*, Series B 34, 187-220.
- * Yu, Q.Q., Wong, G.Y.C, and Ye, L. (1995). Estimation of a survival function with interval-censored data. A simulation study on the redistribution-to-the-inside estimator. *1995 Program of the joint statistical meetings*.

Chapter 3. Parametric Analysis

§3.1. Introduction

Assume that

the failure time X has a df $f = f_o(x; \theta)$,
 where f_o is known, θ is a parameter, $\theta \in \Theta$ (parameter space),
 the censoring vector has a df $g(\cdot)$, which does not depend on θ ,
 (L, R) is an extended observable random vector,
 (L_i, R_i, X_i) , $i = 1, \dots, n$, are i.i.d. copies of (L, R, X) .

Note that $L \leq X \leq R$ and $L_i \leq X_i \leq R_i$,

$$(L, R) \text{ is } \begin{cases} RC & \text{if } R = \infty \\ LC & \text{if } L = -\infty \\ SIC & \text{if } 0 < L < R < \infty \\ \text{exact} & \text{if } L = R. \end{cases}$$

Recall that if we observed $X_1 = x_1, \dots, X_n = x_n$,
 the likelihood of a complete data is

$$l(\phi) = \prod_{i=1}^n f_o(x_i; \phi).$$

Note: θ is a parameter and ϕ is a variable in Θ , as we do not know θ .

In particular,

$$l(\theta) \begin{cases} = \prod_{i=1}^n P(X_i = x_i) & \text{if } X \text{ is discrete} \\ \approx \prod_{i=1}^n \frac{P\{X_i \in [x_i, x_i + \Delta)\}}{\Delta} & \text{if } X \text{ is continuous, where } \Delta \approx 0, \end{cases}$$

as $P\{X \in [x, x + \Delta)\} \approx f(x)\Delta$. The MLE of θ maximizes $l(\phi)$, $\phi \in \Theta$. We shall modify this idea for IC data.

§3.2. Likelihood functions.

If (L, R) is discrete, extending the definition for complete data, we could call $l(\theta)$ the likelihood function of the IC data $(L_i, R_i) = (l_i, r_i)$, $i = 1, \dots, n$, where

$$l(\theta) = \prod_{i=1}^n P((L_i, R_i) = (l_i, r_i))$$

Example 1 (RC model). Denote f_Y , F_Y and S_Y the df, cdf and survival function of Y , respectively.

$$P((L, R) = (l, r)) = \begin{cases} P(X = l \leq Y) & \text{if exact} \\ P(Y = l < X) & \text{if RC.} \end{cases} = \begin{cases} f(l)S_Y(l-) & \text{if exact} \\ S(l)f_Y(l) & \text{if RC.} \end{cases}$$

Let $\mathbf{1}_e$, $\mathbf{1}_r$, $\mathbf{1}_l$ and $\mathbf{1}_s$ be the indicator functions of the events that the observation is exact, RC, LC and SIC, respectively. Denote $\mathbf{1}_{e,i}$, $\mathbf{1}_{r,i}$, $\mathbf{1}_{l,i}$ and $\mathbf{1}_{s,i}$ in an obvious way.

$$P((L, R) = (l, r)) = [f(l)S_Y(l-)]^{\mathbf{1}_e} [S(l)f_Y(l)]^{\mathbf{1}_r}.$$

Note that we ignored the variable θ in the above expression for simplicity,
 as $f = f_o(x; \theta)$, $S = S_o(x; \theta)$ and $F = F_o(x; \theta)$.

Thus, ignoring ϕ , we have

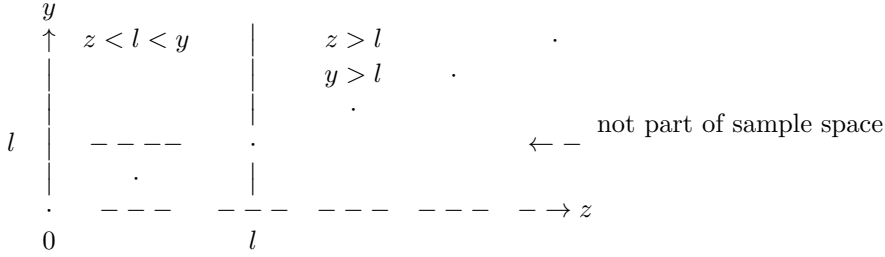
$$\begin{aligned} l(\phi) &= \prod_{i=1}^n \{ [f(l_i)S_Y(l_i-)]^{\mathbf{1}_{e,i}} [S(l_i)f_Y(l_i)]^{\mathbf{1}_{r,i}} \} \\ &= \left\{ \prod_{i=1}^n \{ [f(l_i)]^{\mathbf{1}_{e,i}} [S(l_i)]^{\mathbf{1}_{r,i}} \} \right\} \left\{ \prod_{i=1}^n \{ [(S_Y(l_i-))]^{\mathbf{1}_{e,i}} [f_Y(l_i)]^{\mathbf{1}_{r,i}} \} \right\}. \end{aligned}$$

Example 2 (DC model).

$$P((L, R) = (l, r)) = \begin{cases} P(X = l, Z < l \leq Y) & \text{if exact} \\ P(X > l = Y) & \text{if RC} \\ P(X \leq r = Z) & \text{if LC} \end{cases} = \begin{cases} f(l)P(Z < l \leq Y) & \text{if exact} \\ S(l)f_Y(l) & \text{if RC} \\ F(r)f_Z(r) & \text{if LC.} \end{cases}$$

Note $\{Z \geq l\} \subset \{Y \geq l\}$ as $Z \leq Y$ w.p.1.

Thus $P(Z < l \leq Y) = P(Y \geq l) - P(Z \geq l) = S_Y(l-) - S_Z(l-)$,



$$P((L, R) = (l, r)) = [f(l)(S_Y(l-) - S_Z(l-))] \mathbf{1}_e [S(l)f_Y(l)] \mathbf{1}_r [F(r)f_Z(r)] \mathbf{1}_l.$$

Thus $l(\phi)$

$$\begin{aligned} &= \prod_{i=1}^n \{ [f(l_i)(S_Y(l_i-) - S_Z(l_i-))] \mathbf{1}_{e,i} [S(l_i)f_Y(l_i)] \mathbf{1}_{r,i} [F(r_i)f_Z(r_i)] \mathbf{1}_{l,i} \} \\ &= \left\{ \prod_{i=1}^n ([f(l_i)] \mathbf{1}_{e,i} [S(l_i)] \mathbf{1}_{r,i} [F(r_i)] \mathbf{1}_{l,i}) \right\} \prod_{i=1}^n ([S_Y(l_i-) - S_Z(l_i-)] \mathbf{1}_{e,i} [f_Y(l_i)] \mathbf{1}_{r,i} [f_Z(r_i)] \mathbf{1}_{l,i}). \end{aligned}$$

Example 3 (Case k model, $k \geq 2$).

$$\begin{aligned} P((L, R) = (l, r)) &= \begin{cases} P\{l < X \leq r, (Y_{j-1}, Y_j) = (l, r), j \in \{2, \dots, k\}\} & \text{if SIC} \\ S(l)f_{Y_k}(l) & \text{if RC} \\ F(r)f_{Y_1}(r) & \text{if LC} \end{cases} \\ &= \begin{cases} (S(l) - S(r)) \sum_{j=2}^k f_{Y_{j-1}, Y_j}(l, r) & \text{if it is SIC} \\ S(l)f_{Y_k}(l) & \text{if it is right censored} \\ F(r)f_{Y_1}(r) & \text{if it is left censored} \end{cases} \end{aligned}$$

$$P((L, R) = (l, r)) = [(S(l) - S(r)) \sum_{j=2}^k f_{Y_{j-1}, Y_j}(l, r)] \mathbf{1}_s [S(l)f_{Y_k}(l)] \mathbf{1}_r [F(r)f_{Y_1}(r)] \mathbf{1}_l.$$

Thus $l(\phi)$

$$\begin{aligned} &= \prod_{i=1}^n \{ [(S(l_i) - S(r_i)) \sum_{j=2}^k f_{Y_{j-1}, Y_j}(l_i, r_i)] \mathbf{1}_{s,i} [S(l_i)f_{Y_k}(l_i)] \mathbf{1}_{r,i} [F(r_i)f_{Y_1}(r_i)] \mathbf{1}_{l,i} \} \\ &= \left\{ \prod_{i=1}^n [(S(l_i) - S(r_i)) \mathbf{1}_{s,i} [S(l_i)] \mathbf{1}_{r,i} [F(r_i)] \mathbf{1}_{l,i}] \right\} \\ &\quad \times \left\{ \prod_{i=1}^n \left(\sum_{j=2}^k f_{Y_{j-1}, Y_j}(l_i, r_i) \right) \mathbf{1}_{s,i} [f_{Y_k}(l_i)] \mathbf{1}_{r,i} [f_{Y_1}(r_i)] \mathbf{1}_{l,i} \right\}. \end{aligned}$$

Since the effect of the censoring vector can be factored out separately, and only the first factor in $l(\phi)$ depends on ϕ we can discard the second factor.

Definition. The likelihood function of the IC data is defined to be

$$\mathbf{L}(\phi) = \left\{ \prod_{i=1}^n \{ [f(l_i)]^{\mathbf{1}_{e,i}} [S(l_i) - S(r_i)]^{\mathbf{1}_{s,i}} [S(l_i)]^{\mathbf{1}_{r,i}} [F(r_i)]^{\mathbf{1}_{l,i}} \}, \right.$$

where $f(x) = f_o(x; \phi)$ in the parametric analysis.

Note that while we start our discussion under the assumption that (L, R) is discrete, the definition of \mathbf{L} does not require that the (L, R) be discrete.

In an obvious way, we write

$$\mathbf{L}(\phi) = \prod_{i: ex} f(l_i) \prod_{i: rc} S(l_i) \prod_{i: lc} F(r_i) \prod_{i: ic} (S(l_i) - S(r_i)) = \prod_{i=1}^n [(f(l_i))^{\delta_i} (S(l_i) - S(r_i))^{1-\delta_i}],$$

where $\delta_i = \mathbf{1}(l_i = r_i)$ and $f(x) = f_o(x; \phi)$ in the parametric analysis.

§3.2.2. Homework

1. Mimic examples 1-3 for the mixed case IC model and MIC model (1).
2. If one generates data by $(L_i, R_i) = (X_i - 1, X_i + 1)$, can the likelihood be written as

$$\mathbf{L} = \prod_{i=1}^n (S(x_i - 1) - S(x_i + 1)) ?$$

§3.3. MLE.

Definition. Suppose $\phi = \hat{\theta}$ maximizes $\mathbf{L}(\phi)$, $\phi \in \Theta$. Then $\hat{\theta}$ is called the maximum likelihood estimator (MLE) of θ .

Example. Suppose that X has an exponential distribution, i.e., $f(x) = \rho e^{-\rho x}$, $x > 0$. Find the MLE of ρ under the RC model.

Solution: Observe $(Z_1, \delta_1), \dots, (Z_n, \delta_n)$, where $Z_i = \min\{X_i, Y_i\}$ and $\delta_i = \mathbf{1}_{(X_i \leq Y_i)}$ for all i

$$\begin{aligned} \log \mathbf{L}(\rho) &= \sum_{i=1}^n \log \{ [\rho e^{-\rho Z_i}]^{\delta_i} [e^{-\rho Z_i}]^{1-\delta_i} \} \quad (= \log(\prod_{i: ex} f(l_i) \prod_{i: rc} S(l_i))) \\ &= \sum_{i=1}^n \log \{ \rho^{\delta_i} e^{-\rho Z_i} \} \\ &= \sum_{i=1}^n \delta_i \log \rho - \rho \sum_{i=1}^n Z_i. \end{aligned}$$

Taking derivative and letting it equal 0 yield

$$\hat{\rho} = \frac{\sum_{i=1}^n \delta_i}{\sum_{i=1}^n Z_i}.$$

Since $\frac{\partial^2 \mathcal{L}}{\partial \rho^2} = -\sum_{i=1}^n \delta_i \rho^{-2} < 0$, $\hat{\rho}$ is the MLE of ρ , unless $\sum_{i=1}^n \delta_i = 0$, as $\rho > 0$. In the latter case, $\rho = 0$ uniquely maximizes \mathbf{L} , but it is not the MLE, as $0 \notin (0, \infty)$. We define $\hat{\rho} = 0$ so that $\hat{\rho}$ is well defined.

Remark. Formally, we shall write $\mathbf{L}(p)$ instead of $\mathbf{L}(\rho)$. However, in deriving the MLE, it does not matter.

Recall that if we have complete data, i.e., $\delta_i \equiv 1$, the MLE is $\frac{n}{\sum_{i=1}^n X_i}$. Let $h(x) = 1/x$, then $\frac{n}{\sum_{i=1}^n X_i} = h(\bar{X})$ is strongly consistent, as h is continuous. Furthermore, $\frac{n}{\sum_{i=1}^n X_i}$ is also asymptotically normally distributed as $\sigma_X < \infty$, h' is continuous and $h'(\rho) = -\rho^{-2} \neq 0$. Here we used the following asymptotic results:

1. By the the strong law of large number, if μ_X is finite and h is a continuous function,

$$h(\bar{X}) \text{ converges to } h(\mu_X) \text{ with probability one.}$$

2. By the central limit theorem and Corollary of the Slutsky's Theorem, if σ_X is finite, h' is continuous at μ_X and $h'(\mu_X) \neq 0$,

$$\frac{\sqrt{n}(h(\bar{X}) - h(\mu_X))}{\sigma_X |h'(\mu_X)|} \text{ converges in distribution to } N(0, 1).$$

The above statements are valid even if X is an $m \times 1$ random vector. In particular,

$$\frac{\sqrt{n}(h(\bar{X}) - h(\mu_X))}{\sqrt{\left(\frac{\partial h(\mu_X)}{\partial \mu}\right)^t \Sigma \frac{\partial h(\mu_X)}{\partial \mu}}} \text{ converges in distribution to } N(0, 1).$$

where μ is an $m \times 1$ vector, Σ is the covariance matrix of X and μ^t is the transpose of μ , provided Σ is nonsingular, $\frac{\partial h}{\partial \mu}$ is continuous at $\mu = \mu_X$ and is not a zero vector at $\mu = \mu_X$.

§3.3.2. Homework:

1. Show that $\hat{\rho}$ in the example is consistent under the RC model.
(Hint: $\sum_{i=1}^n Z_i = \sum_{i=1}^n X_i \mathbf{1}_{(X_i \leq Y_i)} + \sum_{i=1}^n Y_i \mathbf{1}_{(X_i > Y_i)}$, or derive the df of Z .)
2. Show that $\hat{\rho}$ is asymptotically normally distributed under the RC model. That is

$$\sqrt{n}(\hat{\rho} - \rho) \text{ converges in distribution to a normal variate.}$$

3. Give a 99% approximate confidence interval for ρ when $n = 100$ and give a 99% confidence interval for ρ when $n = 1$ with $\rho \geq 1/\ln 50$.
4. Show that if X has a continuous cdf F with integrated hazard $H(\cdot)$ and $Z = \min\{X, c\}$, where c is a fixed constant, then $E(H(Z)) = F(c)$.
5. In the above example, if one uses the second incorrect approach to deal with the RC data, the MLE will be $\tilde{\rho} = \frac{n}{\sum_{i=1}^n Z_i}$. Derive the limit of the estimator and show that it is not a consistent estimator of ρ .
6. Suppose that a random sample of size n is generated from a type I RC model with $X \sim U(\theta, 4)$, $\theta \in (0, 4)$ and $P(Y = 3) = 1$. Find the MLE of θ . Derive the mean and variance of the MLE. Can we use the Cramer-Rao Lower Bound as the estimator of the variance of the MLE ?

§3.4. Numerical methods for MLE.

3.4.1. Newton-Raphson method.

If $\mathbf{L}(\phi)$ is continuous and Θ is compact, then the MLE of ϕ exists. However, the MLE may not have closed form solution.

We can try the Newton Raphson method to derive it. Denote

$$\mathcal{L} = \log \mathbf{L}.$$

This is an iterative algorithm:

Step 1. Assign an initial value $\phi_{(1)}$ to ϕ .

Step $k + 1$, $k \geq 1$. Given $\phi_{(k)}$, up-date ϕ by

$$\phi_{(k+1)} = \phi_{(k)} - \left(\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t} \Big|_{\phi=\phi_{(k)}} \right)^{-1} \frac{\partial \mathcal{L}}{\partial \phi} \Big|_{\phi=\phi_{(k)}}. \quad (3.1)$$

Stop when $\|\phi_{(k+1)} - \phi_{(k)}\| < \epsilon$, where ϵ is sufficiently small and

$\|\mathbf{z}\| = \max_i |z_i|$ or $= \sqrt{\sum_i z_i^2}$.

Reason: By the first order Taylor expansion, (**under certain condition on \mathcal{L}**),

$$\frac{\partial \mathcal{L}}{\partial \phi} - \frac{\partial \mathcal{L}}{\partial \phi} \Big|_{\phi=\hat{\theta}} \approx \frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t} \Big|_{\phi=\hat{\theta}} (\phi - \hat{\theta})$$

Since $\frac{\partial \mathcal{L}}{\partial \phi} \Big|_{\phi=\hat{\theta}} = 0$ under certain regularity conditions (as $\hat{\theta}$ maximizes \mathcal{L}),

$$\frac{\partial \mathcal{L}}{\partial \phi} \approx \frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t} \Big|_{\phi=\hat{\theta}} (\phi - \hat{\theta})$$

$$\hat{\theta} \approx \phi - \left(\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t} \Big|_{\phi=\hat{\theta}} \right)^{-1} \frac{\partial \mathcal{L}}{\partial \phi}.$$

Equation (3.1) is based on the last equation with $\left(\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t} \Big|_{\phi=\hat{\theta}}\right)^{-1}$ replaced by $\left(\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t}\right)^{-1}$.

Drawbacks of the algorithm:

1. (Convergence). It may need some regularity assumptions to guarantee the convergence of the algorithm. For example, the algorithm may not converge even in the case that we generate complete data from Weibull(2,2).
2. (Uniqueness). It may converge to a local extreme point, unless $-\mathcal{L}$ is convex in ϕ .
3. (Feasibility). It is often difficult to obtain the inverse matrix $\left(\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t}\right)^{-1}$ if the dimension of ϕ is large.

Example 1 (Log normal under the mixed case IC model)

Let X_i be a survival time and $U_i = \log X_i$ be a $N(\mathbf{z}_i^t \beta, \sigma^2)$ variate, $i = 1, \dots, n$, where \mathbf{z}_i are $p \times 1$ covariate vector (non-random vectors), and are known, and β is a $p \times 1$ parameter vector. Under the mixed case IC model, X_i and thus U_i are not observed and we only observe (L_i, R_i) ($U_i \in (L_i, R_i]$) and \mathbf{z}_i , $i = 1, \dots, n$. We shall estimate β .

We consider U_i rather than X_i because X_i is nonnegative, while a normal variate can be negative. The problem arises from linear regression for complete data. If U_1, \dots, U_n are observed,

$$U_i = \ln X_i = \beta' \mathbf{z}_i + \epsilon_i, \text{ where } \epsilon_i \text{'s are i.i.d. } N(0, \sigma^2),$$

the MLE of β , which is also called the least squares estimator, is

$$\hat{\beta} = (\mathbf{Z}^t \mathbf{Z})^{-1} \mathbf{Z}^t \mathbf{U}, \text{ where } \mathbf{Z} = \begin{pmatrix} \mathbf{z}_1^t \\ \dots \\ \mathbf{z}_n^t \end{pmatrix} \text{ and } \mathbf{U} = \begin{pmatrix} U_1 \\ \dots \\ U_n \end{pmatrix}.$$

Here σ^2 does not matter, even though, it is unknown.

Under the mixed case IC model, $\theta = (\beta^t, \sigma)^t$ and the log likelihood function is

$$\begin{aligned} \mathcal{L}(\theta) &= \log \prod_{i=1}^n (S(l_i; \phi, \mathbf{z}_i) - S(r_i; \phi, \mathbf{z}_i)) \\ &= \log \prod_{i=1}^n \int_{l_i}^{r_i} \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\beta^t \mathbf{z}_i)^2}{2\sigma^2}} dx \text{ (if treating as independent r.v.'s)} \\ &= \log \prod_{i=1}^n \left(\Phi\left(\frac{r_i - \beta^t \mathbf{z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right) \right) \text{ (if treating as i.i.d.),} \end{aligned}$$

where $\Phi(t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$ is the cdf of the $N(0,1)$. Since

$(\int_{-\infty}^{b(x)} f(t) dt)'_x = f(b(x))b'(x)$, taking partial derivatives yields

$$\begin{aligned} \frac{\partial \mathcal{L}(\theta)}{\partial \sigma} &= \sum_{i=1}^n \frac{\frac{\partial}{\partial \sigma} \Phi\left(\frac{r_i - \beta^t \mathbf{z}_i}{\sigma}\right) - \frac{\partial}{\partial \sigma} \Phi\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right)}{\Phi\left(\frac{r_i - \beta^t \mathbf{z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right)} \\ &= -\frac{\sigma^{-2}}{\sqrt{2\pi}} \sum_{i=1}^n \frac{(r_i - \beta^t \mathbf{z}_i) e^{-\frac{(r_i - \beta^t \mathbf{z}_i)^2}{2\sigma^2}} - (l_i - \beta^t \mathbf{z}_i) e^{-\frac{(l_i - \beta^t \mathbf{z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right)} \\ &= -\frac{\sigma^{-1}}{\sqrt{2\pi}} \sum_{i=1}^n \frac{\frac{(r_i - \beta^t \mathbf{z}_i)}{\sigma} e^{-\frac{(r_i - \beta^t \mathbf{z}_i)^2}{2\sigma^2}} - \frac{(l_i - \beta^t \mathbf{z}_i)}{\sigma} e^{-\frac{(l_i - \beta^t \mathbf{z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right)} \end{aligned}$$

If $l_i = -\infty$, then $\Phi\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right) = 0$ and thus $\Phi'_\sigma\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right) = 0$. In such a case, we can nevertheless write $\Phi'_\sigma\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right) = -\frac{(l_i - \beta^t \mathbf{z}_i)}{\sigma^2} e^{-\frac{(l_i - \beta^t \mathbf{z}_i)^2}{2\sigma^2}}$, as it still equals 0 for $l_i = -\infty$. A similar argument can be applied to the case $r_i = \infty$. In this way, $\frac{\partial \mathcal{L}(\theta)}{\partial \sigma}$ can have a simpler

form.

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial \beta} &= - \sum_{i=1}^n \left(\frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \frac{\sigma^{-1}}{\sqrt{2\pi}} \mathbf{Z}_i \right) \\ &= - \frac{\sigma^{-1}}{\sqrt{2\pi}} \sum_{i=1}^n \left(\frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \mathbf{Z}_i \right).\end{aligned}$$

Here

$$\frac{\partial \mathcal{L}}{\partial \beta} = \begin{pmatrix} \frac{\partial \mathcal{L}}{\partial \beta_1} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial \beta_p} \end{pmatrix} \text{ and } \frac{\partial \mathcal{L}}{\partial \beta^t} = \left(\frac{\partial \mathcal{L}}{\partial \beta_1} \quad \cdots \quad \frac{\partial \mathcal{L}}{\partial \beta_p} \right).$$

There is no closed form solution to the equation $\frac{\partial \mathcal{L}}{\partial \beta} = \mathbf{0}$ and $\frac{\partial \mathcal{L}}{\partial \sigma} = 0$. So we shall use the Newton-Raphson method to obtain the MLE of $\begin{pmatrix} \beta \\ \sigma \end{pmatrix}$:

$$\begin{pmatrix} \beta_{(k+1)} \\ \sigma_{(k+1)} \end{pmatrix} = \begin{pmatrix} \beta_{(k)} \\ \sigma_{(k)} \end{pmatrix} - \begin{pmatrix} \frac{\partial^2 \mathcal{L}}{\partial \beta \partial \beta^t} & \frac{\partial^2 \mathcal{L}}{\partial \sigma \partial \beta} \\ \frac{\partial^2 \mathcal{L}}{\partial \beta^t \partial \sigma} & \frac{\partial^2 \mathcal{L}}{\partial \sigma^2} \end{pmatrix}^{-1} \Big|_{\beta=\beta_{(k)}, \sigma=\sigma_{(k)}} \begin{pmatrix} \frac{\partial \mathcal{L}}{\partial \beta} \\ \frac{\partial \mathcal{L}}{\partial \sigma} \end{pmatrix} \Big|_{\beta=\beta_{(k)}, \sigma=\sigma_{(k)}}.$$

Here

$$\begin{aligned}\frac{\partial^2 \mathcal{L}(\theta)}{\partial \beta \partial \beta^t} &= - \sum_{i=1}^n \left[\frac{\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma^2} e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - \frac{l_i - \beta^t \mathbf{Z}_i}{\sigma^2} e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\left(\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)\right) \sqrt{2\pi\sigma^2}} + \left(\frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \right)^2 \frac{1}{2\pi\sigma^2} \right] \mathbf{Z}_i \mathbf{Z}_i^t \\ &= - \sum_{i=1}^n \left[\frac{\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma} e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - \frac{l_i - \beta^t \mathbf{Z}_i}{\sigma} e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\left(\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)\right) \sqrt{2\pi}} + \left(\frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \right)^2 \frac{1}{2\pi} \right] \frac{\mathbf{Z}_i \mathbf{Z}_i^t}{\sigma^2}. \\ \frac{\partial^2 \mathcal{L}(\theta)}{\partial \sigma \partial \beta} &= - \frac{\sigma^{-2}}{\sqrt{2\pi}} \sum_{i=1}^n \mathbf{Z}_i \left\{ \frac{\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{\sigma^2} e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - \frac{(l_i - \beta^t \mathbf{Z}_i)^2}{\sigma^2} e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} - \frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \right. \\ &\quad \left. + \frac{(r_i - \beta^t \mathbf{Z}_i) e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - (l_i - \beta^t \mathbf{Z}_i) e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \cdot \frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \cdot \frac{\sigma^{-1}}{\sqrt{2\pi}} \right\} \\ &= - \frac{\sigma^{-2}}{\sqrt{2\pi}} \sum_{i=1}^n \mathbf{Z}_i \left\{ \frac{\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{\sigma^2} e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - \frac{(l_i - \beta^t \mathbf{Z}_i)^2}{\sigma^2} e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} - \frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \right. \\ &\quad \left. + \frac{\frac{(r_i - \beta^t \mathbf{Z}_i)}{\sigma} e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - \frac{(l_i - \beta^t \mathbf{Z}_i)}{\sigma} e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\left(\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)\right)^2} \cdot \frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\sqrt{2\pi}} \right\}. \\ \frac{\partial^2 \mathcal{L}(\theta)}{\partial \sigma^2} &= 2 \frac{\sigma^{-3}}{\sqrt{2\pi}} \sum_{i=1}^n \frac{(r_i - \beta^t \mathbf{Z}_i) e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - (l_i - \beta^t \mathbf{Z}_i) e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \\ &\quad - \frac{\sigma^{-2}}{\sqrt{2\pi}} \sum_{i=1}^n \frac{\frac{(r_i - \beta^t \mathbf{Z}_i)^3}{\sigma^3} e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - \frac{(l_i - \beta^t \mathbf{Z}_i)^3}{\sigma^3} e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \\ &\quad - \left(\frac{\sigma^{-2}}{\sqrt{2\pi}} \right)^2 \sum_{i=1}^n \left(\frac{(r_i - \beta^t \mathbf{Z}_i) e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - (l_i - \beta^t \mathbf{Z}_i) e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)} \right)^2 \\ &= 2 \frac{\sigma^{-2}}{\sqrt{2\pi}} \sum_{i=1}^n \frac{\frac{(r_i - \beta^t \mathbf{Z}_i)}{\sigma} e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - \frac{(l_i - \beta^t \mathbf{Z}_i)}{\sigma} e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}\right)}\end{aligned}$$

$$\begin{aligned}
& - \frac{\sigma^{-2}}{\sqrt{2\pi}} \sum_{i=1}^n \frac{\frac{(r_i - \beta^t \mathbf{z}_i)^3}{\sigma^3} e^{-\frac{(r_i - \beta^t \mathbf{z}_i)^2}{2\sigma^2}} - \frac{(l_i - \beta^t \mathbf{z}_i)^3}{\sigma^3} e^{-\frac{(l_i - \beta^t \mathbf{z}_i)^2}{2\sigma^2}}}{\Phi\left(\frac{r_i - \beta^t \mathbf{z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right)} \\
& - \frac{\sigma^{-2}}{2\pi} \sum_{i=1}^n \left(\frac{\frac{(r_i - \beta^t \mathbf{z}_i)}{\sigma} e^{-\frac{(r_i - \beta^t \mathbf{z}_i)^2}{2\sigma^2}} - \frac{(l_i - \beta^t \mathbf{z}_i)}{\sigma} e^{-\frac{(l_i - \beta^t \mathbf{z}_i)^2}{2\sigma^2}}}{\left(\Phi\left(\frac{r_i - \beta^t \mathbf{z}_i}{\sigma}\right) - \Phi\left(\frac{l_i - \beta^t \mathbf{z}_i}{\sigma}\right)\right)} \right)^2.
\end{aligned}$$

Some useful asymptotic results.

Under certain regularity conditions, the following asymptotic properties are valid and are used for testing statistical hypotheses and constructing confidence intervals or confidence regions for $\theta \in \Theta \subset R^m$.

- A. $\left(-\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t} \Big|_{\phi=\hat{\theta}}\right)^{1/2} (\hat{\theta} - \theta)$ is approximately $N(0, I_m)$ distributed if n is large, where I_m is a $m \times m$ identity matrix.
- B. $(\hat{\theta} - \theta)^t \left(-\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t} \Big|_{\phi=\hat{\theta}}\right) (\hat{\theta} - \theta)$ is approximately $\chi^2(m)$ distributed if n is large,
- C. For testing $H_0: \theta \in \Theta_0$ v.s. $H_1: \theta \notin \Theta_0$, an asymptotic test is the likelihood ratio test and another test is a score test. Here

the likelihood ratio test is $\mathbf{1}_{\left(-2 \log \frac{\mathbf{L}(\hat{\theta}_0)}{\mathbf{L}(\hat{\theta})} > \chi_{\alpha, d}^2\right)}$ where $\hat{\theta}$ and $\hat{\theta}_0$ are the MLEs of θ in the

space Θ and Θ_0 , respectively, $d = \|\Theta\| - \|\Theta_0\|$ and $\|\Theta_0\|$ is the dimension of Θ_0 ;

The score test is based on $\frac{\partial \mathcal{L}}{\partial \phi} \Big|_{\phi=\hat{\theta}_0}$. In particular, if $\Theta_o = \{\theta_o\}$, then the test is

$\mathbf{1}_{(T > \chi_{\alpha, d}^2)}$, where $T = \left(\frac{\partial \mathcal{L}}{\partial \phi} \Big|_{\phi=\hat{\theta}_0}\right)^t \left(-\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t} \Big|_{\phi=\hat{\theta}_0}\right)^{-1} \frac{\partial \mathcal{L}}{\partial \phi} \Big|_{\phi=\hat{\theta}_0}$.

3.4.2. Monte-Carlo method.

Generate a large number N of potential values of θ , say t_1, \dots, t_N .

Let $\hat{\theta} = \arg \max_{t_i} \{\mathbf{L}(t_i)\}$.

Remark. The drawback:

1. Time-consuming,
2. Do not guarantee to approximate the true MLE.

A simulation example. Generate data from $S_Y(y; \beta) = \exp(-(y/\tau)^\gamma)$, $y > 0$, $\tau = e^6$, $\gamma = 1/3$.

Then pretend we do not know (τ, γ) . Find its MLE.

```

#in R
library(splines)
library(survival)
#RC data
set.seed(1)
b=exp(6)
g=1/3
y=rweibull(100,g,b)           # mean(y) [1] 4145.353 # median(y) [1] 287.3097
c=runif(100,0,780)
d=as.numeric(y<=c)
m=y*d+c*(1-d)
#Monte Carlo
set.seed(1)
N=200
L=0
M=c(0,0)
for(i in 1:N) {
  bb=rnorm(1,b,100)
  gg=rnorm(1,g,1/9)
  a=prod(d*dweibull(m,gg,bb)+(1-d)*(1-pweibull(m,gg,bb)))
  if (L<a) {
    L=a
    M=c(gg,bb)
  }
}

```

```

}
}
> a
[1] 5.802791e-150
> L
[1] 2.207919e-149
> M #Monte Carlo estimate of
[1] 0.3563931 388.9888333
> c(g,b) # true value
[1] 0.3333333 403.4287935

```

3.4.3. R commands.

Probability Distribution Functions in R.

Let X be a random variable (rv).

Its cdf $F(t) = P\{X \leq t\}$,

density function (df) $f(t) = \begin{cases} F'(t) & \text{if } X \text{ is continuous,} \\ F(t) - F(t-) & \text{if } X \text{ is discrete.} \end{cases}$

quartile $Q(u) = F^{-1}(u) = \min\{t : F(t) \geq u\}$.

Example 1. $X \sim$ Weibull distribution with cdf

$$F(x|\gamma, \tau) = 1 - \exp(-(x/\tau)^\gamma), \quad x > 0$$

γ - shape, τ - scale,

`pweibull(x,shape,scale)` — $F(x)$

`qweibull(x,shape,scale)` — $Q(x)$

`dweibull(x ,shape,scale)` — $f(x)$.

`rweibull(10 ,shape,scale)` — 10 observations.

Remark. The list of all distributions is given in Table 5.1.

<i>Dist</i>	<i>S name</i>	<i>parameters</i>
<i>beta</i>	<i>beta</i>	<i>shape1, shape2</i>
<i>binomial</i>	<i>binom</i>	<i>size, prob</i>
<i>Cauchy</i>	<i>cauchy</i>	<i>location, scale</i>
<i>chi - square</i>	<i>chisq</i>	<i>df</i>
<i>exponential</i>	<i>exp</i>	<i>rate</i>
<i>F</i>	<i>f</i>	<i>df1, df2</i>
<i>gamma</i>	<i>gamma</i>	<i>shape, rate</i>
<i>geometric</i>	<i>geom</i>	<i>prob</i>
<i>hypergeometric</i>	<i>hyper</i>	<i>m, n, k</i>
<i>log - normal</i>	<i>lnorm</i>	<i>meanlog, sdlog</i>
<i>logistic</i>	<i>logis</i>	<i>location, scale</i>
<i>negative binomial</i>	<i>nbinom</i>	<i>size, prob</i>
<i>normal</i>	<i>norm</i>	<i>mean, sd</i>
<i>Poisson</i>	<i>pois</i>	<i>lambda</i>
<i>T</i>	<i>t</i>	<i>df</i>
<i>uniform</i>	<i>unif</i>	<i>min, max</i>
<i>Weibull</i>	<i>weibull</i>	<i>shape, scale</i>
<i>Wilcox</i>	<i>wilcox</i>	<i>m, n</i>

The R presents MLE with regression data for additional distributions as follows.

1. weibull distribution

Standard form $S(t) = \exp(-t^\gamma/\theta)$, $t > 0$.

With covariate in R, reparametrization:

$$S_Y(y|\mathbf{x}, \beta, \sigma) = \exp(-(y/e^{\beta'\mathbf{x}})^{1/\sigma}), \quad y > 0, \quad \mathbf{x}' = (1, x_1, \dots, x_p)$$

$$\ln Y = \beta'\mathbf{x} + \sigma \ln T, \quad T \sim \text{Exp}(1).$$

$$T = (Y/e^{\beta'\mathbf{x}})^{1/\sigma}.$$

2. exponential distribution

Standard form $S(t) = \exp(-t/\theta)$, $t > 0$

With covariate in R, reparametrization:

$$S_Y(y|\mathbf{x}, \beta) = \exp(-y/e^{\beta'\mathbf{x}}), y > 0.$$

3. gaussian distribution

Standard form $N(\mu, \sigma^2)$: $f(t) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(t-\mu)^2}{2\sigma^2})$.

With covariate in R, reparametrization:

$$f_Y(y|\mathbf{x}, \beta, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(y-\beta'\mathbf{x})^2}{2\sigma^2})$$

or $Y = \beta'\mathbf{x} + \sigma Z$, $Z \sim N(0, 1)$.

4. logistic distribution

Standard form $logistic(0, 1)$: $S(t) = \frac{1}{1+\exp(t)}$.

With covariate in R, reparametrization:

$$S_Y(y|\mathbf{x}, \beta, \tau) = \frac{1}{1+\exp(\frac{y-\beta'\mathbf{x}}{\tau})}$$

$Y = \beta'\mathbf{x} + \tau Z$, $Z \sim logistic(0, 1)$, with $\sigma_Z = \pi/\sqrt{3}$.

5. lognormal distribution

Assume $\ln Y = \beta'\mathbf{x} + \sigma Z$, where $Z \sim N(0, 1)$.

6. loglogistic distribution

$\ln Y = \beta'\mathbf{x} + \tau Z$, $Z \sim logistic(0, 1)$.

R command:

The parametric MLE is efficient under certain regularity assumptions. In particular, if the residual plot suggests that certain parametric family is plausible, one can apply the Splus functions as follows.

```
zz=survReg(Surv(y)~x, dist="exponential")
dist: (default: weibull), gaussian, logistic, lognormal and loglogistic
predict(zz,data.frame(x=130),se=T)
summary(zz)
```

in R

```
library(splines)
library(survival)
survreg in R =survReg in Splus
```

Example 1. Notice that the scale in `rweibull()` is different from the scale in `survRev()`.

Standard form: $S(x|\gamma, \tau) = \exp(-x^\gamma/\theta)$, $x > 0$.

In `rweibull()`, $S(x|\gamma, \tau) = \exp(-(x/\tau)^\gamma)$, $x > 0$, and $\gamma =$ shape and $\tau =$ scale,

However in `survreg(· ~ 1)`:

$\ln Y = \alpha + \sigma \ln Z$, where $Z \sim Exp(1)$, $\alpha =$ intercept and $\sigma =$ scale

(in `survreg(· ~ x)`: $\ln Y = \alpha + \beta x + \sigma \ln Z$).

Q: Relation between (γ, τ) and (σ, α) ? $Y = \exp(\alpha)Z^\sigma$

$$S_Y(t) = P(e^\alpha Z^\sigma > t) = \exp(-e^{-\alpha/\sigma} t^{1/\sigma}) = \exp(-(t/e^\alpha)^{1/\sigma}) = \exp(-(t/\tau)^\gamma)$$

$$\Rightarrow \gamma = 1/\sigma \text{ and } \tau = e^\alpha.$$

An simulation example with R (a continuation of Monte Carlo example):

```
> library(splines)
> library(survival)

> (zz=survreg(Surv(m,d)~1))
Coefficients:
(Intercept)
5.985651
Scale= 2.781783
Loglik(model)= -342.3 Loglik(intercept only)= -342.3

> summary(zz)
      Value Std.Error      z      p
(Intercept)  5.99    0.371  16.14 1.33e-58
Log(scale)   1.02    0.113   9.02 1.94e-19
Scale= 2.78
Loglik(model)= -342.3 Loglik(intercept only)= -342.3
```



```

> c(1/M[1],log(M[2]))
[1] 2.805890 5.963551
> (x= c(zz$scale, zz$coef))
2.781783 5.985651 # c(1/M[1],log(M[2])) is its approximation.

> gg=1/x[1]
> bb=exp(x[2])
> prod(d*dweibull(m,gg,bb)+(1-d)*(1-pweibull(m,gg,bb)))
[1] 2.221519e-149 # > L = 2.207919e - 149 in Monte Carlo
> x=rpois(100,1)
> (zz=survreg(Surv(m,d)~x+I(x^2))) # log y=6+0*x+0x^2+3*ln z, z~exp(1)
(Intercept)      x      I(x^2)
5.9362703  1.4134716  -0.6472259
Scale= 2.731045
Loglik(model)= -340.1 Loglik(intercept only)= -342.3
Chisq= 4.45 on 2 degrees of freedom, p= 0.11
> summary(zz)

```

	Value	Std.Error	z	p
(Intercept)	5.936	0.533	11.13	8.98e - 29
x	1.413	1.062	1.33	1.83e - 01
I(x ²)	-0.647	0.360	-1.80	7.19e - 02
Log(scale)	1.005	0.113	8.92	4.81e - 19

```

Scale= 2.73
Loglik(model)= -340.1 Loglik(intercept only)= -342.3
Chisq= 4.45 on 2 degrees of freedom, p= 0.11
Number of Newton-Raphson Iterations: 5 n= 100

> z=summary(zz)
> names(z)
[1] "call" "df" "loglik" "iter" "idf"
[6] "scale" "coefficients" "var" "table" "correlation"
[11] "parms" "n" "chi" "robust"
> z$table

```

	Value	Std.Error	z	p
(Intercept)	5.9362703	0.5333632	11.129884	8.975185e - 29
x	1.4134716	1.0615812	1.331478	1.830319e - 01
I(x ²)	-0.6472259	0.3596803	-1.799448	7.194787e - 02
Log(scale)	1.0046844	0.1126759	8.916585	4.808864e - 19

```

> z$var # covariance matrix
(Intercept)      x      I(x^2)      Log(scale)
(Intercept) 0.32635298 -0.38842957 0.095005834 0.016751330
x -0.38842957 1.08207114 -0.345080622 -0.015565844
I(x^2) 0.09500583 -0.34508062 0.125869293 0.004110281
Log(scale) 0.01675133 -0.01556584 0.004110281 0.012808764

> #complete data (uncensored)
> #set.seed(1)
> n=100
> b=exp(6)
> g=1/3
> y=rweibull(n,g,b) # S(y) = exp(-(y/e^6)^1/3)
> mean(y) [1] 2224.555
> set.seed(1)
> z=rexp(n)
> y=exp(6+3*log(z)) # log y=6+3*ln z

```

```

> mean(y) [1] 2117.473
> (zz=survreg(Surv(y)~1, dist="exponential"))
  C coefficients:
  (Intercept)
  7.657979
  Scale fixed at 1
  Loglik(model)= -865.8 Loglik(intercept only)= -865.8
  n= 100
> (zz=survreg(Surv(y)~1))
  Coefficients:
  (Intercept)
  6.259506
  Scale= 2.570908
  Loglik(model)= -733.5 Loglik(intercept only)= -733.5 n= 100
> (zz=survreg(Surv(y)~1, dist="weibull"))
  Coefficients:
  (Intercept)
  6.259506
  Scale= 2.570908
  Loglik(model)= -733.5 Loglik(intercept only)= -733.5 n= 100
> set.seed(1)
> y=rweibull(100,1/3, exp(6))
> (zz=survreg(Surv(y)~1))
  (Intercept)
  5.74657
  Scale= 2.66256
  Loglik(model)= -679.2 Loglik(intercept only)= -679.2 n= 100
> summary(zz)
      Value Std.Error    z      p
(Intercept)  5.747    0.2810  20.5  5.77e - 93
  Log(scale)  0.979    0.0769  12.7  3.94e - 37
Scale= 2.66
Loglik(model)= -679.2 Loglik(intercept only)= -679.2
Number of Newton-Raphson Iterations: 6 n= 100

> y=rnorm(100,8,4)
> (zz=survreg(Surv(y)~1,dist="gaussian"))
  (Intercept)
  7.929682
  Scale= 3.741224
  Loglik(model)= -273.8 Loglik(intercept only)= -273.8 n= 100
> u=exp(y)
> (zz=survreg(Surv(u)~1,dist="lognormal"))
  (Intercept)
  7.929682
  Scale= 3.741224
  Loglik(model)= -1066.8 Loglik(intercept only)= -1066.8
  n= 100

```

§3.4.4. Homework:

1. Generate a RC data set of size 100 from a Weibull distribution with $\kappa = 2$ and $\rho = 1$ (where $S(t) = \exp(-(\rho t)^\kappa)$).
2. Estimate the MLE of the parameter using the Newton-Raphson method (you could use the command in R:
`y=rweibull(100,1/2,1) #kappa = 1/2 scale = 1`
`c=runif(100,0,2)`

```
d=as.numeric(y <= c)
m=d*y+(1-d)*c
yy=survreg(Surv(m,d)~ 1, dist="weibull")
but you need at least to derive the iteration formula.)
```

3. Estimate the covariance matrix of the MLE.
4. Derive a 95% confidence interval for ρ .
5. Test $H_0: \kappa = 1$ v.s. $H_1: \kappa \neq 1$ using the data you generated with size $\alpha = 0.1$. This is to test whether the distribution is actually from an $\text{Exp}(\rho)$ rather from the Weibull distribution.
6. Is the result in problem 5 as what you expected ? Why you say so?
7. If you have a sample of size 4, do you still expect to see what you expect in Problem 6?

§3.5. Consistency and Asymptotic Normality:

Two issues:

1. Is the MLE consistent ?
2. Is the MLE asymptotically normally distributed ?

They need to be verified for each problem. It is easy to verify if the MLE has a closed form expression, e.g., in the example considered in §3.3. In general, it is not so trivial. We shall illustrate the usual approach through the problem of estimating the parameter of an exponential distribution under the C2 model. In particular, we assume:

- 1) $X \sim \text{Exp}(\rho_o)$,
- 2) (U, V) has a joint c.d.f. $G(u, v)$, where $0 \leq U < V$ w.p.1.,
- 3) X and (U, V) are independent;
- 4) $(L, R) = \begin{cases} (U, V) & \text{if } X \in (U, V], \\ (V, +\infty) & \text{if } X > V, \\ (-\infty, U) & \text{if } X \leq U, \end{cases}$

Let $(X_i, U_i, V_i, L_i, R_i)$, $i = 1, \dots, n$, be i.i.d. copies of (X, U, V, L, R) .

The log likelihood function

$$\mathcal{L}(\rho) = \sum_{i:lc} \log F_o(R_i; \rho) + \sum_{i:ic} \log [S_o(L_i; \rho) - S_o(R_i; \rho)] + \sum_{i:rc} \log S_o(L_i; \rho). \quad (1.1)$$

Eq. (1.1) yields that

$$\mathcal{L}(\rho) = \sum_{i,lc} \log(1 - e^{-\rho R_i}) + \sum_{i,ic} \log(e^{-\rho L_i} - e^{-\rho R_i}) - \rho \sum_{i,rc} L_i. \quad (1.2)$$

To find the MLE, we look at the normal equation $\frac{\partial \mathcal{L}}{\partial \rho} = 0$:

$$\sum_{i:lc} \frac{R_i \exp(-\rho R_i)}{1 - \exp(-\rho R_i)} - \sum_{i:rc} L_i - \sum_{i:ic} \frac{L_i \exp(-\rho L_i) - R_i \exp(-\rho R_i)}{\exp(-\rho L_i) - \exp(-\rho R_i)} = 0.$$

There is no closed form expression for the MLE of ρ_o .

Thus we need to use the Newton-Raphson method to derive it.

§3.5.1. Existence of the MLE.

In some extreme case, the MLE may not exist. However, one can always modify the MLE so that it is well defined.

Theorem 1. Suppose that $X \sim \text{Exp}(\rho_o)$. Under the C2 model, the MLE $\hat{\rho}$ exists and $\hat{\rho} \in (0, +\infty)$ unless $\sum_{i,rc} 1 = \sum_i \mathbf{1}_{(R_i=\infty)} = n$ or $\sum_{i,lc} 1 = \sum_i \mathbf{1}_{(L_i=-\infty)} = n$.

Q: $P(\sum_{i,rc} 1 = \sum_i \mathbf{1}_{(R_i=\infty)} = n \text{ or } \sum_{i,lc} 1 = \sum_i \mathbf{1}_{(L_i=-\infty)} = n) = \mathbf{0}$?

Proof. First assume that $\sum_{rc} 1 < n$ and $\sum_{lc} 1 < n$. Verify from (1.2) that

- (i) $\lim_{\rho \rightarrow 0^+} \mathcal{L}(\rho) = -\infty$;
- (ii) $\lim_{\rho \rightarrow +\infty} \mathcal{L}(\rho) = -\infty$;
- (iii) $\mathcal{L}(\rho)$ is continuous in ρ .

It follows that the MLE $\hat{\rho}$ exists and $\hat{\rho} \in (0, +\infty)$.

More specifically, by (i) and (ii), $\exists v > 1$ such that

$$\mathcal{L}(\rho) < \mathcal{L}(1) \text{ if } \rho \notin [1/v, v].$$

Then $\mathcal{L}(\rho)$ is continuous on $[1/v, v]$ and achieves its maximum in $[1/v, v]$ (**how about** $(1/v, v)$?) \square

Remark. If $\sum_{rc} 1 = n$ then $\mathcal{L}(\rho)$ ($= -\rho \sum_{rc} V_i$) is maximized uniquely by $\rho = 0 \notin (0, \infty)$. Thus the MLE does not exist and we define $\hat{\rho} = 0$. If $\sum_{lc} 1 = n$, then $\mathcal{L}(\rho)$ ($= \sum_{lc} \log(1 - e^{-\rho U_i})$) is maximized uniquely by $\rho = \infty \notin (0, \infty)$. Thus the MLE does not exist and we define $\hat{\rho} = +\infty$. In this way, $\hat{\rho}$ is properly defined in the whole sample space and thus we can study its properties.

§3.5.2. Consistency of the MLE.

Hereafter, abusing notations, we write $S = S_o$ etc.. The normalized log likelihood function

$$\begin{aligned} l(\rho) &= \frac{1}{n} \left\{ \sum_{i: ic} \log[S(L_i; \rho) - S(R_i; \rho)] + \sum_{i: rc} \log S(L_i; \rho) + \sum_{i: lc} \log F(R_i; \rho) \right\} \\ &= \frac{1}{n} \left\{ \sum_{i: ic} \log[S(U_i; \rho) - S(V_i; \rho)] + \sum_{i: rc} \log S(V_i; \rho) + \sum_{i: lc} \log F(U_i; \rho) \right\}. \quad (2.1) \\ &= \frac{1}{n} \left(\sum_{ic} \log e^{-\rho U_i} + \sum_{ic} \log(1 - e^{-\rho(V_i - U_i)}) - \rho \sum_{rc} V_i + \sum_{lc} \log(1 - e^{-\rho U_i}) \right) \\ &= \frac{1}{n} \left[\sum_{i=1}^n (\log(1 - e^{-\rho U_i})) \mathbf{1}_{(X_i \leq U_i)} + \sum_{i=1}^n (\log(1 - e^{-\rho(V_i - U_i)})) \mathbf{1}_{(X_i \in (U_i, V_i])} \right. \\ &\quad \left. - \rho \sum_{i=1}^n U_i \mathbf{1}_{(X_i \in (U_i, V_i])} - \rho \sum_{i=1}^n V_i \mathbf{1}_{(X_i > V_i)} \right] \end{aligned}$$

Theorem 2. Suppose that $X \sim \text{Exp}(\rho_o)$. Under the C2 model with $P(0 < U < V) = 1$, the MLE $\hat{\rho}$ of ρ is strongly consistent.

Proof of Theorem 2. Let

$$\mu(\rho) = E(l(\rho)).$$

We shall show

$$\mu(\rho) \leq \mu(\rho_o) \text{ with equality iff } \rho = \rho_o. \quad (2.2)$$

If ρ^* is a limiting point of $\hat{\rho}$, then $\rho^* \in (0, \infty)$ and $\mu(\rho^*) \geq \mu(\rho_o)$ on Ω^* with $P(\Omega^*) = 1$. (2.3)

(2.2) and (2.3) imply that $\mu(\rho^*) = \mu(\rho_o)$ and thus $\rho^* = \rho_o$ by (2.2). Moreover, since ρ^* is an arbitrary limiting point of $\hat{\rho}$, it implies that $\hat{\rho} \rightarrow \rho_o$ w.p.1. In other words, (2.2) and (2.3) imply that $\hat{\rho}$ is strongly consistent. Note that a limiting point can be $\pm\infty$.

Proof of (2.2): Write $l(\rho) = l(\mathbf{L}, \mathbf{R}, \rho)$, where $\mathbf{L} = (L_1, \dots, L_n)$ and $\mathbf{R} = (R_1, \dots, R_n)$. Then

$$\begin{aligned} \mu(\rho) & (= E(l(\mathbf{L}, \mathbf{R}, \rho))) \\ & = E(E(l(\mathbf{L}, \mathbf{R}, \rho) | \mathbf{U}, \mathbf{V})) \\ & = E(E[\mathbf{1}_{(X \leq U)} \log F(U; \rho) + \mathbf{1}_{(X > V)} \log S(V; \rho) + \mathbf{1}_{(U < X \leq V)} \log(S(U; \rho) - S(V; \rho)) | U, V]). \\ & = E[F(U; \rho_o) \log F(U; \rho) + S(V; \rho_o) \log S(V; \rho) + (S(U; \rho_o) - S(V; \rho_o)) \log(S(U; \rho) - S(V; \rho))]. \end{aligned}$$

Define $\log 0 = -\infty$, and $0 \log 0 = 0$. Then $\mu(\rho)$ is uniquely maximized by $\rho = \rho_o$, due to the following three reasons:

(a) For each (u, v) ,

$$F(u; \rho_o) \log F(u; \rho) + (S(u; \rho_o) - S(v; \rho_o)) \log(S(u; \rho) - S(v; \rho)) + S(v; \rho_o) \log S(v; \rho)$$

as a function of $F(\cdot, \rho)$ is uniquely maximized by

$$F(u, \rho) = F(u, \rho_o),$$

$$S(u, \rho) - S(v, \rho) = S(u, \rho_o) - S(v, \rho_o)$$

and $S(v, \rho) = S(v, \rho_o)$,

(b) the last equation in (a) implies $\rho = \rho_o$ as $S(x, \rho) = e^{-\rho x}$.

(c) $|\mu(\rho_o)| < +\infty$.

Remark. (a) is equivalent to say that $g(\cdot, \cdot)$, defined by

$$g(q_1, q_2) = p_1 \log q_1 + p_2 \log q_2 + (1 - p_1 - p_2) \log(1 - q_1 - q_2), \quad q_i, p_i \geq 0, \quad p_1 + p_2 \leq 1, \quad q_1 + q_2 \leq 1,$$

is uniquely maximized by $q_1 = p_1$ and $q_2 = p_2$. (Hint : take partial derivatives and set to zero, etc., or by the Shannon-Kolmogorov inequality).

To prove (c) note that

each summand in $g(p_1, p_2)$ is of the form $x \ln x$,

$0 \geq x \log x \geq -1/e$, for $x \in [0, 1]$,

and thus $g(p_1, p_2) \geq -3/e$,

Thus (c) follows.

(a), (b) and (c) imply that

$$\mu(\rho_o) > \mu(\rho) \text{ if } \rho \neq \rho_o \text{ that is, (2.2) holds.}$$

Prove (2.3). We shall now construct Ω^* .

By definition of the modified MLE $\hat{\rho}$, $l(\hat{\rho}(\omega))(\omega) \geq l(\rho_o)(\omega) \quad \forall \omega$ in the sample space Ω . Thus

$$\liminf_{n \rightarrow \infty} l(\hat{\rho}) \geq \lim_{n \rightarrow \infty} l(\rho_o) = \mu(\rho_o) \quad (2.4)$$

a.s. (by SLLN).

$$\begin{aligned} l(\rho) &= \sum_{i=1}^n \frac{\log(1 - e^{-\rho U_i}) \mathbf{1}_{(X_i \leq U_i)}}{n} + \sum_{i=1}^n \frac{\log(1 - e^{-\rho(V_i - U_i)}) \mathbf{1}_{(X_i \in (U_i, V_i])}}{n} \\ &\quad - \rho \sum_{i=1}^n \frac{U_i \mathbf{1}_{(X_i \in (U_i, V_i])}}{n} - \rho \sum_{i=1}^n \frac{V_i \mathbf{1}_{(X_i > V_i)}}{n}. \end{aligned} \quad (2.5)$$

For each ρ , the four summations in (2.5) all converges a.s. to their means, respectively. However, it is not clear that $P(l(\rho) \rightarrow \mu(\rho) \text{ for all } \rho > 0) = 1$.

Thus, we let K be the set of all positive rational numbers and ρ_o , and

let Ω^* be the event

* such that (2.4) holds, and for all $\rho \in K$, the four summations in (2.5) converges a.s. to their means, respectively.

Then $P(\Omega^*) = 1$ as K is countable.

To emphasize that the MLE $\hat{\rho}$ is a function of n , we write $\hat{\rho} = \hat{\rho}_n$. For each ω in Ω^* , let ρ^* be a limiting point of $\hat{\rho} = \hat{\rho}_n(\omega)$ in the sense that

$\hat{\rho}_{n_j}(\omega) \rightarrow \rho^*$ for a subsequence of $\{\hat{\rho}_n\}$, where ρ^* may be $+\infty$ or 0.

In order to prove inequality (2.3), it suffices to prove

$$\mu(\rho_o) \leq \mu(\rho^*), \quad (2.6)$$

We shall show (2.6) hereafter.

Assume $\omega \in \Omega^*$. If $\rho^* = +\infty$, then Eq. (2.5) implies that

$$l(\hat{\rho}_{n_j}(\omega)) \leq -\hat{\rho}_{n_j}(\omega) \left[\sum_{i=1}^n \frac{V_i \mathbf{1}_{(X_i > V_i)}}{n} \right] (\omega) \rightarrow -\infty \quad (2.7)$$

Then Eq. (2.7) and inequality (2.4) imply that

$$-\infty = \lim_{n_j \rightarrow \infty} l(\hat{\rho}(\omega)) \geq \mu(\rho_o).$$

It reaches a contradiction as $\mu(\rho_o)$ is finite by (2.2). Thus $\rho^* = +\infty$ is impossible.

If $\rho^* = 0$, the first two sums in Eq. (2.5) tend $-\infty$ for large enough n as $\log 0 = -\infty$ by our convention. Thus it can be shown $\rho^* = 0$ is impossible too.

Then for $\omega \in \Omega^*$, $\rho^* \in (0, +\infty)$. For any $m, M \in K$ satisfying $m < \rho^* < M$, if n_j is large enough, then

$$\begin{aligned} \sum_{i=1}^{n_j} \log(1 - e^{-mU_i(\omega)}) \mathbf{1}_{(X_i(\omega) \leq U_i(\omega))} / n_j &\leq \sum_{i=1}^{n_j} \log(1 - e^{-\hat{\rho}_{n_j}(\omega)U_i(\omega)}) \mathbf{1}_{(X_i(\omega) \leq U_i(\omega))} / n_j \\ &\leq \sum_{i=1}^{n_j} \log(1 - e^{-MU_i(\omega)}) \mathbf{1}_{(X_i(\omega) \leq U_i(\omega))} / n_j. \end{aligned}$$

In an obvious way, rewrite the above inequalities as

$$\Psi_{n_j}(m, \omega) \leq \Psi_{n_j}(\hat{\rho}_{n_j}(\omega), \omega) \leq \Psi_{n_j}(M, \omega). \quad (2.8)$$

Then

$$\begin{aligned} \liminf_{j \rightarrow \infty} \Psi_{n_j}(m, \omega) &\leq \liminf_{j \rightarrow \infty} \Psi_{n_j}(\hat{\rho}_{n_j}(\omega), \omega) \\ &\leq \limsup_{j \rightarrow \infty} \Psi_{n_j}(\hat{\rho}_{n_j}(\omega), \omega) \leq \limsup_{j \rightarrow \infty} \Psi_{n_j}(M, \omega). \end{aligned} \quad (2.9)$$

Note that since $m, M \in K$,

$$\liminf_{j \rightarrow \infty} \Psi_{n_j}(m, \omega) = \lim_{n \rightarrow \infty} \Psi_n(m, \omega) = E(\log(1 - e^{-mU}) \mathbf{1}_{(X \leq U)}). \quad (2.10)$$

$$\limsup_{j \rightarrow \infty} \Psi_{n_j}(M, \omega) = \lim_{n \rightarrow \infty} \Psi_n(M, \omega) = E(\log(1 - e^{-MU}) \mathbf{1}_{(X \leq U)}). \quad (2.11)$$

Recall the monotone convergence theorem:

If f_n is a monotone convergent sequence and $f_n \rightarrow f$ and are all integrable, then $\int f_n(x) d\mu(x) \rightarrow \int f(x) d\mu(x)$.

Since $g(\rho, \omega) = \mathbf{1}_{(X_i(\omega) \leq U_i(\omega))} \log(1 - e^{-\rho U_i(\omega)})$ is a monotone function of ρ , $g(m, \cdot)$ (or $g(M, \cdot)$) is an increasing function in m (or M) as $m \uparrow \rho^*$ (or $M \downarrow \rho^*$), and $E(g(m, \cdot)) = \int g(m, \omega) dP$,

by the monotone convergence theorem, taking limits as $m, M \rightarrow \rho^*$ yields

$$E(g(m, \cdot)) \rightarrow E(g(\rho^*, \cdot)) \text{ and } E(g(M, \cdot)) \rightarrow E(g(\rho^*, \cdot)). \quad (2.12)$$

Then it follows from (2.8) through (2.12) that

$$\sum_{i=1}^{n_j} \log(1 - e^{-\hat{\rho}_{n_j}(\omega)U_i(\omega)}) \mathbf{1}_{(X_i(\omega) \leq U_i(\omega))} / n_j \rightarrow E(\log(1 - e^{-\rho^*U}) \mathbf{1}_{(X \leq U)}). \quad (2.13)$$

(Remark: In (2.13), since we are dealing $\lim_{n \rightarrow \infty} \sum_{i=1}^n g_n(\hat{\rho}_n)$, without the arguments from (2.8) through (2.12), we cannot conclude directly,

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n g_n(\hat{\rho}_n) = \lim_{n \rightarrow \infty} \sum_{i=1}^n \lim_{n \rightarrow \infty} g_n(\hat{\rho}_n),$$

even if all the limits make sense.)

Since $\log(1 - e^{-\rho(V_i - U_i)})$ is also a monotone function of ρ , it can be shown that

$$\begin{aligned} \sum_{i=1}^{n_j} \log(1 - e^{-\hat{\rho}_{n_j}(\omega)(V_i - U_i(\omega))}) \mathbf{1}_{(X_i(\omega) \in (U_i(\omega), V_i(\omega)])} / n_j \\ \rightarrow E(\log(1 - e^{-\rho^*(V-U)}) \mathbf{1}_{(X \in (U, V))}). \end{aligned} \quad (2.14)$$

Homework: Prove (2.14) by mimicing the proof of (2.13).

Thus, $\liminf_{n_j \rightarrow \infty} l_{n_j}(\hat{\rho}_{n_j}(\omega)) = \mu(\rho^*)$, as the last two summations in (2.5) do not involve $\hat{\rho}$ and converge a.s. to their means, respectively. Then inequality (2.4) yields

$$\mu(\rho_o) \leq \lim_{n_j \rightarrow \infty} l_{n_j}(\hat{\rho}_{n_j}(\omega)) = \mu(\rho^*), \quad \forall \omega \in \Omega^*,$$

which is Eq. (2.6). This concludes our proof. \square

Remark. The standard theory for proving the consistency of the MLE does not work here, as it usually requires that the MLE belongs to compact set, which is not the case here.

Remark.

1. $P(U < V) = 1 \not\Rightarrow U \not\perp V$.
e.g. If $U \perp V$, $U \sim U(0, 1)$ and $V \sim U(1, 2)$, then $P(U < V) = 1$.
2. $P(U < V) = 1$ and $\sup\{t : F_U(t) < 1\} > \inf\{t : F_V(t) > 0\} \Rightarrow U \not\perp V$.

§3.5.3. Asymptotic Normality of the MLE.

Hereafter, we prove the asymptotic normality under the assumption given in §3.5.

Theorem 3. Suppose that $X \sim \text{Exp}(\rho_o)$. Under the C2 model with $P(0 < U < V) = 1$, the MLE $\hat{\rho}$ of ρ satisfies that $\sqrt{n}(\hat{\rho} - \rho_o) \xrightarrow{D} N(0, \sigma^2)$, where $\sigma^2 = -1/E\left(\frac{\partial^2 l(\rho)}{\partial \rho^2}\right)\Big|_{\rho=\rho_o}$.

Proof. The proof is given in 3 steps.

Step 1 (preliminary). We shall show for each t ,

$$P\left\{\frac{\sqrt{n}(\hat{\rho} - \rho_o)}{\sigma} \leq t\right\} \rightarrow \Phi(t), \quad \text{as } n \rightarrow \infty.$$

Let Ω^o be the event such that $\hat{\rho} \rightarrow \rho_o$, $\lim_{n \rightarrow \infty} \sum_{lc} 1/n < 1$ and $\lim_{n \rightarrow \infty} \sum_{rc} 1/n < 1$. Then $P(\Omega^o) = 1$, as $P(0 < U < V) = 1$.

Now for each $\omega \in \Omega^o$ (and suppress ω in the expressions), Eq. (2.1) yields

$$\begin{aligned} \frac{\partial l(\rho)}{\partial \rho} &= \frac{1}{n} \sum_{i=1}^n \frac{U_i e^{-\rho U_i}}{1 - e^{-\rho U_i}} 1_{(X_i \leq U_i)} + \frac{1}{n} \sum_{i=1}^n \frac{(V_i - U_i) e^{-\rho(V_i - U_i)}}{1 - e^{-\rho(V_i - U_i)}} 1_{(U_i < X_i \leq V_i)} \\ &\quad - \frac{1}{n} \sum_{i=1}^n U_i 1_{(U_i < X_i \leq V_i)} - \frac{1}{n} \sum_{i=1}^n V_i 1_{(X_i > V_i)} \\ &= \frac{1}{n} \sum_{i=1}^n \left(-U_i + \frac{U_i}{1 - e^{-\rho U_i}}\right) 1_{(X_i \leq U_i)} + \frac{1}{n} \sum_{i=1}^n \left(-(V_i - U_i) + \frac{(V_i - U_i)}{1 - e^{-\rho(V_i - U_i)}}\right) 1_{(U_i < X_i \leq V_i)} \\ &\quad - \frac{1}{n} \sum_{i=1}^n U_i 1_{(U_i < X_i \leq V_i)} - \frac{1}{n} \sum_{i=1}^n V_i 1_{(X_i > V_i)} \end{aligned}$$

and

$$-\frac{\partial^2 l(\rho)}{\partial \rho^2} = \frac{1}{n} \sum_{i=1}^n \frac{U_i^2 e^{-\rho U_i}}{(1 - e^{-\rho U_i})^2} 1_{(X_i \leq U_i)} + \frac{1}{n} \sum_{i=1}^n \frac{(V_i - U_i)^2 e^{-\rho(V_i - U_i)}}{(1 - e^{-\rho(V_i - U_i)})^2} 1_{(U_i < X_i \leq V_i)}, \quad (3.1)$$

which are continuous in $\rho \in (0, +\infty)$.

For each $\omega \in \Omega^o$, for n large enough, $\sum_{lc} 1/n < 1$ and $\sum_{rc} 1/n < 1$, thus the MLE $\hat{\rho}(\omega) \in (0, +\infty)$ by Theorem 1. Since $\frac{\partial^2 l(\rho)}{\partial \rho^2}(\omega)$ is continuous in ρ ,

$$\frac{\partial l(\hat{\rho})}{\partial \rho} \left(= \frac{\partial l(\rho)}{\partial \rho}(\omega) \Big|_{\hat{\rho}(\omega)} \right) = 0. \quad (3.2)$$

By the first order Taylor expansion,

$$\frac{\partial l(\rho)}{\partial \rho}(\omega) \Big|_{\rho_o} - \frac{\partial l(\rho)}{\partial \rho}(\omega) \Big|_{\hat{\rho}(\omega)} = \frac{\partial^2 l(\rho)}{\partial \rho^2}(\omega) \Big|_{\rho^*(\omega)} (\rho_o - \hat{\rho}(\omega)), \quad (3.3)$$

provided that $\omega \in \Omega^o$,

where $\rho^*(\omega)$ is between ρ_o and $\hat{\rho}(\omega)$, and thus by the assumption on Ω^o ,

$$\rho^*(\omega) \rightarrow \rho_o, \quad \forall \omega \in \Omega^o. \quad (3.4)$$

Denote

$$p(X, U, V, \rho) = F(U; \rho) \mathbf{1}^{(x \leq v)} \times (S(U; \rho) - S(V; \rho)) \mathbf{1}^{(u < x \leq v)} \times S(V; \rho) \mathbf{1}^{(v < x)}$$

(relation with $l(\rho)$?)

and
$$Z = \frac{\partial}{\partial \rho} \log p(X, U, V, \rho) \Big|_{\rho=\rho_o}.$$

Then
$$l(\rho) = \frac{1}{n} \sum_{i=1}^n \log p(X_i, U_i, V_i, \rho),$$

$$\frac{\partial l}{\partial \rho} \Big|_{\rho=\rho_o} = \bar{Z}.$$

By (3.2) and (3.3),

$$\sqrt{n} \cdot \bar{Z}(\omega) = -\sqrt{n}(\hat{\rho}(\omega) - \rho_o) \frac{\partial^2 l(\rho^*)}{\partial \rho^2}(\omega) \text{ if } \hat{\rho}(\omega) \in (0, \infty) \text{ (or } \omega \in \Omega^o \text{ \& } n \approx \infty). \quad (3.5)$$

By the CLT,

$$\sqrt{n}(\bar{Z} - E(Z)) \xrightarrow{\mathcal{D}} N(0, \sigma_Z^2). \quad (3.6)$$

Step 2. \vdash : In order to prove

$$\sqrt{n}(\hat{\rho} - \rho_o) \xrightarrow{\mathcal{D}} N(0, \tau^2)$$

for some $\tau^2 = \sigma_Z^2 / (E(\frac{\partial^2 l(\rho_o)}{\partial \rho^2}))^2$, it suffices to prove the next two statements:

$$E(Z) = 0 \text{ and } \sigma_Z^2 = E(Z^2) = E\left(\left(\frac{\partial}{\partial \rho} \ln p(X, U, V, \rho_o)\right)^2\right) = -E\left(\frac{\partial^2 \ln p(X, U, V, \rho_o)}{\partial \rho^2}\right). \quad (3.7)$$

$$\frac{\partial^2 l(\rho)}{\partial \rho^2} \Big|_{\rho^*} \rightarrow E\left(\frac{\partial^2 l(\rho_o)}{\partial \rho^2}\right) \text{ a.s.} \quad (3.8)$$

The reason is as follows. (3.5), (3.7) and (3.8) yield

$$\frac{\sqrt{n}(\bar{Z}(\omega) - E(Z))}{\sigma_Z} = \frac{\sqrt{n}(\hat{\rho}(\omega) - \rho_o)}{\frac{\sigma_Z}{-\frac{\partial^2 l(\rho^*)}{\partial \rho^2}(\omega)}} = \frac{\sqrt{n}(\hat{\rho}(\omega) - \rho_o)}{\frac{\sigma_Z}{-\frac{\partial^2 l(\rho_o)}{\partial \rho^2}(\omega)}} \text{ if } \hat{\rho}(\omega) \in (0, \infty)$$

$$\{-\sqrt{n}(\hat{\rho} - \rho_o) \frac{\partial^2 l(\rho^*)}{\partial \rho^2} = \sqrt{n}\bar{Z}\} \supset \{\hat{\rho} \in (0, \infty)\}, \text{ (by (3.5))}$$

$$\text{thus } \{-\sqrt{n}(\hat{\rho} - \rho_o) \frac{\partial^2 l(\rho^*)}{\partial \rho^2} \neq \sqrt{n}\bar{Z}\} \subset \{\hat{\rho} \notin (0, \infty)\}, \quad (3.9)$$

$$\begin{aligned} & |P\{\sqrt{n}(\hat{\rho} - \rho_o)/\tau_n \leq t\} - P\{\sqrt{n}\frac{\bar{Z}}{\sigma_Z} \leq t\}| \quad (\tau_n = -\sigma_Z / \frac{\partial^2 l(\rho^*)}{\partial \rho^2}) \\ &= |P\{\sqrt{n}(\hat{\rho} - \rho_o)/\tau_n \leq t\} \\ &\quad - P\{\sqrt{n}(\hat{\rho} - \rho_o)/\tau_n \leq t, \hat{\rho} \in (0, \infty)\} + P\{\sqrt{n}\frac{\bar{Z}}{\sigma_Z} \leq t, \hat{\rho} \in (0, \infty)\} \quad (= 0 \text{ by (3.5)}) \\ &\quad - P\{\sqrt{n}\frac{\bar{Z}}{\sigma_Z} \leq t\}| \end{aligned}$$

$$\begin{aligned}
&\leq P\{\hat{\rho} \notin (0, \infty)\} + P\{\hat{\rho} \notin (0, \infty)\} \quad (|P(A) - P(A \cap B)| \leq P(B^c)) \quad (\text{by (3.9)}) \\
\text{i.e., } &|P\{\sqrt{n}(\hat{\rho} - \rho_o)/\tau_n \leq t\} - P\{\sqrt{n}\frac{\bar{Z}}{\sigma_Z} \leq t\}| \rightarrow 0, \text{ as } n \rightarrow \infty \text{ due to } \hat{\rho} \rightarrow \rho_o \text{ a.s.} \quad (3.10) \\
&P\{\sqrt{n}\frac{\bar{Z}}{\sigma_Z} \leq t\} = P\{\sqrt{n}\frac{\bar{Z} - E(Z)}{\sigma_Z} \leq t\} \rightarrow \Phi(t) \text{ as } n \rightarrow \infty, \\
&\text{where } \Phi \text{ is the cdf of } N(0, 1) \text{ (due to CLT)} \\
&\text{thus} \\
&P\{\sqrt{n}(\hat{\rho} - \rho_o)/\tau_n \leq t\} \\
&= P\{\underbrace{\sqrt{n}(\hat{\rho} - \rho_o)/\tau_n \leq t - \sqrt{n}\frac{\bar{Z}}{\sigma_Z} \leq t}_{\text{see (3.10)}} + \sqrt{n}\frac{\bar{Z}}{\sigma_Z} \leq t\} \\
&\rightarrow \Phi(t)
\end{aligned}$$

Moreover, by Slutsky's theorem

$$W_n \xrightarrow{\mathcal{D}} W \text{ and } T_n \xrightarrow{\mathcal{D}} b \text{ imply that } W_n T_n \xrightarrow{\mathcal{D}} Wb,$$

$$\text{letting } W_n = \sqrt{n}(\hat{\rho} - \rho_o)/\tau_n \text{ and } T_n = \tau_n/\tau,$$

we have $P\{\sqrt{n}(\hat{\rho} - \rho_o)/\tau \leq t\} \rightarrow \Phi(t)$ or $\sqrt{n}(\hat{\rho} - \rho_o) \xrightarrow{\mathcal{D}} N(0, \tau^2)$, and the claim is proved.

Step 3 (verify (3.7) and (3.8), and conclusion). Verify that

$$\begin{aligned}
E(Z) &= E\left(\frac{\partial \log p(X, U, V, \rho)}{\partial \rho} \Big|_{\rho=\rho_o}\right) \\
&= E\left(\frac{\partial p(X, U, V, \rho)}{\partial \rho} \Big|_{\rho=\rho_o}\right) \\
&= E\left(E\left(\frac{\partial p(X, U, V, \rho)}{\partial \rho} \Big|_{\rho=\rho_o} \mid (U, V)\right)\right) \\
&= E\left\{\frac{\partial F(U; \rho)}{\partial \rho} \Big|_{\rho=\rho_o} F(U; \rho_o) + \frac{\partial S(V; \rho)}{\partial \rho} \Big|_{\rho=\rho_o} S(V; \rho_o)\right. \\
&\quad \left. + \frac{\partial(S(U; \rho) - S(V; \rho))}{\partial \rho} \Big|_{\rho=\rho_o} (S(U; \rho_o) - S(V; \rho_o))\right\} \\
&= E\left[\frac{\partial F(U; \rho)}{\partial \rho} \Big|_{\rho=\rho_o} + \frac{\partial S(V; \rho)}{\partial \rho} \Big|_{\rho=\rho_o} + \frac{\partial(S(U; \rho) - S(V; \rho))}{\partial \rho} \Big|_{\rho=\rho_o}\right] \\
&= E\left[\frac{\partial(F(U; \rho) + S(V; \rho) + S(U; \rho) - S(V; \rho))}{\partial \rho} \Big|_{\rho=\rho_o}\right] \\
&= E\left(\frac{\partial 1}{\partial \rho} \Big|_{\rho_o}\right) \\
&= 0.
\end{aligned}$$

Thus (3.7) holds. Verify that (by (3.1) and formula $\frac{x}{1-x} = -1 + \frac{1}{1-x}$),

$$\begin{aligned}
-\frac{\partial^2 l(\rho)}{\partial \rho^2} \Big|_{\rho^*} &= \frac{1}{n} \sum_{i=1}^n \frac{-U_i^2}{1 - e^{-\rho U_i}} 1_{(X_i \leq U_i)} \Big|_{\rho^*} + \frac{1}{n} \sum_{i=1}^n \frac{U_i^2}{(1 - e^{-\rho U_i})^2} 1_{(X_i \leq U_i)} \Big|_{\rho^*} \quad (3.12) \\
&\quad + \frac{1}{n} \sum_{i=1}^n \frac{-(V_i - U_i)^2}{1 - e^{-\rho(V_i - U_i)}} 1_{(U_i < X_i \leq V_i)} \Big|_{\rho^*} + \frac{1}{n} \sum_{i=1}^n \frac{(V_i - U_i)^2}{(1 - e^{-\rho(V_i - U_i)})^2} 1_{(U_i < X_i \leq V_i)} \Big|_{\rho^*}
\end{aligned}$$

Applying the technique in proving Eq. (2.9) in the proof of Theorem 2 to the four summands in (3.12), we can show

$$\frac{\partial^2 l(\rho)}{\partial \rho^2} \Big|_{\rho^*} \rightarrow E\left(\frac{\partial^2 l(\rho_o)}{\partial \rho^2}\right) \text{ a.s..} \quad (3.13)$$

Thus (3.8) holds. It is worth mentioning that in (3.13),

$$\frac{\partial^2 l}{\partial \rho^2} = \int \int \int \frac{\partial^2}{\partial \rho^2} \ln p(x, u, v, \rho) dF_n(x, u, v) \text{ and}$$

$$E\left(\frac{\partial^2 l}{\partial \rho^2}\right) = \int \int \int \frac{\partial^2}{\partial \rho^2} \ln p(x, u, v, \rho) dF_o(x, u, v),$$

where F_n is the edf of F_o and F_o is the cdf of (X, U, V) .

It is easy to show that

$$E\left(-\frac{\partial^2 l(\rho_o)}{\partial \rho^2}\right) = E\left(\frac{U^2 e^{-\rho_o U}}{(1 - e^{-\rho_o U})^2} 1_{(X \leq U)}\right) + E\left(\frac{(V - U)^2 e^{-\rho_o(V-U)}}{(1 - e^{-\rho_o(V-U)})^2} 1_{(U < X \leq V)}\right).$$

Finally, verify that the Fisher information number

$$\sigma_Z^2 = E\left(\left(\frac{\partial \log p(X, U, V, \rho_o)}{\partial \rho}\right)^2\right) = -E\left(\frac{\partial^2 l(\rho_o)}{\partial \rho^2}\right) = 1/\tau^2. \quad (3.14)$$

Thus

$$\sqrt{n}(\hat{\rho} - \rho_o) \xrightarrow{\mathcal{D}} N(0, \tau^2). \quad \square$$

Comment. This proof can be replaced by the general theory, *e.g.* Cramér's theorem. However, we still need to verify conditions required by the theory. In particular, Cramér's theorem requires that there is a function $K(x, u, v)$ such that $E_{\rho_o}(K(X, U, V)) < \infty$ and $\left|\frac{d^2 l(\rho; x, u, v)}{d\rho^2}\right|$ is bounded by $K(x, u, v)$ uniformly in some neighborhood of ρ_o .

§3.5.4. Homework:

- (1) Prove (3.13).
- (2) Prove Equation (3.14).
- (3) Check the existence of the MLE of ρ of $Exp(\rho)$ under the DC model
- (4) Under the assumption in Theorem 2, compute $\mu(\rho)$ when $(U, V) = (i, i+2)$ w.p.1/2, $i = 1, 2$.

Chapter 4. Univariate nonparametric estimation

§4.1. Introduction.

Suppose that

the failure time $X \sim F_o$ (cdf),

(X_i, L_i, R_i) , $i = 1, \dots, n$ are i.i.d. from an extended random vector (X, L, R) .

Question: Observed (L_i, R_i) s,

$$F_o = ?$$

This is called a nonparametric estimation problem.

The parameter space can be viewed as

$$\Theta_o = \{F : F \text{ is a cdf}\}.$$

However, it is more convenient to viewed the parameter space as

$$\Theta = \{F : F(t) \uparrow, F : [-\infty, \infty] \rightarrow [0, 1], F(-\infty) = 0 \text{ and } F(\infty) = 1. \}. \quad (1.1)$$

Define an interval $I_i = \begin{cases} (L_i, R_i] & \text{if } L_i < R_i \\ [L_i, R_i] & \text{if } L_i = R_i. \end{cases}$

Let $\mu_F(\cdot)$ be the measure induced by F such that

$$\mu_F(I_i) = \begin{cases} F(R_i) - F(L_i) & \text{if } L_i < R_i \\ F(R_i) - F(L_i-) & \text{if } L_i = R_i. \end{cases}$$

Definition. The nonparametric likelihood function based on data (L_i, R_i) , $i = 1, \dots, n$, is

$$\mathbb{L}(F) = \prod_{i=1}^n \mu_F(I_i), \quad F \in \Theta.$$

The above definition is the same as the parametric definition, except that f is replaced by $f(t) = F(t) - F(t-)$. In other words, it is assumed that F is discrete, though F_o maybe continuous.

Definition. The generalized (or nonparametric) maximum likelihood estimator (GMLE or NPMLE), of F_o is an $F = \hat{F} \in \Theta$ such that

$$\hat{F} \text{ maximizes } \mathbb{L}(F) \text{ over } \Theta.$$

Remark.

1. The GMLE is also called the nonparametric MLE.
2. If F is discrete, $\mathbb{L}(F)$ is the same as the definition of parametric likelihood $\mathbb{L}(\phi)$ in §3.2.

§4.1.2. Homework.

Prove the following statement: If the data is complete, but X_i 's are not necessarily distinct (there could be ties), then the GMLE of F_o is given by

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(X_i \leq x)}.$$

\hat{F} above is called the empirical distribution function (edf).

Properties of the edf:

1. Given observations, it is a discrete cdf.
2. It gives equal weight $\frac{1}{n}$ to each observation.
3. The df at t is $\hat{f}(t) = \frac{\sum_{i=1}^n \mathbf{1}_{(X_i=t)}}{n}$ and $\hat{F}(t) = \sum_{x \leq t} \hat{f}(x)$.
4. $\bar{X} = \sum_t t \hat{f}(t)$.

§4.2. The RC model and the product-limit estimator (PLE).

Assume the RC model, i.e.,

- $Y \sim G$,
- X and Y are independent,
- observe $(Z, \delta) = (X \wedge Y, \mathbf{1}_{(X \leq Y)})$. Then

$$\begin{aligned} \mathcal{L}(F) &= \log \mathbb{L}(F) = \log \prod_{i=1}^n \mu_F(I_i) \\ &= \sum_{i: ex} \log f(Z_i) + \sum_{i: rc} \log S(Z_i), \text{ where } f(x) = F(x) - F(x-). \end{aligned}$$

The GMLE of F_o under the RC model is $\hat{F}_{pl} = 1 - \hat{S}_{pl}$, where

$$\hat{S}_{pl}(t) = \prod_{t \geq Z_{(i)}} \left(1 - \frac{1}{n - i + 1}\right)^{\delta_{(i)}},$$

where $Z_{(1)} \leq \dots \leq Z_{(n)}$ are order statistics of Z_i s and

$\delta_{(i)}$ is the δ_j associated with $Z_{(i)}$.

Here we use the convention that $x < x+$ in our ordering.

The GMLE is called the PLE or Kaplan-Meier estimator (Kaplan and Meier (1958)).

4.2.1. Understanding the PLE.

1. Redistribution to the right algorithm.

At each time point t , the PLE redistributes the weight $\hat{S}_{pl}(t-)$

equally to each of the observations to the right of t including t .

In particular, let $a_1 < \dots < a_m$ be all the distinct points among uncensored Z_1, \dots, Z_n , the PLE only puts weights on these a_k s. Let

$$d_k = \sum_{i=1}^n \mathbf{1}_{(L_i=R_i=a_k)}, \# \text{ of deaths at } a_k;$$

$$r_k = \sum_{i=1}^n \mathbf{1}_{(Z_i \geq a_k)}, \# \text{ of people in risk at } a_k;$$

- a. The PLE $\hat{S}_{pl}(t)$ is constant at $[a_{k-1}, a_k)$, $k = 1, \dots, m$, where $a_0 = -\infty$.
 - b. For $t < a_1$, $\hat{S}_{pl}(t) = 1$.
 - c. For $t = a_1$, \hat{S}_{pl} distributes the mass 1 equally to each of the r_1 observations to the right of a_1 (including a_1).
Total r_1 in risk, thus each with probability $\frac{1}{r_1}$.
Since d_1 deaths at a_1 ,
 $\Rightarrow \hat{S}_{pl} = 1 - \frac{d_1}{r_1}$ ($= \hat{P}(X > a_1)$).
 - d. At a_k , total of $\hat{S}_{pl}(a_{k-1})$ mass remains on $[a_k, +\infty)$.
 \hat{S}_{pl} redistributes the mass $\hat{S}_{pl}(a_{k-1})$ equally to each observation in risk to the right, thus each has probability $\frac{\hat{S}_{pl}(a_{k-1})}{r_k}$
 $\Rightarrow \hat{S}_{pl}(a_k) = \hat{S}_{pl}(a_{k-1})(1 - \frac{d_k}{r_k})$.
2. The above algorithm results in the expression

$$\hat{S}_{pl}(t) = \prod_{k: t \geq a_k} (1 - \frac{d_k}{r_k}).$$

Note that

$$\hat{S}_{pl}(t) = \hat{S}_{pl}(a_m), \quad t > a_m. \quad (2.1)$$

Thus $\hat{F}_{pl} = 1 - \hat{S}_{pl}$ may not be a cdf. In particular,

$$\lim_{t \rightarrow \infty} \hat{F}_{pl}(t) = \hat{F}_{pl}(Z_{(n)}) < 1 \text{ if } \delta_{(n)} = 0. \quad (2.2)$$

We define

$$\hat{S}_{pl}(+\infty) = 0,$$

so that $\hat{F}_{pl} \in \Theta$, where $\hat{F}_{pl} = 1 - \hat{S}_{pl}$.

It means that the PLE puts weight $\hat{S}_{pl}(Z_{(n)})$ at ∞ .

3. The PLE $\hat{F}_{pl}(t)$ is nondecreasing in t , but may not be a proper cdf as (2.2) may hold. Some people use the convention that

$$\hat{S}_{pl}(t) = 0 \begin{cases} \text{if } t > Z_{(n)} \\ \text{if } t \geq Z_{(n)} \\ \text{if } t \geq Z_{(n)} + c \text{ where } c > 0. \end{cases}$$

However, it can be shown that only the last definition has optimal asymptotic properties.

The following table calculates the PLE using the Leukaemia data group 0, as in Table 1.1

of §1.2.

Remission Time	Reverse Order (K)	$(1 - \frac{1}{n-i+1})^{\delta_{(i)}}$	$(1 - \frac{d_k}{r_k})$	$\hat{S}_{pl}(a_k)$
6	21	20/21		
6	20	19/20		
6	19	18/19	18/21	18/21
6+	18	1		18/21
7	17	16/17	16/17	$\frac{18}{21} \cdot \frac{16}{17}$
9+	16	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{16}{16}$
10	15	14/15	14/15	$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15}$
10+	14	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{14}{14}$
11+	13	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{14}{14} \cdot \frac{14}{14}$
13	12	11/12	11/12	$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{11}{12}$
16	11	10/11	10/11	$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{10}{11}$
17+	10	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{10}{11} \cdot \frac{10}{10}$
19+	9	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{10}{11} \cdot \frac{10}{10} \cdot \frac{10}{10}$
20+	8	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{10}{11} \cdot \frac{10}{10} \cdot \frac{10}{10} \cdot \frac{10}{10}$
22	7	6/7	6/7	$\frac{6 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{21 \cdot 17 \cdot 15 \cdot 11 \cdot 12}$
23	6	5/6	5/6	$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{21 \cdot 17 \cdot 15 \cdot 11 \cdot 12}$
25+	5	1		$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{21 \cdot 17 \cdot 15 \cdot 11 \cdot 12} \cdot \frac{18}{18}$
32+	4	1		$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{21 \cdot 17 \cdot 15 \cdot 11 \cdot 12} \cdot \frac{18}{18} \cdot \frac{18}{18}$
32+	3	1		$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{21 \cdot 17 \cdot 15 \cdot 11 \cdot 12} \cdot \frac{18}{18} \cdot \frac{18}{18} \cdot \frac{18}{18}$
34+	2	1		$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{21 \cdot 17 \cdot 15 \cdot 11 \cdot 12} \cdot \frac{18}{18} \cdot \frac{18}{18} \cdot \frac{18}{18} \cdot \frac{18}{18}$
35+	1	1		$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{21 \cdot 17 \cdot 15 \cdot 11 \cdot 12} \cdot \frac{18}{18} \cdot \frac{18}{18} \cdot \frac{18}{18} \cdot \frac{18}{18} \cdot \frac{18}{18}$

$$\hat{S}_{pl}(t) = \begin{cases} 1 & \text{if } t < 6, \\ 18/21 & \text{if } t \in [6, 7) \\ 96/119 & \text{if } t \in [7, 10) \\ \dots & \dots \end{cases}$$

§4.2.1.2. Homework:

1. Apply the redistribution-to-the-right (RTR) method to the last 13 data from Table 1 and compute by hand the PLE. (Now $n = 13$.)
2. Suppose that $X \sim Bin(3, 1/3)$, $Y \sim Bin(1, 0.4)$. There are two observations under the RC model. Compute the mean of $\hat{S}_{pl}(t)$ for $t \leq 1$.

Remark: Since we do not have observation beyond 1, we do not expect that we can make decent inference on $S_o(t)$ for $t > 1$. However, we can estimate S_o for $t \leq 1$.

3. In # 2, is the PLE of S_o an unbiased estimator for $t \leq 1$?

Theorem 1. (Johanson (1978)). *The PLE \hat{S}_{pl} is a GMLE of $S_o (= 1 - F_o)$, that is, it maximizes $\prod_{i=1}^n (S(Z_i -) - S(Z_i))^{\delta_i} (S(Z_i))^{1 - \delta_i}$, $F = 1 - S \in \Theta$.*

We only prove the theorem under a special case.

Suppose that there are only 3 distinct Z_i s say $a_1 < a_2 < a_3$, among n observations. Denote

$c_k = \sum_{i=1}^n \mathbf{1}_{(X_i \neq Z_i = a_k)} \geq 1$, # of people censored at a_k (may all be right censored),

$d_k = \sum_{i=1}^n \mathbf{1}_{(Z_i = X_i = a_k)}$ and

$r_k = \sum_{i=1}^n \mathbf{1}_{(Z_i \geq a_k)}$. Then

$[-\infty, a_1], [a_1, a_1], (a_1, a_2), [a_2, a_2], (a_2, a_3), [a_3, a_3]$ and $(a_3, \infty]$ is a partition of $[-\infty, +\infty]$, with the measure assigned by an $F \in \Theta$ to these intervals being p_1, \dots, p_7 ($\sum_{i=1}^7 p_i = 1$).

$$\begin{aligned} \mathbb{L}(F) &= \prod_{k=1}^3 (\mu_F([a_k, a_k])^{\sum_{i=1}^n \mathbf{1}_{(Z_i = X_i = a_k)}} \prod_{k=1}^3 (\mu_F((a_k, +\infty]))^{\sum_{i=1}^n \mathbf{1}_{(X_i \neq Z_i = a_k)}} \\ &= \prod_{k=1}^3 p_{2k}^{d_k} \prod_{k=1}^3 (p_{2k+1} + \dots + p_7)^{c_k} \end{aligned}$$

$$\begin{aligned}
&= \prod_{k=1}^3 p_{2k}^{d_k} (p_3 + p_4 + p_5 + p_6 + p_7)^{c_1} (p_5 + p_6 + p_7)^{c_2} (p_7)^{c_3} \\
&\leq \prod_{k=1}^3 s_{2k}^{d_k} (0 + s_4 + 0 + s_6 + s_7)^{c_1} (0 + s_6 + s_7)^{c_2} (s_7)^{c_3} \quad (\text{if } p_1 + p_2 = s_2, \\
&\quad p_3 + p_4 = s_4, p_5 + p_6 = s_6, p_7 = s_7 \text{ and } s_i = 0, i = 1, 3, 5) \\
&= \left(\prod_{k=1}^3 s_{2k}^{d_k} \right) (s_4 + s_6 + s_7)^{c_1} (s_6 + s_7)^{c_2} s_7^{c_3} \quad (\text{note } s_2 + s_4 + s_6 + s_7 = 1) \\
&\quad (\text{it can be replaced by } \prod_{k=1}^m \dots)
\end{aligned}$$

$$\begin{aligned}
\mathbf{L}(F) &\leq [s_2^{d_1} (s_4 + s_6 + s_7)^{c_1}] [s_4^{d_2} (s_6 + s_7)^{c_2}] [s_6^{d_3} s_7^{c_3}] \\
&= s_2^{d_1} (1 - s_2)^{c_1} (s_4)^{d_2} (s_6 + s_7)^{c_2} (s_6)^{d_3} (s_7)^{c_3} \\
&= s_2^{d_1} (1 - s_2)^{c_1 + d_2 + c_2 + d_3 + c_3} \left(\frac{s_4}{1 - s_2} \right)^{d_2} \left(\frac{s_6 + s_7}{1 - s_2} \right)^{c_2} \left(\frac{s_6}{1 - s_2} \right)^{d_3} \left(\frac{s_7}{1 - s_2} \right)^{c_3} \\
&= s_2^{d_1} (1 - s_2)^{c_1 + d_2 + c_2 + d_3 + c_3} \left(\frac{s_4}{1 - s_2} \right)^{d_2} \left(1 - \frac{s_4}{1 - s_2} \right)^{c_2} \left(\frac{s_6}{1 - s_2} \right)^{d_3} \left(\frac{s_7}{1 - s_2} \right)^{c_3} \\
&= s_2^{d_1} (1 - s_2)^{c_1 + d_2 + c_2 + d_3 + c_3} \left(\frac{s_4}{1 - s_2} \right)^{d_2} \left(1 - \frac{s_4}{1 - s_2} \right)^{c_2} \\
&\quad \times \left(\frac{1 - s_2 - s_4}{1 - s_2} \right)^{d_3 + c_3} \left(\frac{s_6}{1 - s_2 - s_4} \right)^{d_3} \left(\frac{s_7}{1 - s_2 - s_4} \right)^{c_3} \\
&= s_2^{d_1} (1 - s_2)^{c_1 + d_2 + c_2 + d_3 + c_3} \left(\frac{s_4}{1 - s_2} \right)^{d_2} \left(1 - \frac{s_4}{1 - s_2} \right)^{c_2 + d_3 + c_3} \\
&\quad \times \left(\frac{s_6}{1 - s_2 - s_4} \right)^{d_3} \left(1 - \frac{s_6}{1 - s_2 - s_4} \right)^{r_3 - d_3} \\
&= s_2^{d_1} (1 - s_2)^{r_1 - d_1} \left(\frac{s_4}{1 - s_2} \right)^{d_2} \left(1 - \frac{s_4}{1 - s_2} \right)^{r_2 - d_2} \\
&\quad \times \left(\frac{s_6}{1 - s_2 - s_4} \right)^{d_3} \left(1 - \frac{s_6}{1 - s_2 - s_4} \right)^{r_3 - d_3} \tag{2.3}
\end{aligned}$$

(as $r_1 = n$, $r_2 = d_2 + c_2 + d_3 + c_3$ and $r_3 = d_3 + c_3$),

where $s_2 + s_4 + s_6 + s_7 = 1$ and $s_i \geq 0$. Verify that

(a) (2.3) equals $(s_2^*)^{d_1} (1 - s_2^*)^{r_1 - d_1} (s_4^*)^{d_2} (1 - s_4^*)^{r_2 - d_2} (s_6^*)^{d_3} (1 - s_6^*)^{r_3 - d_3}$,
where s_i^* 's are defined in an obvious way, satisfying $s_i^* \in [0, 1]$.

(b) The transformation from (s_2, s_4, s_6) to (s_2^*, s_4^*, s_6^*) is one-to-one and onto.

Thus, $\mathbf{L}(F)$ is maximized by setting $s_{2i}^* = d_i / r_i$ for $i = 1, 2, 3$.

step	$(-\infty, a_1)$	$\{a_1\}$	(a_1, a_2)	$\{a_2\}$	(a_2, a_3)	$\{a_3\}$	(a_3, ∞)	
1	p_1	p_2	p_3	p_4	p_5	p_6	p_7	$\sum_i p_i = 1$
2	0	s_2	0	s_4	0	s_6	s_7	$\sum_i s_i = 1$
3	0	s_2^*	0	s_4^*	0	s_6^*	s_7^*	$s_{2k}^* \in (0, 1)$

Relation between s_i and s_i^* :

$$\begin{aligned}
s_2 &= s_2^*, \\
\frac{s_4}{1 - s_2} &= s_4^*, \\
\frac{s_6}{1 - s_2 - s_4} &= s_6^*.
\end{aligned}$$

It follows from the relation between s_i and s_i^* that

$\mathbf{L}(F)$ is maximized by

$$\begin{aligned}
s_2 &= d_1 / r_1 = \frac{d_1}{n}, \\
s_4 &= (1 - s_2) s_4^* = \frac{r_1 - d_1}{r_1} \frac{d_2}{r_2}, \\
s_6 &= (1 - s_2 - s_4) s_6^* = \dots
\end{aligned}$$

Consequently, $\hat{S}_{pl}(a_1) = \frac{r_1 - d_1}{r_1}$, $\hat{S}_{pl}(a_2) = \frac{r_1 - d_1}{r_1} \frac{r_2 - d_2}{r_2}$, $\hat{S}_{pl}(a_3) = \frac{r_1 - d_1}{r_1} \frac{r_2 - d_2}{r_2} \frac{r_3 - d_3}{r_3}$.

$$\hat{S}_{pl}(t) = \hat{S}_{pl}(a_i) \text{ if } t \in [a_i, a_{i+1}), i = 0, \dots, 3. \quad a_0 = ? \quad a_4 = ?$$

§4.2.1.3. Homework:

4. Extend the proof of the GMLE from the case $a_m = a_3$ to the general case by induction on the number of distinct Z_i s. Notice that $d + k + c_k \geq 1$, but $c_k = 0$ is possible now.

§4.2.2. Properties of the PLE \hat{F} .

We shall first state the main results on the properties of the PLE and present some simpler proofs under the assumption that the random variables take on finitely many values. Let $\tau = \sup\{t : F_Z(t) < 1\}$, where $Z = X \wedge Y$, and $D_Z^* = \begin{cases} \{t : t \leq \tau\} & \text{if } P(X = \tau \leq Y) > 0 \\ \{t : t < \tau\} & \text{otherwise.} \end{cases}$

Theorem 1. (Yu, Ai and Yu (2012)) *Suppose that either under the RC model, or the assumption $X \perp Y$ in the RC model is weakened by the next two assumptions:*

- (1) *Given r , $\exists G_1(r)$ such that $F_{Y|X}(r|t) = G_1(r)$ a.e. in t on (r, ∞) (w.r.t. μ_{F_o}).*
- (2) *$G_1(\cdot)$ does not depend on $F_o(\cdot)$.*

Then $\sup_{t \in D_Z^} |\hat{F}(t) - F_o(t)| \xrightarrow{a.s.} 0$ and $\sup_t |\hat{F}(t) - F_*(t)| \xrightarrow{a.s.} 0$ where*

$$F_*(t) = \begin{cases} F_o(t) & \text{if } t \in D_Z^* \cup \{\infty\} \\ F_o(\tau) & \text{if } t \in (\tau, \infty) \text{ and } P(X = \tau \leq Y) > 0. \\ F_o(\tau-) & \text{if } t \in [\tau, \infty) \text{ and } P(X = \tau \leq Y) = 0. \end{cases}$$

Note. $F_o(t)$, $t > \tau$, is not estimable as there is no observation beyond τ , unless $F_o(\tau-) = 1$.

Examples that G depends on F_o :

- (1) $G = 1 - (S_o)^r$, where $r > 0$;
 - (2) $G(t) = [G_o(t) + F_o(t)]/2$, where G_o is a cdf.
- For clarification, two instances of discrete $f_{Y|X}$ are given as follows:

$$\begin{array}{l} \text{case (1)} \quad \begin{array}{c} t \text{ value :} \\ f_{Y|X}(1|t) \\ f_{Y|X}(2|t) \\ f_{Y|X}(3|t) \end{array} \begin{array}{cc} 2 & 3 \\ \left(\begin{array}{cc} 1/3 & 1/2 \\ 1/3 & 1/6 \\ 1/3 & 1/3 \end{array} \right) \end{array} \text{ and case (2)} \quad \begin{array}{c} t \text{ value :} \\ f_{Y|X}(1|t) \\ f_{Y|X}(2|t) \\ f_{Y|X}(3|t) \end{array} \begin{array}{ccc} 2 & 3 & 4 \\ \left(\begin{array}{ccc} 1/5 & 1/5 & 1/5 \\ 1/5 & 3/5 & 3/5 \\ 3/5 & 1/5 & 1/5 \end{array} \right) \end{array} \end{array}$$

4.2.2.1. Homework. Verify that case (2) satisfies assumptions (1) and (2) but not case (1) and the PLE is not consistent in case (1).

Several weaker results on the consistency were established earlier by Peterson (1977), Phadia and Van Ryzin (1980), Shorack and Wellner (1986), Wang (1987), Stute and Wang (1993) and Yu and Li (1994), among others.

In particular, Under the standard RC model,

- * Peterson (1977), Phadia and Van Ryzin (1980), Shorack and Wellner (1986) showed that the PLE $\hat{S}(t)$ is consistent if $t < \tau$ and if S_o is discrete, or S_o is continuous;
- * Wang (1987) showed that the PLE $\hat{S}(t)$ is consistent if $t \leq Z_{(n)}$;
- * Stute and Wang (1993) showed that the PLE $\hat{S}(t)$ is consistent in the set

$$D_Z = \begin{cases} (-\infty, \tau] & \text{if } P(Y \geq \tau) > 0 \text{ or } P(X = \tau) = 0 \\ (-\infty, \tau) & \text{otherwise.} \end{cases}$$

but F_X and F_Y do not have jumps in common;

- * Yu and Li (1994) show that the PLE $\hat{S}(t)$ is consistent in \mathcal{D}_Z .

Remark. If $X \perp Y$, then $P(X = \tau \leq Y) = P(X = \tau)P(\tau \leq Y)$.

Otherwise, $P(X = \tau \leq Y) \neq P(X = \tau)P(\tau \leq Y)$.

Theorem 2. (Breslow and Crowley (1974), Gill (1983), Gu and Zhang (1993), Stute (1995), Yu and Hsu (2015)). *Suppose that the assumptions in Theorem 1 all holds,*

$$U_n(t) = \sqrt{n} \left(\frac{\hat{S}_{pl}(t) - S_*(t)}{S_*(t)} \right) \xrightarrow{D} N(0, \sigma^2) \text{ for } t < \tau, .$$

The asymptotic covariance of $\hat{S}_{pl}(t)$ and $\hat{S}_{pl}(s)$ is

$$nCov(\hat{S}_{pl}(t), \hat{S}_{pl}(s)) \approx S_*(t)S_*(s) \int_0^{t \wedge s} \frac{1}{S_*(x-)S_{Y|X}(x - |\tau)S_*(x)} dF_*(x), \quad t, s < \tau,$$

where $S_{Y|X}(y|x) = P(Y > y|X = x)$. *The above two statements also hold $\forall t, s < \infty$, iff either (1) $\tau = \infty$, or (2) $S_o(\tau-) > 0$ or (3) $S_o(\tau-) = 0$, $\tau < \infty$ and $\sigma_\tau = 0$, where*

$$\sigma_t^2 = \frac{1}{n} (S_*(t))^2 \int_0^t \frac{1}{S_Z(x-)S_*(x)} dF_*(x).$$

It follows that $\sigma_{\hat{S}_{pl}(t)}^2 \approx \sigma_t^2/n$, which can be estimated by

$$\begin{aligned} n(\hat{\sigma}_{\hat{S}_{pl}(t)})^2 &= (\hat{S}_{pl}(t))^2 \int_0^t \frac{1}{\hat{S}_Z(x-)\hat{S}_{pl}(x)} d\hat{F}_{pl}(x) \\ &= (\hat{S}_{pl}(t))^2 \sum_{k: a_k \leq t} \frac{\hat{F}_{pl}(a_k) - \hat{F}_{pl}(a_k-)}{\hat{S}_Z(a_k-)\hat{S}_{pl}(a_k)} \quad (a_k \text{'s are distinct exact observations}) \\ &= (\hat{S}_{pl}(t))^2 \sum_{k: a_k \leq t} \frac{\hat{f}(a_k)}{\hat{S}_Z(a_k-)\hat{S}_{pl}(a_k)}, \end{aligned}$$

Note

$$\begin{aligned} \hat{S}_{pl}(t) &= \prod_{i: Z_{(i)} \leq t} \left(1 - \frac{1}{n-i+1}\right)^{\delta_{(i)}} = \prod_{i: Z_{(i)} \leq t} \left(1 - \frac{\delta_{(i)}}{n-i+1}\right). \\ \hat{S}_Z(t) = \hat{S}_Y(t)\hat{S}_{pl}(t) &= \prod_{i: Z_{(i)} \leq t} \left(1 - \frac{1}{n-i+1}\right) = \sum_{i=1}^n \frac{\mathbf{1}_{\{Z_i > t\}}}{n}. \end{aligned}$$

The GMLE of S_Y is also a PLE of S_Y .

$$\hat{S}_Y(t) = \prod_{i: Z_{(i)} \leq t} \left(1 - \frac{1}{n-i+1}\right)^{1-\delta_{(i)}} = \prod_{i: Z_{(i)} \leq t} \left(1 - \frac{1-\delta_{(i)}}{n-i+1}\right).$$

Explain using RTR method.

R Commands:

```
b=exp(6)
g=1/3
y=rweibull(100,g,b)
c=runif(100,0,780)
d=as.numeric(y<=c)
# d=ifelse(y>c, 0,1)
m=y*d+c*(1-d)
aa=survfit(Surv(m,d)~1)
plot(aa)
aa$time
aa$surv
summary(aa)
summary(aa)$time
summary(aa)$surv
  time      n.risk  n.event  survival  std.err  lower95%CI  upper95%CI
1.17e-04    100      1      0.990    0.00995    0.971      1.000
2.62e-02     99      1      0.980    0.01400    0.953      1.000
7.44e-02     98      1      0.970    0.01706    0.937      1.000
```

§4.2.2.2. Homework:

1. Suppose that $X \sim \text{Bin}(3, 1/3)$, $Y \sim \text{Bin}(1, 0.4)$. Under the RC model, what is the limit of $\hat{S}_{pl}(2)$?
2. Suppose $X \sim \text{Exp}(\rho)$, where $E(X) = 1/\rho$, and $Y \sim U(0, 4)$. Select a ρ_o . Generate 100 RC data, plot, on the same figure, (1) the PLE \hat{S}_{pl} , (2) the true survival function $S(\cdot, \rho_o)$ and (3) the $S(\cdot, \hat{\rho})$, where $\hat{\rho}$ is the parametric MLE of ρ . **Moreover, add (4) the parametric MLE of the survival function of a normal distribution $N(\mu, 1)$ using Monte Carlo method (or `survreg(dist="gaussian")`) and (5) the parametric MLE of the survival function of a $U(\theta, 5)$.** Make comments on their deviations from the true $S(t, \rho)$ at each time point t .

3. Using the data in problem #2, we can derive the confidence intervals for $S_o(t)$ based on $\hat{S}_{pl}(t)$ and based on the other two estimates of S_o . They are

$$(\hat{S}(t) - z_{\alpha/2}\hat{\sigma}_{\hat{S}(t)}, \hat{S}(t) + z_{\alpha/2}\hat{\sigma}_{\hat{S}(t)})$$

and

$$(S(t, \hat{\rho} + z_{\alpha/2}\hat{\sigma}_{\hat{\rho}}), S(t, \hat{\rho} - z_{\alpha/2}\hat{\sigma}_{\hat{\rho}})),$$

etc., respectively, where $\Phi(-z_\alpha) = \alpha$ and Φ is the cdf of $N(0, 1)$ (Actually, for the $U(\theta, 5)$, one can use the Bootstrap method to get SD of the MLE). Both the ends are curves and each pair of the curves induced by the ends of the confidence interval is called a confidence band. Plot on the same figure the **three** confidence bands and the true survival function $S(\cdot, \rho)$.

R program for putting the two graphs on one paper.

```
par(mfrow=c(3,1))
x=(-35:45)/10
y=1-pnorm(1.65-x)
plot(x,y,type="l",lty=1 ,xlim=c(-3.5,4.5), ylim=c(0,1.0))
y=1-pnorm(1.65-2*x)
lines(x,y,lty=2)
y=1-pnorm(1.65-3*x)
lines(x,y,lty=3,col=1)
y=1-pnorm(1.65-4*x)
lines(x,y,lty=4,col=2)
leg.names<-c("x1", "x2", "x3", "x4")
legend(-2.5, 0.77, leg.names, lty=c(1,2,3,4),cex=1.5)
#-----
x=rexp(100,0.5)
mean(x)
[1] 1.503257
z=rnorm(200,1.5,3)
u=rnorm(200,1.5,3)
v=3*z+4*u
u[max(v)]
[1] 2.364634
z[max(v)]
[1] -0.1354150
```

Idea for Monte Carlo Method in computing MLE of μ assuming $N(\mu, 1)$ and based on some data:

1. Guess the range of μ and then generate 1000-10000 number in that range.
2. Compare the likelihoods with these numbers being μ

R code for Monte Carlo method:

```
x=rexp(100)
d=rbinom(100,1,0.9)
n=1000
m=runif(n,-1,4)
l=-10000
j=-1
for(i in 1:n){
y=m[i]
L=sum(d*log(dnorm(x,y,1)) + (1-d)*log(1-pnorm(x,y,1)))
if (L>l) {
l=L
j=i
}
```

}
m[j]

4. Suppose there are 6 right-censored observations (Z_i, δ_i) : (1,1), (2,0), (3,1), (2,1), (2,1), (4,1). Compute the PLE \hat{S}_X and \hat{S}_Y of S_X and S_Y respectively, on $[0, \infty)$. Show their product is the empirical survival function $\tilde{S}(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(t < Z_i)}$. That is, $\tilde{S}(t) = \hat{S}_X(t)\hat{S}_Y(t)$ for all t .
5. Suppose that (Z_i, δ_i) , $i = 1, \dots, 100$, is a random sample of right-censored data and suppose that the failure time X has an exponential distributions with the parameter ρ and the censoring time $Y = 2$ with probability 1. The PLE \hat{S}_{pl} and $\hat{S}_1 = e^{-\hat{\rho}t}$ are both estimators of S_X . Compute their asymptotic standard deviations explicitly. Which is bigger? **Hint: you may draw their graphs first.**
6. If you want to estimate a cdf F based on a random sample of right-censored data (Z_i, δ_i) , $i = 1, \dots, n$, and you believe the cdf is an exponential distribution, what is your estimator?
7. Suppose $X \sim Exp(\rho)$, where $E(X) = 1/\rho$, and $Y \equiv 2$. Do a simulation study as follows.
Select a ρ_o . Generate 3 RC data. Compute the PLE of $S(t)$ and its Variance. Construct a CI for $S(t)$ based on the PLE. with the significance level you chosen and the 3 data you generated. Notice that your CI is a statistic, depends only on the data.

We shall give a proof of Theorems 1 and 2 under the assumption that $X \perp Y$ and Z takes on finitely many values, with the largest values of Y is τ .
Let $a_1 < \dots < a_m$ be all the possible values of $X \leq \tau$ and $\tau = a_m$.
Assume that Y can only take values among a_j s.
Let d_k, c_k and r_k be defined as before (though these a_k 's are defined differently from before).
The likelihood function becomes

$$L(\vec{p}) = \prod_{i=1}^m p_i^{d_i} \left(\sum_{j>i} p_j \right)^{c_i},$$

Then the problem reduces to a parametric problem of multinomial distribution

$$(Y_1, \dots, Y_{2m-1}) \sim M(n, c\theta_1, \dots, c\theta_{2m-1}),$$

where for $i = 1, \dots, m$, $Y_i = d_i$, $Y_{m+i} = c_i$, $\theta_i = p_i$, $\theta_{m+i} = \sum_{j>i} p_j$, $p_i = P\{X = a_i\}$ and $c = 1/\sum_{i=1}^{2m} \theta_i$. Here θ_i 's are function of $\vec{p} = (p_1, \dots, p_{m-1})$. The multinomial distribution belongs to the exponential family and its consistency and asymptotic normality can be proved easily by standard approach such as Cramér's theorem (see, e.g., Ferguson (1996)) However, it is not easy to verify the expression for the variance of \hat{S}_{pl} from the inverse of the Fisher information matrix. Thus we use a different approach.

The following consistency proof helps understanding the PLE.

A simple proof of Theorem 1. Under the given assumptions, we need to prove the statement as follows.

$$\sup_{t \leq a_m} |\hat{F}(t) - F_*(t)| \xrightarrow{a.s.} 0, \text{ where } F_*(t) = \sum_{i=1}^m F_o(a_i) \mathbf{1}(t \in [a_i, a_{i+1})).$$

The PLE can be written as

$$\hat{S}_{pl}(t) = \prod_{k: t \geq a_k} \left(1 - \frac{d_k}{r_k}\right), .$$

Since m is finite,

$$\frac{d_k}{n} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(Z_i = X_i = a_k)}$$

$$\begin{aligned} &\rightarrow E(\mathbf{1}_{(Z=X=a_k)}) \quad \text{a.s. by SLLN.} \\ &= P(X = a_k \leq Y) = P(X = a_k)P(Y \geq a_k), \end{aligned} \quad (2.4)$$

$$\begin{aligned} \frac{r_k}{n} &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(Z_i \geq a_k)} \\ &\rightarrow E(\mathbf{1}_{(Z \geq a_k)}) \quad \text{a.s. by SLLN.} \\ &= P(Y \geq a_k, X \geq a_k) = P(X \geq a_k)P(Y \geq a_k). \end{aligned} \quad (2.5)$$

Prove by induction on $k = 1, \dots, m$, the statement

$$\hat{S}_{pl}(t) \rightarrow S_o(t) \text{ for } t \in [a_{k-1}, a_k) \text{ a.s. where } a_0 = -\infty \text{ and } a_{m+1} = \infty. \quad (2.6)$$

(Case $k = 1$). $S_o(t) = 1$ and $\hat{S}_{pl}(t) = 1$ for $t < a_1$.

Thus $\hat{S}_{pl}(t) \rightarrow S_o(t)$ for $t \in [a_0, a_1)$ a.s..

(Case $k < m$). Assume that (2.6) holds.

(Case $k + 1$). For $t \in [a_k, a_{k+1})$,

$$\begin{aligned} \hat{S}_{pl}(t) &= \hat{S}_{pl}(a_k) \\ &= \prod_{j: a_k \geq a_j} \left(1 - \frac{d_j}{r_j}\right) \\ &= \prod_{j: j < k} \left(1 - \frac{d_j}{r_j}\right) \left(1 - \frac{d_k}{r_k}\right) \\ &= \hat{S}_{pl}(a_{k-1}) \left(1 - \frac{d_k}{r_k}\right) \\ &= \hat{S}_{pl}(a_{k-1}) \left(1 - \frac{d_k/n}{r_k/n}\right) \\ &\rightarrow S_o(a_{k-1}) \left(1 - \frac{P(X = a_k)P(Y \geq a_k)}{P(X \geq a_k)P(Y \geq a_k)}\right) \text{ a.s.} \\ &\quad \text{(by induction assumption on } k \text{ and by (2.4) and (2.5))} \\ &= S_o(a_{k-1}) \left(1 - \frac{P(X = a_k)}{P(X \geq a_k)}\right) \\ &= P(X > a_{k-1}) \frac{P(X > a_k)}{P(X \geq a_k)} \\ &= P(X > a_k) = S_o(a_k) = S_o(t). \end{aligned} \quad (2.7)$$

Thus (2.6) holds for $k + 1$ as well.

This completes the induction proof.

Since S_o takes finitely many values, point-wise strong consistent implies uniform strong consistency. Thus it completes the proof of Theorem 1 under the finite assumption. \square

The following proof makes the expression of $\sigma_{\hat{S}_{pl}(t)}^2$ more explicit than the inverse of the Fisher information matrix.

Recall $\hat{S}_{pl}(t) = \prod_{k: a_k \leq t} \left(1 - \frac{d_k/n}{r_k/n}\right)$.

A simple proof of Theorem 2. We shall now give a proof under the simple assumption that $X \perp Y$, Y and X take on finitely many values before τ , $F_Y(\tau) = 1$ and $r_1 = n$. Assume $f(a_k) > 0$, $k = 1, \dots, m + 1$, where $f(a_{m+1}) = P(X > a_m)$. Then $\hat{S}_{pl}(t)$ is a function of d_1, \dots, d_m and r_1, \dots, r_m . Under these assumptions,

Theorem 2 becomes:

$$U_n(t) = \sqrt{n} \left(\frac{\hat{S}_{pl}(t) - S_o(t)}{S_o(t)} \right) \xrightarrow{D} N(0, \sigma^2) \text{ for } t < \tau, .$$

The asymptotic covariance of $\hat{S}_{pl}(t)$ and $\hat{S}_{pl}(s)$ is

$$nCov(\hat{S}_{pl}(t), \hat{S}_{pl}(s)) \approx S_o(t)S_o(s) \int_0^{t \wedge s} \frac{1}{S_o(x-)S_Y(x-)S_o(x)} dF_o(x), \quad t, s \leq \tau.$$

In particular, for $t = a_j$, $j < m$, we can write

$$\ln \hat{S}_{pl}(t) = g(d_1/n, r_1/n, \dots, d_j/n, r_j/n) = \sum_{i=1}^j \ln(1 - \frac{d_i/n}{r_i/n}).$$

That is

$$g(\mathbf{w}) = \ln(1 - \frac{w_1}{w_2}) + \dots + \ln(1 - \frac{w_{2j-1}}{w_{2j}}),$$

where $\mathbf{w}^t = (w_1, w_3, w_4, w_5, \dots, w_{2j})$,

as $w_2 = r_1/n = n/n = 1$ is a constant (note that $w_{2m-1}/w_{2m} = 1$ if $f(a_{m+1}) = 0$). Write

$$\bar{\mathbf{W}} = (\bar{W}_1, \bar{W}_3, \bar{W}_4, \dots, \bar{W}_{2j-1}, \bar{W}_{2j})^t = (d_1/n, d_2/n, r_2/n, \dots, d_j/n, r_j/n)^t,$$

as

$$\begin{aligned} \bar{W}_1 &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}(X_i = a_1 \leq Y_i), \\ \bar{W}_3 &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}(X_i = a_2 \leq Y_i), \\ \bar{W}_4 &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}(Z_i \geq a_2), \dots \end{aligned}$$

Then

$$g(\bar{\mathbf{W}}) = \ln \hat{S}_o(t) \text{ and } g(E(\bar{\mathbf{W}})) = \ln S_o(t) \text{ (by (2.7)).}$$

Then it follows from a corollary of Slutsky's theorem that

the $\ln \hat{S}_{pl}(t)$ is asymptotically normally distributed with an asymptotic variance

$$\sigma_{\ln \hat{S}_{pl}(t)}^2 \approx \frac{\partial g}{\partial \mathbf{w}^t} \Big|_{\mathbf{w}=E(\bar{\mathbf{W}})} \Sigma \frac{\partial g}{\partial \mathbf{w}} \Big|_{\mathbf{w}=E(\bar{\mathbf{W}})},$$

where Σ is the covariance matrix of $\bar{\mathbf{W}}$.

$$\sqrt{n}(g(\bar{\mathbf{W}}) - g(E(\bar{\mathbf{W}}))) \xrightarrow{D} N(0, \sigma^2).$$

$$\sqrt{n}(\ln \hat{S}_{pl}(t) - \ln S_o(t)) \xrightarrow{D} N(0, \sigma^2).$$

$$\sqrt{n}(h(\ln \hat{S}_{pl}(t)) - h(\ln S_o(t))) \xrightarrow{D} N(0, \sigma_2^2) \text{ where } h(t) = e^t.$$

$$\sqrt{n}(\hat{S}_{pl}(t) - S_o(t)) \xrightarrow{D} N(0, \sigma_2^2).$$

That is, $\hat{S}_{pl}(t)$ is asymptotically normally distributed with an asymptotic variance

$$\begin{aligned} \sigma_2^2/n &\approx \sigma_{\hat{S}_{pl}(t)}^2 \approx \frac{\partial e^x}{\partial x} \Big|_{x=\ln S_o(t)} \sigma_{\ln \hat{S}_{pl}(t)}^2 \frac{\partial e^x}{\partial x} \Big|_{x=\ln S_o(t)} \\ &\approx (S_o(t))^2 \sigma_{\ln \hat{S}_{pl}(t)}^2, \end{aligned}$$

as $\hat{S}_{pl}(t) = e^{\ln(\hat{S}_{pl}(t))}$.

For simplicity, for now, we let $j = 2$. Then

$$\mathbf{w}^t = (w_1, w_3, w_4),$$

$$g(\mathbf{w}) = \ln(1 - \frac{w_1}{w_2}) + \ln(1 - \frac{w_3}{w_4}), \text{ where } w_2 = 1.$$

$$\frac{\partial g}{\partial \mathbf{w}} = \left(\frac{1}{-w_2(1-\frac{w_1}{w_2})} \quad \frac{1}{-w_{2j}(1-\frac{w_{2j-1}}{w_{2j}})} \quad \frac{w_{2j-1}}{(1-\frac{w_{2j-1}}{w_{2j}})} \right)^t = \left(\frac{-1}{w_2-w_1} \quad \frac{-1}{w_{2j}-w_{2j-1}} \quad \frac{w_{2j-1}}{w_{2j}(w_{2j}-w_{2j-1})} \right)^t$$

$$\frac{\partial g}{\partial \mathbf{w}} \Big|_{\mathbf{w}=E(\bar{\mathbf{W}})} = \left(\frac{-1}{S_o(a_1)S_Y(a_1-)} \quad \frac{-1}{S_o(a_j)S_Y(a_j-)} \quad \frac{f_o(a_j)}{S_o(a_j-)S_o(a_j)S_Y(a_j-)} \right)^t \text{ (by (2.4) \& (2.5)),}$$

$$\sum = Cov(\bar{\mathbf{W}}) = \frac{1}{n} Cov(\mathbf{W}) = \frac{1}{n} Cov(\mathbf{1}_{(X=Z=a_1)}, \mathbf{1}_{(X=Z=a_2)}, \mathbf{1}_{(Z \geq a_2)}),$$

Denote

$$\frac{\partial g}{\partial \mathbf{w}} \Big|_{\mathbf{w}=E(\overline{\mathbf{W}})} = (p_1, p_2, p_3)^t \quad (= \mathbf{p}^t).$$

Denote

$$n\Sigma = (s_{ih})_{(2j-1) \times (2j-1)}.$$

$s_{ij} = ???$

Note that $s_{ij} = s_{ji}$.

$$\begin{aligned} s_{11} &= E([\mathbf{1}_{(X=Z=a_1)}]^2) - [E(\mathbf{1}_{(X=Z=a_1)})]^2 = f_o(a_1)S_Y(a_1-)(1 - f_o(a_1)S_Y(a_1-)) \\ &\begin{pmatrix} s_{12} & s_{13} \\ s_{22} & s_{23} \\ \cdot & s_{33} \end{pmatrix} \\ &= \begin{pmatrix} -E(\mathbf{1}_{(X=Z=a_1)})E(\mathbf{1}_{(X=Z=a_2)}) & -E(\mathbf{1}_{(X=Z=a_1)})E(\mathbf{1}_{(Z \geq a_2)}) \\ E(\mathbf{1}_{(X=Z=a_2)}) - [E(\mathbf{1}_{(X=Z=a_2)})]^2 & E(\mathbf{1}_{(X=Z=a_2)}) - E(\mathbf{1}_{(X=Z=a_2)})E(\mathbf{1}_{(Z \geq a_2)}) \\ \cdot & E(\mathbf{1}_{(Z \geq a_2)}) - [E(\mathbf{1}_{(Z \geq a_2)})]^2 \end{pmatrix} \\ &= \begin{pmatrix} -f_o(a_1)S_Y(a_1-)f_o(a_2)S_Y(a_2-) & -f_o(a_1)S_Y(a_1-)S_o(a_2-)S_Y(a_2-) \\ f_o(a_2)S_Y(a_2-)(1 - f_o(a_2)S_Y(a_2-)) & f_o(a_2)S_Y(a_2-)(1 - S_o(a_2-)S_Y(a_2-)) \\ \cdot & S_o(a_2-)S_Y(a_2-)(1 - S_o(a_2-)S_Y(a_2-)) \end{pmatrix} \end{aligned}$$

Then

$$n\sigma_{\ln \hat{S}_{pl}(t)}^2 \approx n\mathbf{p}^t \sum \mathbf{p} = (1, 1) \begin{pmatrix} p_1 & 0 & 0 \\ 0 & p_2 & p_3 \end{pmatrix} \begin{pmatrix} s_{11} & s_{12} & s_{13} \\ s_{12} & s_{22} & s_{23} \\ s_{13} & s_{23} & s_{33} \end{pmatrix} \begin{pmatrix} p_1 & 0 \\ 0 & p_2 \\ 0 & p_3 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

Verify that

$$(p_1 \ 0 \ 0) \begin{pmatrix} s_{11} & s_{12} & s_{13} \\ s_{12} & s_{22} & s_{23} \\ s_{13} & s_{23} & s_{33} \end{pmatrix} \begin{pmatrix} 0 \\ p_2 \\ p_3 \end{pmatrix} = 0. \quad (2.8)$$

By symmetry

$$\begin{aligned} n\sigma_{\ln \hat{S}_{pl}(t)}^2 &\approx (1, 1) \begin{pmatrix} p_1 & 0 & 0 \\ 0 & p_2 & p_3 \end{pmatrix} \begin{pmatrix} s_{11} & 0 & 0 \\ 0 & s_{22} & s_{23} \\ 0 & s_{32} & s_{33} \end{pmatrix} \begin{pmatrix} p_1 & 0 \\ 0 & p_2 \\ 0 & p_3 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\ &= (1, 1) \begin{pmatrix} p_1 p_1 s_{11} & 0 \\ 0 & p_2 p_2 s_{22} + p_2 p_3 s_{23} + p_3 p_2 s_{23} + p_3 p_3 s_{33} \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\ &= p_1 p_1 s_{11} + p_2 p_2 s_{22} + p_2 p_3 s_{23} + p_3 p_2 s_{23} + p_3 p_3 s_{33} \end{aligned} \quad (2.9)$$

Then it yields

$$\sigma_{\ln \hat{S}_{pl}(t)}^2 \approx \frac{1}{n} \sum_{a_k \leq t} \frac{f_o(a_k)}{S_o(a_k-)S_Y(a_k-)S_o(a_k)}. \quad (2.10)$$

§4.2.2.3. Homework:

8. Verify (2.8).
9. Using (2.9) to verify (2.10).

References.

- [*] Breslow, N.E. and Crowley, J. (1974). A large-sample study of the life table and product limit estimates under random censorship. *Ann. Statist.* 2 437-453.
- [*] Ferguson, T. S. (1996). A course in large sample theory. p. 119, Cahpman & Hall. New York.
- [*] Gill, R. (1983). Convergence of the product limit estimator on the entire half line. *Ann. Statist.* 11 49-59.
- [*] Johanson, S. (1978). The product limit estimator as maximum likelihood estimator. *Scan. J. Statist.*, 5, 195-199.

- [*] Kaplan, E. L. and Meier, P. (1958). Nonparametric estimation from incomplete observations. *J. Amer. Stat. Assoc.*, 53, 457-481.
- [*] Miller Jr., R. G. (1981). Survival analysis. *Wiley*. p. 62.
- [*] Peterson, A. V. (1977). Expressing the Kaplan-Meier estimator as a function of the empirical subsurvival functions. *JASA*. 72, 854-858.
- [*] Phadia, E.G. and Van Ryzin, J (1980). A note on the convergence rates for the product limit estimator. *Ann. Statist.* 8, 673-678.
- [*] Shorack, G. R. and Wellner, J. A. (1986). Empirical Processes with applications to statistics. *Wiley*, New York.
- [*] Stute, W. and Wang J.L. (1993). The strong law under random censorship. *Ann. Statist.* 21 1591-1607.
- [*] Wang, J. G. (1987). A note on the uniform consistency of the Kaplan-Meier estimator. *Ann. Statist.* 15, 1313-1316.
- [*] Yu, Q.Q, Ai, X.S. and Yu, K. (2012). Asymptotic Properties Of The Product-Limit-Estimator With Dependent Right Censoring *International Journal of Statistics and Management Systems* (Accepted).

§4.3. C1 model and GMLE.

Assume the case 1 interval censorship model, i.e.,
 Each patient is followed once at time $Y \sim G$;
 X and Y are independent;
 Observable random vector is

$$(L, R) = \begin{cases} (-\infty, Y) & \text{if } X \leq Y \\ (Y, +\infty) & \text{if } X > Y, \end{cases}$$

or

$$(Y, \delta) \text{ where } \delta = \mathbf{1}_{(X \leq Y)}.$$

It is more convenient to use (Y, δ) . It corresponds to an interval

$$I = \begin{cases} (-\infty, Y] & \text{if } X \leq Y \\ (Y, +\infty) & \text{if } X > Y. \end{cases}$$

Given a random sample of size n , say (Y_i, δ_i) , $i = 1, \dots, n$, let $a_1 < \dots < a_m$ be all the distinct values of Y_i s. Let

$$\begin{aligned} N_-(a_j) &= \sum_{i=1}^n \mathbf{1}_{(X_i \leq Y_i = a_j)}, \\ N_+(a_j) &= \sum_{i=1}^n \mathbf{1}_{(X_i > Y_i = a_j)}, \\ N(a_j) &= \sum_{i=1}^n \mathbf{1}_{(Y_i = a_j)}. \end{aligned}$$

The likelihood function is

$$\begin{aligned} \mathbb{L}(F) &= \prod_{i=1}^n \mu_F(I_i) \\ &= \prod_{j=1}^m [(F(a_j))^{N_-(a_j)} (S(a_j))^{N_+(a_j)}], \quad F \in \Theta, \text{ where} \end{aligned}$$

$\Theta = \{F : F \text{ is a nondecreasing function on } [-\infty, +\infty], F(-\infty) = 0 \text{ and } F(+\infty) = 1 \}$.
 That is

$$F(a_1) \leq F(a_2) \leq \dots \leq F(a_m).$$

Let $s_j = F(a_j)$, the log likelihood is

$$\mathcal{L}(F) = \sum_{j=1}^m [N_-(a_j) \log s_j + N_+(a_j) \log(1 - s_j)], \quad 0 \leq s_1 \leq \dots \leq s_m \leq 1. \quad (3.1)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial s_j} &= \frac{N_-(a_j)}{s_j} - \frac{N_+(a_j)}{1 - s_j} = 0 \\ \Rightarrow s_j &= \tilde{F}(a_j) = \frac{N_-(a_j)}{N(a_j)}. \end{aligned} \quad (3.2)$$

If $\tilde{F} \in \Theta$, or $0 \leq s_1 \leq \dots \leq s_m \leq 1$, it is the GMLE (in fact, each summand looks like the log likelihood of a binomial distribution $\text{Bin}(N(a_j), s_j)$). Otherwise, it is not and the GMLE will be on the boundary $s_i = s_j$, $i \neq j$.

Example 1. A random sample of $n = 4$. The observations (Y_i, δ_i) s are

$(1, 1), (2, 0), (2, 1), (3, 1)$.

$$\text{Then } \tilde{F}(t) = \begin{cases} 0 & \text{if } t < 1 \\ 1 & \text{if } t \in [1, 2) \\ \frac{1}{2} & \text{if } t \in [2, 3) \\ 1 & \text{if } t \geq 3. \end{cases}$$

$F \notin \Theta$. Thus \tilde{F} is not a GMLE of F . To be completed in a homework.

Theorem 1. A GMLE of F_o under C1 model is

$$\hat{F}(t) = \hat{F}(Y_{(j)}) \quad \text{if } t \in [Y_{(j)}, Y_{(j+1)}), \quad j = 0, 1, \dots, n,$$

where $Y_{(0)} = -\infty$, $Y_{(n+1)} = +\infty$, and

$$\hat{F}(Y_{(j)}) = \max_{i \leq j} \min_{k \geq j} \frac{\sum_{i \leq h \leq k} \delta_{(h)}}{k - i + 1}, \quad j = 1, \dots, n, \quad (3.3)$$

and $\delta_{(h)}$ is the δ_i that is associated with the order statistic $Y_{(h)}$, i.e. $Y_i = Y_{(h)}$, provided $\hat{F} \in \Theta$.

For proof, we refer to Ayer *et al.* (1955).

Remark.

1. The GMLE $\hat{F}(t)$ is uniquely determined at each Y_i .
2. The GMLE is not uniquely determined in the interval $(-\infty, a_1)$ and (a_i, a_{i+1}) , where $i \geq 0$, unless $\hat{F}(a_i) = \hat{F}(a_{i+1})$.
3. It can be shown that (homework) (3.3) is the same as

$$\hat{F}(Y_{(j)}) = \min_{k \geq j} \max_{i \leq j} \frac{\sum_{i \leq h \leq k} \delta_{(h)}}{k - i + 1}, \quad (3.4)$$

$$\hat{F}(Y_{(j)}) = \max_{i \leq j} \min_{k \geq i} \frac{\sum_{i \leq h \leq k} \delta_{(h)}}{k - i + 1}, \quad (3.5)$$

$$\hat{F}(Y_{(j)}) = \min_{k \geq j} \max_{i \leq k} \frac{\sum_{i \leq h \leq k} \delta_{(h)}}{k - i + 1}. \quad (3.6)$$

Example 2. For a random sample of 4 C1 observations, $(1,0)$, $(2,1)$, $(3,0)$ and $(3,1)$.

$$\hat{F}(Y_{(j)}) = \max_{i \leq j} \min_{k \geq j} \frac{\sum_{i \leq h \leq k} \delta_{(h)}}{k - i + 1}, \quad j = 1, \dots, n, \quad (3.3)$$

Let $A_{ik} = \frac{\delta_{(i)} + \dots + \delta_{(k)}}{k - i + 1}$.

$B_{ij} = \min_{k \geq j} \frac{\delta_{(i)} + \dots + \delta_{(k)}}{k - i + 1}$, $i \leq j$.

Order data: $(1,0)$, $(2,1)$, $(3,1)$, $(3,0)$.

$i \backslash k$	1	2	3	4	B_{ij}	$\hat{F}(Y_{(j)})$
1	$\frac{\delta_{(1)}}{1}$	$\frac{\delta_{(1)}+\delta_{(2)}}{2}$	$\frac{\delta_{(1)}+\delta_{(2)}+\delta_{(3)}}{3}$	$\frac{\delta_{(1)}+\delta_{(2)}+\delta_{(3)}+\delta_{(4)}}{4}$		
2		$\frac{\delta_{(2)}}{1}$	$\frac{\delta_{(2)}+\delta_{(3)}}{2}$	$\frac{\delta_{(2)}+\delta_{(3)}+\delta_{(4)}}{3}$		
3			$\frac{\delta_{(3)}}{1}$	$\frac{\delta_{(3)}+\delta_{(4)}}{2}$		
4				$\frac{\delta_{(4)}}{1}$		

$\max_{i \leq j} B_{ij}$

$i \backslash k$	1	2	3	4	B_{i1}	$\hat{F}(Y_{(1)})$	B_{i2}	$\hat{F}(Y_{(2)})$	B_{i3}	$\hat{F}(Y_{(3)})$	B_{i4}	$\hat{F}(Y_{(4)})$
1	0	1/2	2/3	2/4	0	0	1/2		1/2		1/2	
2		1	2/2	2/3			2/3	2/3	2/3		2/3	
3			1/1	1/2					1/2	2/3	1/2	
4				0							0	2/3

The GMLE is uniquely determined on the set $(-\infty, 1] \cup [2, 3]$. $\hat{F}(i) = 0, 2/3, 2/3$, respectively, and it arbitrary on $(1, 2) \cup (3, \infty)$.

In fact, it is more convenient to use formula (3.5) (or (3.6)) which results in the following matrix.

$$\hat{F}(Y_{(j)}) = \max_{i \leq j} \min_{k \geq i} \frac{\sum_{i \leq h \leq k} \delta_{(h)}}{k - i + 1}, \quad (3.5)$$

$$\hat{F}(Y_{(j)}) = \min_{k \geq j} \max_{i \leq k} \frac{\sum_{i \leq h \leq k} \delta_{(h)}}{k - i + 1}. \quad (3.6)$$

Order data: (1,0), (2,1), (3,1), (3,0).

$i \backslash k$	1	2	3	4	$\min_k A_{ik}$	$\hat{F}(Y_{(i)})$	<i>Eq. (3.5)</i>
1	$\frac{\delta_{(1)}}{1}$	$\frac{\delta_{(1)}+\delta_{(2)}}{2}$	$\frac{\delta_{(1)}+\delta_{(2)}+\delta_{(3)}}{3}$	$\frac{\delta_{(1)}+\dots+\delta_{(4)}}{4}$	m_1	m_1	
2		$\frac{\delta_{(2)}}{1}$	$\frac{\delta_{(2)}+\delta_{(3)}}{2}$	$\frac{\delta_{(2)}+\delta_{(3)}+\delta_{(4)}}{3}$	m_2	$m_1 \vee m_2$	
3			$\frac{\delta_{(3)}}{1}$	$\frac{\delta_{(3)}+\delta_{(4)}}{2}$	m_3	$\max_{i \leq 3} m_i$	
4				$\frac{\delta_{(4)}}{1}$	m_4	$\max_{i \leq 4} m_i$	
$\max_i A_{ik}$	M_1	M_2	M_3	M_4			
$\hat{F}(Y_{(k)})$	$\min_k M_k$	$\min_{k \geq 2} \{M_k\}$	$\min\{M_3, M_4\}$	M_4			

Eq. (3.6)

$\delta :$	0	1	1	0			
$i \backslash k$	1	2	3	4	$\min_k A_{ik}$	$\hat{F}(Y_{(i)})$	<i>Eq. (3.5)</i>
1	0	1/2	2/3	2/4	0	0	
2		1	2/2	2/3	2/3	2/3	
3			1/1	1/2	1/2	2/3	
4				0	0	2/3	
$\max_i A_{ik}$	0	1	1	2/3			
$\hat{F}(Y_{(k)})$	0	2/3	2/3	2/3			

Eq. (3.6)

Obtain the GMLE directly from \mathbf{L} in (3.1):

Order data: (1,0), (2,1), (3,1), (3,0).

$$\mathcal{L} = \ln(1 - s_1) + \ln s_2 + \ln s_3 + \ln(1 - s_3), \quad 0 \leq s_1 \leq s_2 \leq s_3 \leq 1$$

Boundary: $s_1 = 0$; ($\Leftarrow s_2 = 0$; $s_3 = 0$);

$$s_3 = 1; (\Leftarrow s_2 = 1; s_1 = 1);$$

$$s_1 = s_2; s_2 = s_3; (\Leftarrow (s_1 = s_3)).$$

$\mathbf{s} :$ *critical pt* (0, 1, 1/2) $s_1 = 0$ $s_3 = 1$ $s_1 = s_2$ *others*
 $\mathcal{L} :$ **violating the constraint** no need to check

$$\begin{array}{ll}
s_1 = 0: & \mathbf{s}: \text{critical point } (0, 1, 1/2) \quad s_2 = s_3 \quad s_2 = 0 \quad s_3 = 1 \\
& \mathcal{L}: \text{violating the constraint} \quad \text{no need to check} \\
s_2 = s_3: & \mathbf{s}: \text{critical point } (0, 2/3, 2/3) \quad s_1 = 0 \quad s_3 = 1 \quad s_1 = s_2 = s_3 \\
& \mathcal{L}: \quad \quad \quad \ln(\frac{2^2}{3^3}) \quad \ln(\frac{2^2}{3^3}) \quad -\infty \quad \ln\frac{1}{2^4} \\
s_1 = s_2: & \mathbf{s}: (1/2, 1/2, 1/2) \quad s_1 = 0 \quad s_3 = 1 \\
& \mathcal{L}: \quad \quad \quad \ln\frac{1}{2^4} \quad -\infty \quad -\infty
\end{array}$$

Other boundaries can be skipped.

Obtain the GMLE through $L(F)$ directly:

The distinct observations Y_i 's, 1,2,3, partition $(-\infty, \infty)$ as 4 disjoint intervals.

Let the weight assigns by an F to these intervals as p_1, \dots, p_4 . Then the likelihood function

$$\begin{aligned}
L &= (p_2 + p_3 + p_4)(p_1 + p_2)(p_4)(p_1 + p_2 + p_3) \\
&\leq (q_2 + q_4)(q_2)(q_4)(q_2) && \text{setting } p_1 + p_2 + p_3 = q_2, p_4 = q_4 \\
&= (q_2 + q_4)(q_2)^2(q_4) && q_2 + q_4 = 1 \\
&= (q_2)^2(q_4)
\end{aligned}$$

Thus $q_2 = 2/3, q_4 = 1/3$.

program to compute GMLE with C1 data

```
x=c(1,2,3,4,3,0,1,1,0,1) # (Y[1:5], delta[1:5])
```

```
x
```

```
dim(x)=c(5,2)
```

```
x
```

```

c1=function(data) {
  ord = order(data[, 1])
  data = data[ord, ]
  s = table(data[, 1], data[, 2])
  Nminus = s[, 2]
  N = s[, 1] + s[, 2]
  L = length(N)
  pt = unique(data[, 1])
  r = rep(0, L)
  mat = matrix(0, L, L)
  for(i in 1:L) {
    for(j in i:L)
      mat[i, j] = sum(Nminus[i:j])/sum(N[i:j])
  }
  r[1] = min(mat[1, ])
  r[L] = max(mat[, L])
  for(i in 2:(L - 1))
    r[i] = max(apply(mat[1:i, i:L], 1, min))
  cbind(pt, r)
}

```

```
}
```

```
c1(x)
```

§4.3.1.2. Homework.

1. Derive the GMLE of F using data from Example 1. Use two approaches: (1) use formula (3.3); (2) derive directly from the likelihood function in (3.1).
2. Assuming Y is discrete with finitely many values, show that the estimator \tilde{F} in (3.2) is a consistent estimator of $F_o(a_j)$ and find its asymptotic variance.
3. Show that the definitions (3.4), (3.5) and (3.6) of $\hat{F}(Y_{(j)})$ are equivalent to (3.3). **Hint:**

Inspect the matrix $(A_{ik})_{j \times (n-j+1)}$, where $A_{ik} = \frac{\sum_{i \leq h \leq k} \delta_{(h)}}{k-i+1}$. Try first $n = 3$ or 4 and use induction argument on j .

In many practical situations, consistency of the GMLE is only relevant for $t \leq \tau$, where

$$\tau = \sup\{t : G(t) < 1\}.$$

Consistency of the GMLE has been investigated by Ayer *et al.* (1955), Groeneboom and Wellner (1992), Yu *et al.* (1998), and Schick and Yu (2000).

Theorem 2. Under the C1 model, the GMLE \hat{F} satisfies

- (1) $\lim_{n \rightarrow \infty} \int |\hat{F}(t) - F_o(t)| dG(t) = 0$ a.s. (Schick and Yu (2000));
- (2) $\hat{F}(a) \rightarrow F_o(a)$ a.s. for each $a \in \mathcal{A}$, where $\mathcal{A} = \{a : P(Y = a) > 0\}$ (Yu *et al.* (1998));
- (3) If F_o is continuous in $(0, \tau]$, $P\{Y = \tau\} > 0$ or $F_o(\tau) = 1$, and the range of Y is dense in $[0, \tau]$, then \hat{F} is uniformly strongly consistent on $(-\infty, \tau]$, i.e., $\sup_{x \leq \tau} |\hat{F}(x) - F_o(x)| \rightarrow 0$ a.s. (Schick and Yu (2000)).

In clinical follow-ups, the studies typically last for a certain period of time, say $[0, \tau]$. It is often that $F_o(\tau) < 1$.

Gentleman and Geyer (1994) claimed a vague convergence result in their Theorem 2 and Huang (1996) claimed a uniform strong consistency result in his Theorem 3.1.

Both of their results as stated imply

- the uniform strong consistency of the GMLE on $[0, \tau]$ in the case 1 model,
- if F_o is continuous and $G'(x) > 0$ on $[0, \tau]$.

The following example shows that this is not true.

Example 3. Consider C1 data $(Y_1, \mathbf{1}_{(X_1 \leq Y_1)}), \dots, (Y_n, \mathbf{1}_{(X_n \leq Y_n)})$, where the survival times X_1, \dots, X_n are i.i.d. $\sim U(0, 3)$ and the inspection times Y_1, \dots, Y_n are i.i.d. $\sim U(0, 2)$. $\tau = ?$ Note that

on the event $\cup_{j=1}^n B_j$, where $B_j = \{X_j \leq 1 \leq Y_j, Y_j > Y_i, i = 1, \dots, n, i \neq j\}$

we have $\hat{F}_n(2) = \hat{F}_n(2-) = \hat{F}(Y_{(n)}) = 1$ (as $Y_{(n)} < 2$ and

$$\hat{F}(Y_{(n)}) = \max_{i \leq n} \min_{k \geq n} \frac{\sum_{i \leq h \leq k} \delta^{(h)}}{k-i+1} = \max_{i \leq n} \frac{\sum_{i \leq h \leq n} \delta^{(h)}}{n-i+1} = 1 \text{ (see (3.3)).}$$

The event has probability $1/3 - \frac{1}{3 \cdot 2^n} \geq \frac{1}{6}$, as it equals

$$\begin{aligned} & P\{\cup_{j=1}^n B_j\} \\ &= nP\{X_1 \leq 1 \leq Y_1, Y_1 > Y_j, j = 2, \dots, n\} \\ & \quad (B_j \text{s are disjoint events with the same probability}) \\ &= nP\{X_1 \leq 1\}P\{1 \leq Y_1, Y_1 > Y_j, j = 2, \dots, n\} \\ &= \frac{n}{3} \int_1^2 \frac{1}{2} P\{Y_j < y_1, j = 2, \dots, n\} dy_1 \\ &= \frac{n}{3} \frac{1}{n} \left(\frac{y}{2}\right)^n \Big|_1^2 \\ &= \frac{1}{3} (1 - 2^{-n}). \end{aligned}$$

That is $P(\hat{F}(2-) = 1) \geq P(\cup_j B_j) = \frac{1}{3} - \frac{1}{32^n}$.

Since $F_o(2) = F_o(2-) = 2/3$, we see that the following two statements are false:

$\hat{F}(2-)$ converges to $F_o(2-)$ a.s. ($P(\hat{F}(2-) \rightarrow F_o(2-)) = 1$),

and $\hat{F}(x)$ converges to $F_o(x)$ a.s. for $x = 2$ ($P(\hat{F}(2) \rightarrow F_o(2)) = 1$).

This shows that point-wise convergence on the closed interval $[0, 2]$ to a continuous F_o is not implied by the condition: $\frac{dG}{dx} > 0$ for all $x \in [0, 2]$.

Remark. The GMLE under the assumptions in Example 2 is consistent for all $t < \tau$. It follows from Proposition 3.2 in Schick and Yu (2000).

Question: Is $\hat{F}(t)$ consistent for all $t < \tau$?

Example 4. Suppose that $X \sim Exp(1)$ and Y has a Poisson distribution with mean 1. Then $\tau = +\infty$.

$$\lim_{n \rightarrow \infty} \hat{S}(1) = S_o(1) \text{ a.s. ? } (S_o(1) = e^{-1})$$

$$\lim_{n \rightarrow \infty} \hat{S}(1.5) = S_o(1.5) \text{ a.s. ? } (S_o(1.5) = e^{-1.5})$$

Answer:

$$\lim_{n \rightarrow \infty} \hat{S}(1) = S_o(1) \text{ a.s. !}$$

$$\lim_{n \rightarrow \infty} \hat{S}(1.5) \neq S_o(1.5) \text{ a.s..}$$

$\lim_{n \rightarrow \infty} \hat{S}(1.5) = S_o(1)$ a.s..

In fact, the GMLE is not consistent at $t \in (i, i + 1)$ for all integers i .

Note that under the RC model, the PLE $\hat{S}_{pl}(t)$ is always consistent for $t < \tau$!

Peto (1973) and Turnbull (1976) conjectured that

for arbitrary F_o and G , the GMLE is asymptotically normal at the usual $n^{1/2}$ rate.

It was, however, shown by Groeneboom and Wellner (1992) that

this conjecture is false if F_o and G satisfy certain smoothness assumptions.

Indeed, their Theorem 3 below establishes that under differentiability assumptions on F_o and G the convergence is at the slower $n^{1/3}$ rate and the limiting distribution is not normal.

Theorem 3. *Let t_o be such that $0 < F_o(t_o), G(t_o) < 1$, and let F_o and G be differentiable at t_o , with strictly positive derivatives $f_o(t_o)$ and $g(t_o)$, respectively. Then*

$$n^{1/3} \frac{\hat{F}(t_o) - F_o(t_o)}{\{\frac{1}{2}F_o(t_o)S_o(t_o)f_o(t_o)/g(t_o)\}^{1/3}} \xrightarrow{D} 2Z, \text{ as } n \rightarrow \infty,$$

where $Z \equiv \operatorname{argmin}(W(t) + t^2)$ and W is the two-sided Brownian motion starting from 0. i.e., $\forall \omega$ in the sample space, $Z(\omega) = t_o$, where $W(t_o)(\omega) + t_o^2 \leq W(t)(\omega) + t^2$ for $t \geq 0$.

Definition. A real-valued continuous-time process $W(t)$ is called a **Gaussian process** if each finite-dimensional vector $(W(t_1), \dots, W(t_m)) \sim N(\bar{\mu}(\mathbf{t}), \Sigma(\mathbf{t}))$, where $\Sigma(\mathbf{t})$ can be singular and $\mathbf{t} = (t_1, \dots, t_m)$, $m \geq 1$. If W has independent increments and $W(s+t) - W(s) \sim N(0, \sigma^2 t)$, where $t > 0$, it is called a **Wiener process**, or a **Brownian Motion**.

Definition. A 3-dimensional stochastic process $\mathbf{X} = (X_1, X_2, X_3)$ is called a **Brownian Motion** if it satisfies

1. X_i 's are i.i.d. Wiener processes,
2. $\mathbf{X}(0) = (0, 0, 0)$,

Remark. The main assumption in Theorem 3 is that G is differentiable, i.e. Y is continuous.

There are many practical situations in which Y is discrete. In medical research, for example, the data are often recorded as integers (to represent number of days, weeks etc). Then the conclusion is different.

Let $\mathcal{A}_* = \mathcal{A} \cup \{-\infty, \infty\}$. For $x \in (-\infty, +\infty)$, set

$$x_- := \sup\{a \in \mathcal{A}_* : a < x\} \text{ and } x_+ := \inf\{a \in \mathcal{A}_* : a > x\}.$$

We say x is a *regular point*, if x belongs to \mathcal{A} , x_- and x_+ belong to \mathcal{A}_* , $x_- < x < x_+$ and $F_0(x_-) < F_0(x) < F_0(x_+)$. It is worth mentioning that there may be infinitely many regular points. For example, if F_0 is strictly increasing and \mathcal{A} is the set of all positive integers, then every positive integer is a regular point.

Theorem 4. (Yu *et al.* (1998)). *Let x be a regular point. Then $n^{1/2}(\hat{F}(x) - F_o(x))$ is asymptotically normal with mean 0 and variance $F_o(x)(1 - F_o(x))/g(x)$. This asymptotic variance can be consistently estimated by $\hat{F}(x)(1 - \hat{F}(x))/N(x)$. Also, if $x_1 < \dots < x_m$ are regular points, then $n^{1/2}(\hat{F}(x_1) - F_o(x_1), \dots, \hat{F}(x_m) - F_o(x_m))$ is asymptotically normal with mean vector 0 and diagonal covariance matrix.*

Remark. Theorems 3 and 4 both allow X takes on infinitely many values.

Remark. Suppose that F_o is strictly monotone and Y takes on finitely many values, say at $a_1 < \dots < a_m$, with $F(a_1) > 0$ and $F(a_m) < 1$. For n large enough,

$$\hat{F}(a_i) = \tilde{F}(a_i), \quad i = 1, \dots, m - 1.$$

Reason: Let $\epsilon = \min\{F_o(a_1), F_o(a_i) - F_o(a_{i-1}), i = 2, \dots, m, S_o(a_m)\}$. Then $\epsilon > 0$. Let Ω be the event that $N_-(a_j)/n$ and $N(a_j)/n$ converge for $j = 1, \dots, m$. Then for each $\omega \in \Omega$, for n large enough, $|\tilde{F}(a_j) - F_o(a_j)| < \epsilon/3$ for all j . It follows that

$$0 \leq \tilde{F}(a_1) < F_o(a_1) + \epsilon/3 < F_o(a_2) - \epsilon/3 < \tilde{F}(a_2) < \dots < \tilde{F}(a_m) \leq 1. \quad (3.7)$$

As a consequence, \tilde{F} is a GMLE.

In fact, the asymptotic covariance matrix of $(\hat{F}(a_1), \dots, \hat{F}(a_m))$ can be estimated by the expression

$$\hat{\Sigma} = \left(-\frac{\partial^2 \mathcal{L}}{\partial \mathbf{s} \partial \mathbf{s}^t} \Big|_{\mathbf{s}^t = (\hat{F}(a_1), \dots, \hat{F}(a_m))} \right)^{-1}.$$

Now $\mathcal{L} = \ln \prod_{i=1}^m s_i^{N_-(a_i)} (1 - s_i)^{N_+(a_i)}$.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{s}^t} = \left(\frac{N_-(a_1)}{s_1} - \frac{N_+(a_1)}{1 - s_1}, \dots, \frac{N_-(a_m)}{s_m} - \frac{N_+(a_m)}{1 - s_m} \right).$$

$$\frac{\partial^2 \mathcal{L}}{\partial \mathbf{s} \partial \mathbf{s}^t} = - \begin{pmatrix} \frac{N_-(a_1)}{s_1^2} + \frac{N_+(a_1)}{(1-s_1)^2} & 0 & \cdot & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \frac{N_-(a_m)}{s_m^2} + \frac{N_+(a_m)}{(1-s_m)^2} \end{pmatrix}.$$

$$\begin{aligned} & \frac{1}{n} \left(-\frac{\partial^2 \mathcal{L}}{\partial \mathbf{s} \partial \mathbf{s}^t} \Big|_{\mathbf{s}^t = (\hat{F}(a_1), \dots, \hat{F}(a_m))} \right) \\ \rightarrow & \begin{pmatrix} \frac{F_o(a_1)g(a_1)}{F_o^2(a_1)} + \frac{S_o(a_1)g(a_1)}{(S_o(a_1))^2} & 0 & \cdot & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \frac{F_o(a_m)g(a_m)}{F_o^2(a_m)} + \frac{S_o(a_m)g(a_m)}{(S_o(a_m))^2} \end{pmatrix} \text{ a.s.} \\ = & \begin{pmatrix} \frac{g(a_1)}{F_o(a_1)} + \frac{g(a_1)}{S_o(a_1)} & 0 & \cdot & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \frac{g(a_m)}{F_o(a_m)} + \frac{g(a_m)}{S_o(a_m)} \end{pmatrix} \\ = & \begin{pmatrix} \frac{g(a_1)}{F_o(a_1)S_o(a_1)} & 0 & \cdot & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \frac{g(a_m)}{F_o(a_m)S_o(a_m)} \end{pmatrix} \\ n\hat{\Sigma} \rightarrow & \begin{pmatrix} \frac{F_o(a_1)S_o(a_1)}{g(a_1)} & 0 & \cdot & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \frac{F_o(a_m)S_o(a_m)}{g(a_m)} \end{pmatrix} \text{ a.s.} \end{aligned}$$

Theorem 4 is an extension of the above observation.

However, if F_o is not strictly increasing on $\{a_1, \dots, a_m\}$, it is not true that

$$\tilde{F}(a_i) = \hat{F}(a_i) \quad i = 2, \dots, m - 1, \text{ if } n \text{ is large enough,}$$

as (3.7) does not hold. In such a case, Theorem 4 is false.

Question. In application, how to determine when the censoring distribution G is continuous or discrete ?

The important feature of a discrete Y is that there are ties among Y_i 's.

If there are few ties, then one should consider Y is continuous.

Otherwise, discrete.

Question. If Y is continuous, how to construct a confidence interval (CI) for $F(t_o)$?

Note that by Theorem 3 and Theorem 2.4 of Banerjee and Wellner (2001),

$$n^{1/3}(\hat{F}(t_o) - F_o(t_o)) \xrightarrow{D} hZ,$$

where

$$h = \{4F_o(t_o)S_o(t_o)f_o(t_o)/g(t_o)\}^{1/3}.$$

Thus, a 95% CI for $F_o(t_o)$ is

$$\hat{F}(t_o) \pm n^{-1/3} \hat{h} Q_{0.025} = \hat{F}(t_o) \pm n^{-1/3} \hat{h} 0.99818$$

where Q_α is the $100(1 - \alpha)$ quantile of the distribution of Z ($P\{Z > Q_\alpha\} = \alpha$), which is provided in Groeneboom and Wellner (2001), and \hat{h} is an estimate of h , e.g.,

$$\hat{h} = \{4\hat{F}(t_o)\hat{S}(t_o)\hat{f}(t_o)/\hat{g}(t_o)\}^{-1/3}$$

$$\hat{f}(t) = \int \frac{1}{w_n} K\left(\frac{x-t}{w_n}\right) d\hat{F}(x),$$

w_n is the window width, $w_n \rightarrow 0$, $K(\cdot)$ is a kernel (e.g., $K(x) = \frac{1}{2}\mathbf{1}_{(-1 < x \leq 1)}$) and thus

$$\hat{f}(t) = \frac{\hat{F}(t + w_n) - \hat{F}(t - w_n)}{2w_n},$$

$$\hat{g}(t) = \int \frac{1}{w_n} K\left(\frac{x-t}{w_n}\right) d\hat{G}(x),$$

$$\hat{G}(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(Y_i \leq x)}.$$

§4.3.1.3. Homework.

4. Let F_o be $Exp(\rho)$ (with mean $1/\rho$) and $Y \sim U(0, 3/\rho)$.
 - 4.a. Generate $n = 400$ case 1 data and compute $\hat{F}(1/\rho)$.
 - 4.b. Repeat 4.a 50 times and compute the sample mean $\overline{\hat{F}}(1/\rho)$ and sample variance S_n^2 .
 - 4.c. Discuss the region on which the GMLE is consistent.
 - 4.d. Repeat 4.a and 4.b, with $n = 100$ and compare $\frac{S_{100}}{S_{400}}$ to $\sqrt{400/100}$ and to $(\frac{400}{100})^{1/3}$, which is closer ?
5. Let Y be a discrete random variable taking values 1, 2, 3 and 4; and let $X \sim U(0, 5)$.
 - 5.a. Generate 400 case 1 data and compute $\hat{F}(3)$.
 - 5.b. Repeat 4.a 50 times and compute the sample mean $\overline{\hat{F}}(3)$ and sample variance S_n^2 .
 - 5.c. Discuss the region on which the GMLE is consistent.
 - 5.d. Repeat 5.a and 5.b, with $n = 100$ and compare $\frac{S_{100}}{S_{400}}$ to $\sqrt{400/100}$ and to $(\frac{400}{100})^{1/3}$, which is closer ?
6. Under the assumption in Problem # 4, generate a random sample of $n = 400$ C1 data from $Exp(\rho)$. Plot the survival curves of $S(t; \rho_o)$, the MLE and the GMLE on the same graph. Now pretend the data is RC data, plot the survival curves of $S(\cdot; \rho_o)$, the PLE and MLE curves. Make comments on the plots.

References.

- * Ayer, M., Brunk, H.D., Ewing, G.M., Reid, W.T. and Silverman, E. (1955). An empirical distribution function for sampling incomplete information. *Ann. Math. Statist.* 26, 641-647.
- * Bannerjee, M and Wellner, J.A. (2001). Likelihood ratio tests for monotone functions. *Ann. Statist.* 29 1699-1731.
- * Gentleman, R. and Geyer, C. J. (1994). Maximum likelihood for interval censored data: Consistency and computation. *Biometrika*, 81, 618-623.
- * Groeneboom, P. and Wellner, J.A. (2001). Computing Chernoff's distribution. *J. Comp. and Graphical Statist.* 10, 388-400.
- * Huang, J. (1996). Efficient estimation for proportional hazards models with interval censoring. *Ann. Statist.*, 24, 540-568.
- * Schick, A. and Yu, Q. Q. (2000). Consistency of the GMLE with mixed case interval-censored data. *Scan. J. of Statist.* Vol. 27 45-55.
- * Yu, Q. Q., Schick, A., Li, L. X. and Wong, G. Y. C. (1998). Asymptotic properties of the GMLE in the case 1 interval-censorship model with discrete inspection times. *Canadian Journal of Statistics*, 26, No. 4, 619-627.

§4.4. Self-consistent (SC) algorithm.

Under the RC model, C1 model and LC model, the GMLE of F_o has a closed form solution. However, this is not so in general.

There are several numerical methods to compute the GMLE:

1. The Newton-Raphson algorithm (Peto (1973));
2. The self-consistent algorithm (Turnbull (1976));
3. The convex minorant algorithm (Groeneboom and Wellner (1992)).

The self-consistent algorithm is easy to implemented, so we introduce it in this section.

Let I_1, \dots, I_n denote the observed intervals.

Define *innermost intervals* $A_j, j = 1, 2, \dots, m$, induced by I_1, \dots, I_n to be all the disjoint intervals which are non-empty intersections of these I_i 's

(e.g. $I_k = I_k \cap I_k$ is an intersection of I_i 's) such that

$$A_j \cap I_i = \emptyset \text{ or } A_j \quad \forall i \text{ and } j.$$

Let the endpoints of the innermost intervals be a_j and $b_j, j = 1, \dots, m$, where

$$a_1 \leq b_1 \leq a_2 \leq b_2 \leq \dots \leq a_m \leq b_m.$$

Let $\delta_{ij} = \mathbf{1}_{(A_j \subset I_i)}$.

The following example illustrates the procedure for finding innermost intervals.

Example 2.1. Suppose that there are five observed intervals

$(1, 4], [2, 2], (2, 6], [5, 5]$, and $(1, 6]$.

Then there are two exact observations,

$I_1 = [2, 2]$ and $I_2 = [5, 5]$, and

three censored intervals,

$I_3 = (1, 4], I_4 = (2, 6]$ and $I_5 = (1, 6]$.

Furthermore, there are three innermost intervals,

$A_1 = [2, 2], A_2 = [5, 5]$, and $A_3 = (2, 4]$.

Peto (1973) shows that the GMLE of F_o only

assigns weights, say s_1, \dots, s_m , to the corresponding innermost intervals A_1, \dots, A_m .

The generalized likelihood function \mathbf{L} can be simplified as

$$\mathbf{L} = \mathbf{L}(s_1, \dots, s_m) = \prod_{i=1}^n \left[\sum_j \delta_{ij} s_j \right], \quad (4.1)$$

where $\mathbf{s}^t = (s_1, \dots, s_{m-1}) \in D_s$ is the transpose of \mathbf{s} ,

$D_s = \{ \mathbf{s} : s_i \geq 0, s_1 + \dots + s_{m-1} \leq 1 \}$ and

$s_m = 1 - s_1 - \dots - s_{m-1}$.

Turnbull (1976) proposes a self-consistent algorithm for obtaining the GMLE via an iterative procedure as follows.

At step 1, let $s_j^{(1)} = 1/m$ for $j = 1, \dots, m$.

At step h , $s_j^{(h)} = \sum_{i=1}^n \frac{1}{n} \frac{\delta_{ij} s_j^{(h-1)}}{\sum_{k=1}^m \delta_{ik} s_k^{(h-1)}}, \quad j = 1, \dots, m, \quad h \geq 2.$

Stop when $\mathbf{s}^{(h)}$ converges, i.e., $\|\mathbf{s}^{(h)} - \mathbf{s}^{(h-1)}\|$ is small enough.

He shows that, as $h \rightarrow \infty$, $s_j^{(h)}$ converges to the GMLE, \hat{s}_j , which maximizes \mathbf{L} and satisfies the system of self-consistent equations

$$s_j = \sum_{i=1}^n \frac{1}{n} \frac{\delta_{ij} s_j}{\sum_{k=1}^m \delta_{ik} s_k}, \quad j = 1, \dots, m, \quad \mathbf{s} \in D_s. \quad (4.2)$$

A solution $\mathbf{s} = \hat{\mathbf{s}}$ to (4.2) is called a self-consistent estimator (SCE) of \mathbf{s} .

An estimate $\hat{F}(t)$ of $F(t)$ can be uniquely defined for $t \in [b_i, a_{i+1}]$ by

$$\hat{F}(b_i) = \hat{F}(a_{i+1}-) = \hat{s}_1 + \cdots + \hat{s}_i,$$

but is not uniquely defined for t being in a non-singleton innermost interval (Peto, 1973; Turnbull, 1976).

To avoid this ambiguity we define

$$\hat{F}(t) = \sum_{j: A_j \subset (-\infty, t]} \hat{s}_j = \sum_{b_j \leq t} \hat{f}(b_j). \quad (4.3)$$

where $\hat{f}(b_j) = \hat{s}_j$ and $(a, b]$ is an empty set if $a = b$.

Under the RC model, the definition of \hat{F} given by (4.3) reduces to the PLE.

Under the C1 model, it reduces to the max-min solution.

Remark. An SCE of \mathbf{s} can be viewed as

a critical point of $L(\mathbf{s})$, subject to the constraint on \mathbf{s} .

The reason is as follows.

Given $j \in \{1, \dots, m\}$, let $\mathbf{s}(\epsilon) = (s_1(\epsilon), \dots, s_{m-1}(\epsilon))$ be defined by

$$s_k(\epsilon) = \begin{cases} \frac{s_k}{1+\epsilon} & \text{if } k \neq j \\ \frac{s_j + \epsilon}{1+\epsilon} & \text{if } k = j, \end{cases} = \begin{cases} \frac{s_k}{1+\epsilon} & \text{if } k \neq j \\ \frac{s_j - 1}{1+\epsilon} + 1 & \text{if } k = j, \end{cases}$$

Write $\Lambda_j(\epsilon) = \ln L(\mathbf{s}(\epsilon)) = \sum_{i=1}^n \ln \sum_{k=1}^m \delta_{ik} s_k(\epsilon)$, then verify that

$$\left. \frac{\partial \Lambda_j(\epsilon)}{\partial \epsilon} \right|_{\epsilon=0} = -n + \sum_{i=1}^n \frac{\delta_{ij}}{\sum_{k=1}^m \delta_{ik} s_k}.$$

If \mathbf{s} is the GMLE, then the value $\epsilon = 0$ maximizes $\Lambda_j(\epsilon)$ for each j subject to $\epsilon \geq 0$.

In other words, we have

$$\left. \frac{\partial \Lambda_j(\epsilon)}{\partial \epsilon} \right|_{\epsilon=0} = -n \left(1 - \frac{1}{n} \sum_{i=1}^n \frac{\delta_{ij}}{\sum_{k=1}^m \delta_{ik} s_k} \right) \leq 0$$

with equality unless $s_j = 0$. (**compare to (4.2)**).

Consequently, (4.2) holds. That is, the GMLE satisfies (4.2).

Of course the fact that (4.2) holds for \mathbf{s} does not implies that \mathbf{s} is the GMLE, but it implies that \mathbf{s} is a critical point of the likelihood.

§4.4.2. Homework.

1. Suppose our data consist of 5 intervals: $(-\infty, 5]$, $[2, 2]$, $(4, 7]$, $(3, 8]$, $(9, +\infty)$. Find all the innermost intervals. Compute the GMLEs of \mathbf{s} and F by two methods: directly from differentiation and by the SC algorithm. Verify that it is a solution of the self-consistent equation (4.2).

Eq. (4.2) is in the form of density function as $\hat{s}_i = \hat{f}_X(b_i)$. Its cumulative form is

$$\begin{aligned} H(x) &= \sum_{i=1}^n \frac{1}{n} \frac{\mu_H(I_i \cap (-\infty, x])}{\mu_H(I_i)} & (s_j &= \frac{1}{n} \sum_{i=1}^n \frac{\delta_{ij} s_j}{\sum_{k=1}^m \delta_{ik} s_k}) \\ &= \sum_{i=1}^n \frac{1}{n} \mathbf{1}_{(l_i \leq x < r_i)} \frac{H(x) - H(l_i)}{H(r_i) - H(l_i)} + \sum_{i=1}^n \frac{\mathbf{1}_{(r_i \leq x, r_i \in \mathcal{S}_H)}}{n}, \quad H \in \Theta, \end{aligned}$$

which is equivalent to

$$H(x) = \int_{l \leq x < r} \frac{H(x) - H(l)}{H(r) - H(l)} d\hat{Q}(l, r) + \int_{r \leq x} \mathbf{1}_{(r \in \mathcal{S}_H)} d\hat{Q}(l, r), \quad H \in \Theta, \quad (4.4)$$

where $Q(l, r) = P(L \leq l, R \leq r)$,

$$\hat{Q}(l, r) = \sum_{i=1}^n \frac{\mathbf{1}_{(l_i \leq l, r_i \leq r)}}{n} \quad (\text{the edf of } Q),$$

\mathcal{S}_H is the support set of H , namely, $x \in \mathcal{S}_H$ if $|H(x + \epsilon) - H(x - \epsilon)| > 0 \forall \epsilon > 0$.

Definition A solution H (or \mathbf{s}) to (4.4) (or (4.2)) is called a self-consistent estimator (SCE) of F_o (or \mathbf{s}).

Remark GMLE \Rightarrow SCE, but SCE $\not\Rightarrow$ GMLE.

Example 1. Six case 2 IC data: two $(1,5]$ and $(3,7]$, one $(-\infty, 3]$ and $(5, +\infty)$.

3 IIs: $(1, 3]$, $(3,5]$ and $(5,7]$.

Solving equation (4.2) yields two solutions for (s_1, s_2, s_3) :

$(1/3, 1/3, 1/3)$ and $(1/2, 0, 1/2)$.

These two SCEs of \mathbf{s} yield two SCEs of F_o :

$$\hat{F}(t) = \begin{cases} 0 & \text{if } t < 3 \\ 1/3 & \text{if } 3 \leq t < 5 \\ 2/3 & \text{if } 5 \leq t < 7 \\ 1 & \text{if } t \geq 7; \end{cases} \quad \text{and} \quad \tilde{F}(t) = \begin{cases} 0 & \text{if } t < 3 \\ 1/2 & \text{if } 3 \leq t < 7 \\ 1 & \text{if } t \geq 7. \end{cases}$$

In this example, without using any algorithm, there are **two ways to find the GMLE**:

- (1) derive directly from $L(F)$, or
- (2) find the SCE that has the largest $L(F)$. Verify that

$$L(\hat{F}) > L(\tilde{F}). \quad (4.5)$$

Thus \hat{F} is a GMLE but \tilde{F} is not a GMLE.

\hat{F} and \tilde{F} both satisfy the summation equation in (4.4).

Thus \hat{F} and \tilde{F}_2 are SCEs.

Remark Gu and Zhang (1993) defined an SCE H of F_o to be a solution to the equation in (4.4) without the restriction $H \in \Theta$. They thought that the solution of (4.4) will belong to Θ , as their proofs need the restriction. The following example show that their thought is wrong and there is an $H \notin \Theta$ but the H is a solution to (4.4).

Example 2 Six DC data.

$$N_1 = \sum_{i=1}^n \mathbf{1}_{(L_i=R_i=2)} = 1, \quad N_2 = \sum_{i=1}^n \mathbf{1}_{(L_i=R_i=5)} = 3,$$

$$N_3 = \sum_{i=1}^n \mathbf{1}_{((L_i, R_i)=(-\infty, 4])} = 1 \text{ and } N_4 = \sum_{i=1}^n \mathbf{1}_{((L_i, R_i)=(2, +\infty))} = 1.$$

A solution H to the SC equation in (4.4) gives mass $(1/2, -1/4, 3/4)$ to points 2, 4 and 5.

Since H has jumps only at 2, 4 and 5, it suffices to verify (4.4) at $-\infty, 2, 4$ and 5.

Verify as follows.

$$H(-\infty) = 0, \quad \text{RHD of (4.4)} = \frac{N_1}{n}0 + \frac{N_2}{n}0 + \frac{N_3}{n}0 + \frac{N_4}{n}0 = 0 \text{ trivially.}$$

$$H(2) = 1/2, \quad \text{RHD of (4.4)} = \frac{1}{6} + \frac{3}{6} \cdot 0 + \frac{1}{6} \frac{1/2 - 0}{1/4 - 0} = 1/2.$$

$$H(4) = 1/4, \quad \text{RHD of (4.4)} = \frac{1}{6} + \frac{3}{6} \cdot 0 + \frac{1}{6} \cdot 1 + \frac{1}{6} \frac{1/4 - 1/2}{1 - 1/2} = 1/4.$$

$$H(5) = 1, \quad \text{RHD of (4.4)} = \frac{1}{6} + \frac{3}{6} + \frac{1}{6} + \frac{1}{6} \frac{1 - 1/2}{1 - 1/2} = 1.$$

Remark. The above examples illustrate that

1. A solution to the SC equation is a critical point under the constraint that $\sum_{i=1}^m s_i = 1$.
2. Moreover, the solution to the SC algorithm with initial point $s_i \geq 0$ and $\sum_{i=1}^m s_i = 1$ is a critical point under the constraint that $s_i \geq 0$ and $\sum_{i=1}^m s_i = 1$.
3. Finally, the solution to the SC algorithm with initial point $s_i = 1/m$ is the GMLE.

§4.4.3. Homework.

2. Solve the two SCEs directly from Eq. (4.2) and (4.3) using the data in Example 1. Derive the GMLE in two ways.
3. Define a second GMLE of F_o in Example 1. Where are the GMLE uniquely defined in the example ?
4. Derive an SCE in Example 2 using SC algorithm (using some program).

§4.5. SCE under DC model

Assume the DC model. That is,

1. X and the censoring vector (Z, Y) are independent;

2. $Z < Y$ w.p.1.;
3. Observe $(L, R) = (-\infty, Z)\mathbf{1}_{(X \leq Z)} + (X, X)\mathbf{1}_{(Z < X \leq Y)} + (Y, +\infty)\mathbf{1}_{(X > Y)}$.

Define $\tau_l = \sup\{x : \max(S_X(x), S_Z(x)) = 1\}$ and $\tau_r = \inf\{x : \min(S_X(x), S_Y(x)) = 0\}$.

It is obvious that

if $F_o(\tau_l) > 0$ (or $F_o(\tau_r) < 1$), $F_o(x)$ is not identifiable for $x < \tau_l$ (or $x > \tau_r$).

Thus in such cases, we only consider the estimation of $F_o(x)$ for $x \in [\tau_l, \tau_r]$.

Denote $P_c(x) = P\{X \text{ is not censored} | X = x\}$ and

$K(x) = P\{Z < x \leq Y\}$.

$K(x) = P_c(x)$ if $P_c(x)$ exists.

Turnbull (1974), Chang and Yang (1987), Chang (1990), Gu and Zhang (1993), and Yu and Li (2001) established consistency and asymptotic normality of the SCE with DC data under various sets of regularity conditions such as follows.

- (AS1) The probability $P\{X \in (\tau_l, \tau_r] \text{ and } K(X) = 0\} = 0$.
- (AS2) $P\{L = \tau_r\} > 0$ if $F_o(\tau_r-) < 1$ and $P\{R = \tau_l\} > 0$ if $F_o(\tau_l) > 0$.
- (AS3) $\int_{\tau_l < u < \tau_r} \frac{-d[S_Y(u) + S_Z(u)]}{S_Y(u) - S_Z(u)} < \infty$;
- (AS4) $K(x) = P_c(x) > 0$ for all $x \in (\tau_l, \tau_r)$;
- (AS5) Z and Y take on finitely many values, say $b_1 < \dots < b_N$, and $0 < F_o(b_1) < \dots < F_o(b_N) < 1$.
- (AS6) There are at most m IIs for each sample size n and $\mu_{F_o}(A_j) > 0$ for each II A_j

Theorem 1 (Yu and Li (2001)) *Suppose that (AS1) and (AS2) hold, the SCE \hat{F} satisfies*

$$\lim_{n \rightarrow \infty} \sup_{x \in [\tau_l, \tau_r]} |\hat{F}(x) - F_o(x)| = 0 \text{ almost surely,}$$

$$\lim_{n \rightarrow \infty} \sup_x |\hat{F}(x) - F_o(x)| = 0 \text{ almost surely if } F_o(\tau_l) = 0 \text{ and } F_o(\tau_r) = 1.$$

Theorem 2 (Yu and Li (2001)) *Suppose that AS1 holds. Moreover either AS5 or AS6 holds; or AS2, AS3 and AS4 hold, $\tau_l < \tau_r$. Then the GMLE \hat{F} satisfies that $\sqrt{n}(\hat{F} - F_o)$ converges in distribution to a Gaussian process $Z(x)$ on $[\tau_l, \tau_r]$. Furthermore, the GMLE is asymptotically efficient.*

Theorem 3 (Turnbull (1974)) *Suppose that the assumptions in Theorem 2 hold. Let $b_1 < \dots < b_{k+1}$ be the distinct right endpoints of all the innermost intervals induced by the observed intervals. Then the covariance matrix of the SCE $(\hat{F}(b_1), \dots, \hat{F}(b_k))$ can be estimated by $(J(\hat{F}))^{-1}$, where $J(F) = -\left(\frac{\partial^2 \mathcal{L}}{\partial F(b_i) \partial F(b_j)}\right)_{k \times k}$ and $\mathcal{L} = \ln L$. Furthermore, J^{-1} is of the*

$$\text{form } J^{-1} = \begin{pmatrix} c_1 & d_1 & 0 & 0 & \dots & 0 & 0 \\ d_1 & c_2 & d_2 & 0 & \dots & 0 & 0 \\ 0 & d_2 & c_3 & d_3 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & d_{k-1} & c_k \end{pmatrix}.$$

Remark Under AS6 and under all the IC models discussed so far, the GMLE of the weights on the II's are consistent and let $b_1 < \dots < b_{k+1}$ be the right endpoints of the all possible II's, the $(\hat{F}(b_1), \dots, \hat{F}(b_k))$ are asymptotically normal and their asymptotic covariance matrix is the inverse of the Fisher information matrix $(E(J(F_o)))^{-1}$ and can be estimated by the inverse of the empirical Fisher information matrix $(J(\hat{F}))^{-1}$.

We shall refer the proofs of the theorems to the literature. By means of the following example, we explain the assumptions AS1-AS6.

Example 1 Let (Z, Y) takes values $(0.5, 2)$, $(2, 4)$ and $(4, 8)$, with equal probabilities g_1 , g_2 and g_3 , respectively, where $g_1 + g_2 + g_3 = 1$, and let $F_o(x) = p_1 \mathbf{1}(x \geq 1) + p_2 \mathbf{1}(x \geq 5)$,

where $p_1 = p_2$. Derive all possible SCEs of F_o based on Eq. (4.2) in §4.4. Derive the GMLE of F_o . Are they consistent and asymptotically normal ?

Solution We first solve for SCE of \mathbf{s} . Example 2 in §4.4 is a special case.

$\tau_l = 1$ and $\tau_r = 5$.

Possible observed intervals: $[1,1]$, $(-\infty, 2]$, $[5,5]$, $(4, +\infty)$, $(-\infty, 4]$, $(2, +\infty)$.

Possible IIs: $[1,1]$, $(2,4]$, $[5,5]$ (when n is large enough), with weights s_1, s_2, s_3 .

Suppose that there are N_1 $(1, 1)$'s or $(-\infty, 2)$'s, N_2 $(5, 5)$'s or $(4, +\infty)$'s, N_3 $(2, +\infty)$'s and N_4 $(-\infty, 4)$'s among n doubly-censored observations.

Eq. (4.2) of §4.4 becomes

$$\begin{aligned} s_1 &= \frac{N_1 s_1}{n s_1} + \frac{N_2}{n} \cdot 0 + \frac{N_3}{n} \cdot 0 + \frac{N_4}{n} \frac{s_1}{s_1 + s_2} && \text{Can } s_1 = 0 ? \\ s_2 &= \frac{N_1}{n} \cdot 0 + \frac{N_2}{n} \cdot 0 + \frac{N_3}{n} \cdot \frac{s_2}{s_2 + s_3} + \frac{N_4}{n} \frac{s_2}{s_1 + s_2} && \text{Can } s_2 = 0 ? \end{aligned}$$

(we skip 1 equation in (4.2) as $m = 3$ but $s_1 + s_2 + s_m = 1$).

Solving the equations yields two solutions of \mathbf{s} : $(s_1, s_2) = (\frac{N_1+N_4}{n}, 0)$ and $(s_1, s_2) = (\frac{N_1}{N_1+N_3}, U_n)$, where

$$U_n = \frac{N_4}{N_2 + N_4} - \frac{N_1}{N_1 + N_3}.$$

The first one is an SCE of \mathbf{s} as $s_i \geq 0$ and $s_1 + s_2 \leq 1$. The second one is an SCE if and only if $U_n \geq 0$.

Let

1. $H_1 = \frac{N_1+N_4}{n} \mathbf{1}_{(x \geq 1)} + \frac{N_2+N_3}{n} \mathbf{1}_{(x \geq 5)}$.
2. $H_2 = \frac{N_1}{N_1+N_3} \mathbf{1}_{(x \geq 1)} + U_n \mathbf{1}_{(x \geq 4)} + \frac{N_2}{N_2+N_4} \mathbf{1}_{(x \geq 5)}$.

Thus there are possible two SCEs of F_o :

H_1 or

H_2 if $U_n \geq 0$.

The GMLE is $\hat{F} = \begin{cases} H_1 & \text{if } U_n < 0 \\ H_2 & \text{if } U_n \geq 0. \end{cases}$

The reason is as follows.

Theorem 4 *Given interval-censored data (C1, C2, DC or MIC data), let $A_i, i = 1, \dots, m$ be all the distinct IIs. Suppose $\hat{s}_i, i = 1, \dots, m$ are weights assigned to these IIs by an SCE, respectively, and satisfies $\hat{s}_j = 0 \Rightarrow \tilde{s}_j = 0$ for any other SCE with weights $\tilde{s}_i, i = 1, \dots, m$, then $\hat{\mathbf{s}}$ is a GMLE.*

The proof of the theorem utilizes the convexity of $-\ln \mathbf{L}$ in (s_1, \dots, s_k) .

Note that GMLE is also an SCE.

If $U_n < 0$, there is only one SCE of \mathbf{s} , which induces H_1 ,

thus H_1 is a GMLE if $U_n < 0$.

If $U_n \geq 0$, there are exactly two SCEs, which induce H_1 and H_2 .

Then it follows

either from Theorem 4

or from $\mathbf{L}(H_2) > \mathbf{L}(H_1)$

that H_2 is a GMLE if $U_n \geq 0$.

To prove consistency. Verify

AS1: $K(x) = 1/3 \forall x \in (0.5, 2] \cup (2, 4] \cup (4, 8]$, which contains $[1, 5] = [\tau_l, \tau_r]$.

Thus AS1 holds *i.e.* $P(X \in [\tau_l, \tau_r] \text{ and } K(X) = 0) = 0$.

AS2: $F(\tau_l) = 1/2 > 0$ and $P\{L = \tau_l\} = P\{X = 1, (Z, Y) = (0.5, 2)\} = 1/6 > 0$.

$F(\tau_r-) = 1/2 < 1$ and $P\{R = \tau_r\} = P\{X = 5, (Z, Y) = (4, 8)\} = 1/6 > 0$.

Thus AS2 holds.

Thus the SCEs are consistent.

To discuss the asymptotic normality, verify that

(1) AS5 and AS6 fail, as $F_o(2) = p_1 \not\prec F_o(4) = p_1$ and

(2) AS4 fails, as $K(3) = P\{Z < 3 \leq Y\} = g_2 > 0$, but $P_c(3)$ is not defined.

Thus Theorem 2 is not valid. We will verify that the asymptotic covariance matrix of \hat{F} is not J^{-1} given in Theorem 3 with $k = 2$. We now establish the asymptotic normality directly. Note that

$$\begin{aligned} E(N_1/n) &= P\{X = 1, (Z, Y) = (0.5, 2) \text{ or } (2, 4)\} = p_1(g_1 + g_2) \\ E(N_2/n) &= P\{X = 5, (Z, Y) = (4, 8) \text{ or } (2, 4)\} = p_2(g_2 + g_3) \\ E(N_3/n) &= P\{X = 5, (Z, Y) = (0.5, 2)\} = p_2g_1 \\ E(N_4/n) &= P\{X = 1, (Z, Y) = (4, 8)\} = p_1g_3. \end{aligned}$$

We now derive the asymptotic covariance matrix of $(\hat{F}(1), \hat{F}(4))$. By SLLN,

$$U_n \xrightarrow{a.s.} \frac{p_1g_3}{p_1g_3 + p_2(g_2 + g_3)} - \frac{p_1(g_1 + g_2)}{p_1(g_1 + g_2) + p_2g_1},$$

which < 0 as $p_1 = p_2$ and $g_1 = g_2 = g_3$. Thus for n large enough, there is only one SCE and one GMLE. That is H_1 . It is easy to show by SLLN or Theorem 1 that H_1 is consistent. It is also easy to show that

$$\frac{H_1(t) - F_o(t)}{\sigma_{H_1(t)}} \xrightarrow{D} N(0, 1), \quad t \in [1, 5], \quad (5.1)$$

by the corollary of Slutsky's theorem or the CLT.

$$\begin{aligned} H_1(t) &= \begin{cases} (N_1 + N_4)/n & \text{if } t \in [1, 5) \\ 1 & \text{if } t \geq 5 \end{cases} \\ N_1 &= \sum_{i=1}^n \mathbf{1}((L_i, R_i) = (1, 1) \text{ or } (-\infty, 2)) \text{ and } N_4 = \sum_{i=1}^n \mathbf{1}((L_i, R_i) = (-\infty, 4)) \\ H_1(t) &= \begin{cases} \mathbf{1}(X = 1) & \text{if } t \in [1, 5) \\ 1 & \text{if } t \geq 5 \end{cases} \end{aligned}$$

Thus

$$\sigma_{H_1(t)}^2 = \text{Var}\left(\sum_{i=1}^n \mathbf{1}_{(X_i=1)}/n\right) = p_1p_2/n, \quad t \in [1, 5). \quad (5.2)$$

$$\mathcal{L} = N_1 \ln F(1) + N_2 \ln(1 - F(4)) + N_3 \ln(1 - F(1)) + N_4 \ln F(4), \quad \text{where } F(5) = 1.$$

Verify the Fisher information matrix is

$$\begin{aligned} J/n &= -\frac{1}{n} \frac{\partial^2 \mathcal{L}}{\partial(F(1), F(4)) \partial(F(1), F(4))^t} = \begin{pmatrix} \frac{N_1}{nF^2(1)} + \frac{N_3}{n(1-F(1))^2} & 0 \\ 0 & \frac{N_4}{nF^2(4)} + \frac{N_2}{n(1-F(4))^2} \end{pmatrix} \\ &\xrightarrow{a.s.} \begin{pmatrix} \frac{g_1+g_2}{p_1} + \frac{g_1}{p_2} & 0 \\ 0 & \frac{g_3}{p_1} + \frac{g_2+g_3}{p_2} \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \end{aligned}$$

If $(E(J))^{-1}$ is the covariance matrix of $(H_1(1), H_1(4))$, it contradicts (5.2) **Why ??**

§4.5.2. Homework

1. Let (Z, Y) takes values $(0.5, 2)$, $(2, 4)$ and $(4, 8)$, with positive probabilities g_1, g_2 and g_3 , respectively, where $g_1 = g_2 = g_3 = 1/3$, and let $F_o(x) = p_1 \mathbf{1}_{(x \geq 1)} + p_2 \mathbf{1}_{(x \geq 3)} + p_3 \mathbf{1}_{(x \geq 5)}$, where $p_1 = p_2 = p_3 = 1/3$. In the following, you may assume n is sufficiently large.
 - 1.a. Derive all possible SCEs of F_o based on Eq. (4.2) in §4.4.
 - 1.b. Find the limits of the SCEs directly.
 - 1.c. Derive the GMLE of F_o .
 - 1.d. Are the SCEs consistent and asymptotic normal? (Prove or disprove them).
 - 1.e. Derive the asymptotic covariance matrix of the SCEs $(\hat{F}(1), \hat{F}(3))$.
2. Consider Example 1.
 - 2.a. What are the limits of H_1 and H_2 ? (Note that H_i is a function of both sample size n and $t \in R^1$).
 - 2.b. Are they consistent estimators of F_o ?

2.c. Check whether AS1, AS2 and AS3. hold. AS3 can be interpreted as

$$\sum_{i: u_i \in (\tau_l, \tau_r)} \frac{f_Y(u_i) + f_Z(u_i)}{S_Y(u_i) - S_Z(u_i)} < \infty.$$

2.d. Compute the Fisher information matrix $-H = -E \left(\frac{\partial^2 \mathcal{L}}{\partial F(b_i) \partial F(b_j)} \Big|_{F=F_o} \right)_{2 \times 2}$, where $(b_1, b_2) = (1, 4)$ and $\mathcal{L} = \ln L$.

2.e. Compute the covariance matrix of the vector $(H_1(1), H_1(4))$, denoted by Σ . Does $\Sigma = -H^{-1}$? Does it contradicts Theorem 3 ?

3. Prove that under the DC model, $H = F_o$ satisfies the self-consistent equation

$$H(x) = \int_{l \leq x < r} \frac{H(x) - H(l)}{H(r) - H(l)} dQ(l, r) + P\{R \leq x\}, \quad H \in \Theta.$$

4. Under double censoring with $Y = -X$ and $Z = Y - 1$, derive the function $E(\mathbf{1}(X \leq t, L \leq t < R) | (L, R) = (l, r))$ of (l, r, t) for an F_X you select, and show that $H = F_o$ does not satisfies the SE equation.

References

- * Chang, M. N. (1990). Weak convergence of a self-consistent estimator of the survival function with doubly censored data. *Ann. Statist.*, 18, 391-404.
- * Chang, M.N. and Yang, G. (1987). Strong consistency of a self-consistent estimator of the survival function with doubly censored data. *Ann. Statist.*, 15, 1536-1547.
- * Gu, M.G. and Zhang, C-H. (1993). Asymptotic properties of self-consistent estimator based on doubly censored data. *Ann. Statist.*, 21, 611-624.
- * Tsai, W. and Crowley, J. (1985). A large sample study of the generalized maximum likelihood estimators from incomplete data via self-consistency. *Ann. Statist.*, 13, 1317-1334.
- * Turnbull, B. W. (1974). Nonparametric estimation of a survivorship function with doubly censored data. *JASA.* 69, 169-173.
- * Yu, Q.Q. and Li, L.X. (1994). On the strong consistency of the product limit estimator. *Sankhya*, A. 56, 416-430.
- * Yu, Q. Q. and Li, L.X. (2001). Asymptotic properties of the GMLE of self-consistent estimators with doubly-censored data. *Acta Math. Sinica.* 17, 581-594.

§4.6. **SCE under MIC model** Consider MIC model (2). Assume:

- (1) N is a random integer;
- (2) $T, Y_1 < Y_2 < \dots < Y_k < \dots$ are inspection times;
- (3) X and $(N, T, Y_i, i \geq 1)$ are independent;
- (4) $P(N = 0) > 0$ and $P(N > 1) > 0$;
- (5) The observable random vector is

$$(L, R) = \begin{cases} (X, X) & \text{if } X \leq T \text{ and } N = 0 \\ (T, \infty) & \text{if } X > T \text{ and } N = 0 \\ (-\infty, Y_1) & \text{if } X \leq Y_1 \text{ and } N \geq 1 \\ (Y_i, Y_{i+1}) & \text{if } Y_i < X \leq Y_{i+1}, i = 1, \dots, N - 1 \text{ and } N \geq 1 \\ (Y_N, \infty) & \text{if } X > Y_N \text{ and } N \geq 1. \end{cases}$$

Let $(L_i, R_i), i = 1, \dots, n$, be an i.i.d. copies of (L, R) .

Remark If one replaces (4) by $N = 0$ or 1 w.p.1, it can be used to formulate the DC data.

Hereafter, denote H_n an SCE of F_o . Define

$$\tau = \sup\{x : F_o(x) < 1 \text{ and } F_T(x) < 1\},$$

$$\tau_Y = \sup\{t : F_{Y_N}(t) < 1\}.$$

$$\tau_N = \sup\{i : f_N(i) > 0\}.$$

A point x is called a *support point* of a cdf F if $|F(x + \epsilon) - F(x - \epsilon)| > 0 \forall \epsilon > 0$. Denote \mathcal{S}_F the set of all support points of F . People make use of the assumptions as follows:

(AS1) $\tau_Y \leq \tau$;

(AS2) $P\{T \text{ or } Y_N = \tau\} > 0$ if $F_o(\tau-) < 1$.

(AS3) H_n is right continuous and $\mathcal{S}_{H_n} \subset \{R_1, \dots, R_n\}$;

(AS4) $F_o(\tau) > 0$ and $\cup_{i \leq \tau_N} \mathcal{S}_{F_{Y_i}} \subset \mathcal{S}_{F_o}$;

(AS5) $\cup_{i \leq \tau_N} \mathcal{S}_{F_{Y_i}}$ is a finite set and F_o is strictly increasing on $\cup_{i \leq \tau_N} \mathcal{S}_{F_{Y_i}}$.

(AS6) There are at most m IIs for each sample size n and $\mu_{F_o}(A_j) > 0$ for each II A_j .

Theorem 1 (Yu, Li and Wong (1998,2000)) *Suppose that AS1 and AS2 hold. Then the SCE H_n satisfies*

$$\lim_{n \rightarrow \infty} \sup_x |H_n(x) - F_o(x)| = 0 \text{ a.s. if } F_o(\tau) = 1 \text{ and } \lim_{n \rightarrow \infty} \sup_{x \leq \tau} |H_n(x) - F_o(x)| = 0 \text{ a.s..}$$

Theorem 2 (Yu, Li and Wong (1998,2000)) *Suppose that AS1 and AS2 hold. Moreover, either AS5 holds or AS6 holds or AS3 and AS4 hold. Then*

$$\sqrt{n}(H_n(x) - F_o(x)) \xrightarrow{D} Z(x) \text{ for } x \leq \tau,$$

where Z is a Gaussian process on $[0, \tau]$. Let $b_1 < \dots < b_{k+1}$ be the right endpoints of all the distinct IIs induced by the observed intervals. The covariance matrix of $(H_n(b_1), \dots, H_n(b_k))$

can be estimated by $-\left(\frac{\partial^2 \mathcal{L}}{\partial F(b_i) \partial F(b_j)}\right)_{k \times k}^{-1} \Big|_{F=\hat{F}}$ and $\mathcal{L} = \ln L$.

Remark. The assumptions AS1, AS4, AS5 and AS6 eliminate the case in which $\mu_{F_o}(A_j) = 0$ for an II A_j . The following case is a situation when AS4-AS6 are violated.

$X = 3, 6, N = 0, 1, T = 8, Y_1 = 4, 5,$

Then an II is $(4, 5]$ and $\mu_{F_o}((4, 5]) = 0$.

Remark. For IC data, $-\ln L$ is strictly convex in the interior of $\mathbf{s} \in D_s$, where

$D_s = \{\mathbf{s} : \vec{s} = (s_2, \dots, s_m), s_i \geq 0, \sum_{i=2}^m s_i \leq 1\}, s_1 = 1 - s_2 - \dots - s_m.$

Thus there is a unique GMLE of \mathbf{s} . Moreover, if $\hat{\mathbf{s}}$ is the GMLE and $\tilde{\mathbf{s}} \neq \hat{\mathbf{s}}$ is an SCE, then

(1) $\hat{s}_j = 0$ implies $\tilde{s}_j = 0$ for each j ;

(2) there is at least one j such that $\tilde{s}_j = 0 < \hat{s}_j$.

A function g is strictly convex if $g(ax + (1-a)y) < ag(x) + (1-a)g(y)$ for each possible pair of (x, y) and $a \in (0, 1)$. Since $\sum_{j=1}^m \delta_{ij} s_j$ is linear in $\mathbf{s} \in D_s$ and $-\ln x$ is strictly convex

in $x \in R^1$, as $\frac{d^2(-\ln x)}{dx^2} = \frac{1}{x^2} > 0$, $-\ln L$ is strictly convex.

§4.6.2. Homework

1. Consider the MIC model (2). Suppose that

X takes values 1, 3 and 5 w.p.1.;

N takes values 0, 1 and 2 w.p.1.;

$T = 7,$

Y_1 takes values 2 and 4 w.p.1;

$Y_2 = Y_1 + 2.$

1.a. Describe all possible observed intervals.

1.b. Derive all possible SCEs of F_o . (**Hint:** If you properly group the observations, the estimators are similar to H_1 and H_2 in Example 1 or Homework problem 1 of §4.5.)

§4.7. Case 2 data and the GMLE

The mixed case IC model is used to formulate IC data which do not involve exact observations. Let $\mathbf{Y} = \{Y_j : j \geq 1\}$ be a sequence of random variables such that $Y_1 < Y_2 < \dots$.

K be a positive random integer such that $(K, \mathbf{Y}) \perp X$;

On the event $\{K = k\}$, let (L, R) denote the endpoints of that random interval among $(-\infty, Y_1], (Y_1, Y_2], \dots, (Y_k, \infty)$ which contains X .

Let ν be the measure that is the sum of the two measures induced by the marginal distributions of observable extended random variables L and R .

Theorem 1. (Schick and Yu (2000)) *Suppose $E(K) < \infty$. Let $A = \cup_{i \leq \tau_K} \mathcal{S}_{F_{Y_i}}$, where $\tau_K = \sup\{i : P\{K = i\} > 0\}$, and \hat{F} be a GMLE of F_o . Then*

1. $\int |\hat{F} - F_o| d\nu \rightarrow 0$ almost surely;
2. If $P(Y_i = a) > 0$ for some $i \leq \tau_K$, then $\hat{F}(a) \rightarrow F_o(a)$ a.s.;
3. If A is dense in $[0, \infty)$ and F_o is continuous, then $\sup_x |\hat{F}(x) - F_o(x)| \rightarrow 0$ a.s..

Question: Relation between the 3 statements in the theorem ?

Question: Why not SCE ?

With case 2 IC data, the SCE may not be consistent even at discrete inspection times. See the following example:

Example 1. Let X take on values 2, 4, 6 with probability 1/3 for each value. Let $K = 2$ w.p.1 and (Y_1, Y_2) take on values (1,5) and (3,7) with probabilities 1/2 and 1/2, respectively. If n is large enough, the observations are

$$N_1 (-\infty, 3], N_2 (5, \infty), N_3 (3, 7] \text{ and } N_4 (1, 5].$$

$$\text{IIs: } (1, 3], (3, 5] \text{ and } (5, 7]$$

As in Example 1 of §4.5, there are two solutions to equation (4.2):

They induce two SCEs H_1 and \hat{F} almost the same as in Example 1 of §4.5.

$$H_1(t) = \frac{N_1 + N_4}{n} \mathbf{1}_{(t \geq 3)} + \frac{N_2 + N_3}{n} \mathbf{1}_{(t \geq 7)}$$

$$\hat{F}(t) = \frac{N_1}{N_1 + N_3} \mathbf{1}_{(t \geq 3)} + U_n \mathbf{1}_{(t \geq 5)} + \frac{N_2}{N_2 + N_4} \mathbf{1}_{(t \geq 7)}$$

Note that under the assumption there $\hat{F} = H_1$ if n is large enough.

However, in the current situation, $\hat{F} \neq H_1$ if n is large enough, as pointed out next.

Note Y_i 's are discrete, thus we expect a GMLE is consistent at 1, 3, 5, 7.

The limit of $H_1(3)$:

$$\lim_{n \rightarrow \infty} \frac{N_1 + N_4}{n} = P\{X = 2, (Y_1, Y_2) = (3, 7)\} + P\{X \leq 4, (Y_1, Y_2) = (1, 5)\} = 1/2 \neq 1/3.$$

Thus H_1 is not consistent at 3.

Theorem 2. (Yu, Schick, Li and Wong (1998)) *Under the mixed case model, if there are only $k + 1$ IIs with right endpoints b_i for each sample size, F_o is strictly increasing on b_i 's, then*

$$\sqrt{n} \begin{pmatrix} \hat{F}(b_1) - F_o(b_1) \\ \dots \\ \hat{F}(b_k) - F_o(b_k) \end{pmatrix} \xrightarrow{\mathcal{D}} N(0_k, \Sigma_k)$$

as $n \rightarrow \infty$, where

$$\Sigma_k = -n \left(E \left(\frac{\partial^2 \mathcal{L}}{\partial F(b_i) \partial F(b_j)} \Big|_{F=F_o} \right) \right)_{k \times k}^{-1}$$

and the asymptotic covariance matrix of \hat{F} can be estimated by $-\left(\frac{\partial^2 \mathcal{L}}{\partial F(b_i) \partial F(b_j)} \right)_{k \times k}^{-1} \Big|_{F=\hat{F}}$.

The convergence rate for the GMLE

is $n^{1/2}$ under the assumption of Theorem 2;

is $n^{1/3}$ under another set of continuity assumptions given by Groeneboom (1996), and

is conjectured to be $(n \log n)^{1/3}$ under a set of continuity assumptions given by Groeneboom and Wellner (1992).

Under the case 2 model, Groeneboom and Wellner (1992) establish the result as follows.

Suppose that F_o and G have continuous derivatives,

with their densities $f_o(x_o) > 0$ and $g(x_o, x_o) > 0$, and

let \tilde{F} be the estimator of F_o , obtained at the first step of the iterative convex minorant algorithm, starting the iterations with F_o . Then the statistic

$$(\ln n)^{1/3} \frac{\tilde{F}(x_o) - F_o(x_o)}{\{\frac{3}{4}(f_o(x_o))^2 / g(x_o, x_o)\}^{1/3}} \xrightarrow{\mathcal{D}} 2Z^*,$$

where Z^* is the last time where standard two - sided Brownian motion minus the parabola $y(t) = t^2$ reaches its maximum.

Conjecture (G&W (1992, p. 108)): The GMLE \hat{F} has the same asymptotic distribution as \tilde{F} . Thus, the convergence rate of the GMLE is conjectured to be $(n \log n)^{1/3}$ under the same conditions mentioned above.

Remark. \tilde{F} is not an estimator, because F_o is unknown and thus it is impossible to start the iterations from F_o except in simulation.

Theorem 3. (Groeneboom (1996)). Let F_o be continuous with a bounded derivative f_o on $[0, M]$, satisfying $f_o(x) \geq c_o > 0$, $x \in (0, M)$, for some constant $c_o > 0$. Let (Y_1, Y_2) be the two continuous random inspection times in the Case 2 model, with $df g(\cdot, \cdot)$. Let g_1 and g_2 be the first and second marginal density of g , respectively. Suppose that the following conditions are satisfied

- (S1) g_1 and g_2 are continuous, with $g_1(x) + g_2(x) > 0 \forall x \in [0, M]$;
- (S2) $g(\cdot, \cdot)$ is continuous, with uniformly bounded partial derivatives, except at a finite number of points, where left and right (partial) derivatives exist;
- (S3) $P\{Y_2 - Y_1 < \epsilon_o\} = 0$ for some ϵ_o with $0 < \epsilon_o \leq 1/(2M)$, so g does not have mass close to the diagonal.

Then we have at each point $t_o \in (0, M)$

$$n^{1/3}\{2a(t_o)/f_o(t_o)\}^{1/3}\{\hat{F}(t_o) - F_o(t_o)\} \xrightarrow{\mathcal{D}} 2Z^*,$$

where Z^* is defined as in Theorem 2, and

$$a(t_o) = \frac{g_1(t_o)}{F_o(t_o)} + k_1(t_o) + k_2(t_o) + \frac{g_2(t_o)}{1 - F_o(t_o)},$$

$$k_1(u) = \int_u^M \frac{g(u, v)}{F_o(v) - F_o(u)} dv \text{ and } k_2(v) = \int_0^v \frac{g(u, v)}{F_o(v) - F_o(u)} du.$$

Remark. The main differences between the assumptions in the above theorems that the convergence rate varies are as follows.

1. The main assumption in Theorem 2 is that K is finite and Y_i , $i = 1, \dots, K$, takes on finitely many values. Then the convergence rate is $n^{1/2}$.
2. The main assumption in the conjecture of G&W (1992) is that $P\{(Y_1, Y_2) \in N(x_o, x_o, \epsilon)\} > 0$ for each neighborhood $N(x_o, x_o, \epsilon)$ of (x_o, x_o) , in addition to smoothness.

Then the convergence rate is $(n \ln n)^{1/3}$.

3. The main assumption in Theorem 3 is that (Y_1, Y_2) does not fall along a strip near the diagonal $y_1 = y_2$, in addition to smoothness. Then the convergence rate is $n^{1/3}$.

§4.7.2. Homework

1. a. Derive the limits of the SCEs in Example 1 and compare to F_o .
b. Derive the asymptotic variance of the GMLE in Example 1.
c. Give an estimate of the variance in part b.
2. Suppose that $F \sim Exp(\rho)$, $K = 2$ w.p.1., $Y_1 \sim Exp(\rho)$ and $Y_2 = Y_1 + Z$, where $Z \sim Exp(\rho)$ and $Z \perp Y_1$. let S_n^2 be the sample variance of the GMLE $\hat{F}(2)$ based on 1000 simulations of random samples of size n . What do you expect S_{100}^2/S_{400}^2 to be? State your reasoning.
3. Suppose that $F \sim Exp(\rho)$, $K = 2$ w.p.1., Y_1 has a discrete uniform distribution on the set $\{1, 2, \dots, 9\}$ and $Y_2 = Y_1 + 1$. Let S_n^2 be the sample variance of the GMLE $\hat{F}(2)$ based on 1000 simulations of random samples of size n . What do you expect S_{100}^2/S_{400}^2 to be? State your reasoning.
4. Try simulation to # 2 and # 3 with hw10.r
5. Consider the model in Example 1. Try to construct two solutions that satisfy the (population) self-consistent equation in Homework # 3 in §4.5.2. This example shows that unlike the DC model, the solution to the population SE equation is not unique. Moreover, the solution $H(t)$ at $t = R$ is also not unique.

Reference.

- * Groeneboom, P. and Wellner, J. A. (1992). Information bounds and nonparametric maximum likelihood estimation. *Birkhäuser Verlag, Basel*.
- * Groeneboom, P. (1996). Lecture on inverse problems. In P. Bernard. (Ed.), *Lectures on probability and statistics*. p. 157. Berlin, Springer-Verlag.
- * Schick, A. and Yu, Q.Q. (2000). Consistency of the GMLE with mixed case interval-censored data. *Scand. J. Statist.*, 27, 45-55.
- * Yu, Q. Q., Schick, A., Li, L. X. and Wong, G. Y. C. (1998). Asymptotic properties of the GMLE in the case 1 interval-censorship model with discrete inspection times. *Canadian Journal of Statistics*, 26, No. 4, 619-627.

§4.8. Application to model justification.

It is desirable to use parametric models instead of nonparametric models. However, one needs to justify why a certain type of model can be applied to the given data, either from science (physics, biology, etc.) or from empirical data.

A naive method is to plot the empirical cdf against the targetting parametric cdf, replacing the parameters by their MLEs (see ch4.r):

1. Plot the GMLE of $S_o(t)$ and its 95% CI along $t < \tau$ based on the IC data.
2. Plot the MLE of $S_o(t)$ for the potential parametric model based on the IC data.

If the MLE curve lies within the confidence band, the parametric model may fit the data. Here, if the data is complete, GMLE=edf; if data are RC, GMLE=PLE; if data are C1, GMLE is given by the max-min form; otherwise, use the SC-algorithm.

We shall introduce several other methods here.

§4.8.1. Q-Q plot

For a complete data set

X_1, \dots, X_n from a cdf F_o .

If we suspect F_o belongs to a parametric distribution family, say

$$H_0: F_o = F(\cdot, \theta),$$

we can use probability plot to check.

First estimate the parameter θ by $\hat{\theta}$ (MLE etc.).

Plot of sample quantiles v.s. quantiles of $F(\cdot, \hat{\theta})$ (Q-Q plot).

Let $X_{(1)} \leq \dots \leq X_{(n)}$ be order statistics;

They are $100 \frac{1}{n+1}, \dots, 100 \frac{n}{n+1}$ -th sample percentiles (quantiles) of the sample.

Let $y_i = \sup\{u : F(u, \hat{\theta}) < \frac{i}{n+1}\}, i = 1, \dots, n.$

— estimated population percentiles (quantiles).

Plot $(y_i, X_{(i)}), i = 1, \dots, n.$

If the plot is close to the line $y = x$, then $F \approx F_o$.

qqplot() is a function in R.

It only needs the sample and F_o .

Recall the the $100p$ -th quantile of a cdf F is a

$$q = F^{-1}(p), \text{ where } F^{-1}(p) = \sup\{q : F(q) < p\}.$$

We first explain the reasoning of QQ-plot with complete data.

Complete data. Let $a_1 < \dots < a_m$ be distinct points among X_1, \dots, X_n . Since the edf $\hat{F} \rightarrow F_o$ a.s. on $(0, \infty)$, we expect the quantile functions $(\hat{F})^{-1}$ and F_o^{-1} are close. These $(a_i, F_o^{-1}(\hat{F}(a_i))), i = 1, \dots, m$ are around the line $y = x$ or at least around $y = \sigma x + \mu$ if $\sigma \neq 1$ or $\mu \neq 0$ under H_0 .

The following are two QQ-plot graphs. The first one is 100 observations from Exp(1) plotting against normal distribution using qqnorm(). The second is a QQ-plot of 100 observations from normal distribution plotting against normal distribution.

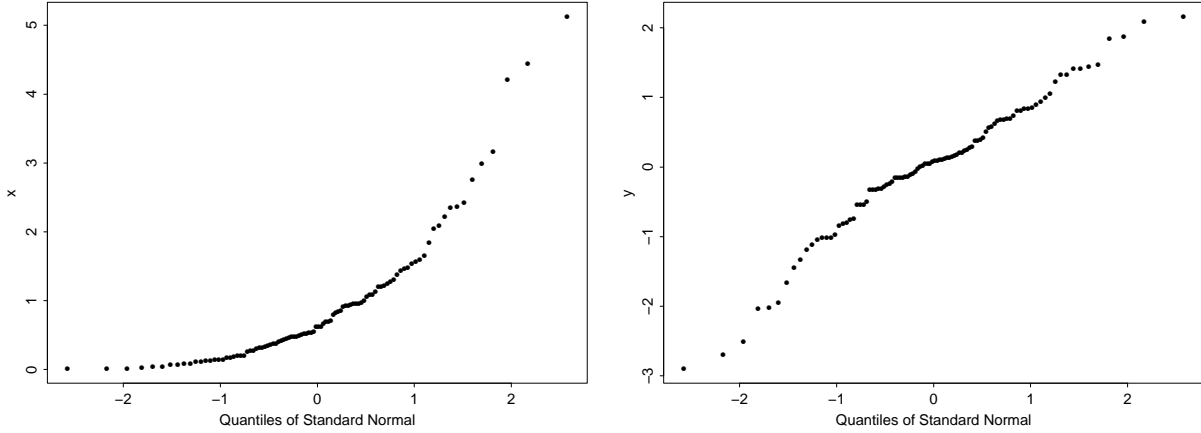


Fig. 1. QQ-plot

For RC or IC data Suppose that $X \sim F_o$ and (L_i, R_i) 's are iid IC observations. If we suspect that F_o is some of given form, say $F_o = F_*(\cdot, \theta)$, we can also do QQ-plot as follows.

1. Obtain the IIs based on observations (L_i, R_i) , $i = 1, \dots, n$, denoted by A_j , $j = 1, \dots, m$.
2. Obtain the GMLE of the distribution function F_o , denoted by \hat{F}_* . The GMLE \hat{F}_* is redefined to be linear on each non-singleton A_j rather than a right continuous step function with jumps only at the right endpoints of the A_j 's, provided the A_j is a finite set.
3. Obtain the MLE of the parameter θ , say $\hat{\theta}$, in F_* .
4. Denote the midpoint of the finite A_j 's by m_i 's. Plot $(m_i, F_*^{-1}(\hat{F}_*(m_i), \hat{\theta}))$'s for all i that A_i is finite, which ideally should also be around $y = x$ or $y = \sigma x + \mu$, in the latter case, it suggests that $F_o(t) = F_*(\frac{t-\mu}{\sigma}, \theta)$.

§4.8.2. Hazard plot for RC data

(Z_i, δ_i) , $i = 1, \dots, n$.

Plot sample integrated hazard v.s. integrated hazards of $F_o(\cdot, \hat{\theta})$.

To get sample integrated hazards:

- a. Order Z_i 's as $Z_{[n]} \leq \dots \leq Z_{[1]}$
(reverse order);
- b. Denote $\delta_{[i]}$, $i = 1, \dots, n$, correspondingly;
- c. For each $\delta_{[k]} = 0$ (RC ones), the sample hazard is 0;
- d. For each $\delta_{[k]} = 1$ (exact ones),
the sample hazard $\hat{h}(Z_{[k]})$ is $1/k$;
- e. The sample integrated hazard at $Z_{[k]}$ is

$$\hat{H}(Z_{[k]}) = - \sum_{i: Z_{[i]} \leq Z_{[k]}, \delta_{[i]}=1} \log[1 - \hat{h}(Z_{[i]})].$$

$$(\hat{S}(t) = \prod_{i: Z_{(i)} < t} (1 - \frac{\delta_{(i)}}{n+1-i}))$$

- f. Plot $(\hat{H}(Z_{[k]}), H_o(Z_{[k]}))$, where $\delta_{[k]} = 1$ and $H_o(x) = -\log S_o(x, \hat{\theta})$.

Example. Hazard plot. Use Sample 0 in Table 1.1 v.s. exponential distribution.

Table 1.2. Calculation of sample hazards				
Patient Number	Remission Time	Reverse Order (K)	$h(x) = f(x)/S(x-)$	Integrated Hazard
2	6	21	1/21	0.15
3	6	20	1/20	0.15
4	6	19	1/19	0.15
1	6+	18		
5	7	17	1/17	0.21
6	9+	16		
8	10	15	1/15	0.28
7	10+	14		
9	11+	13		
10	13	12	1/12	0.37
11	16	11	1/11	0.47
12	17+	10		
13	19+	9		
14	20+	8		
15	22	7	1/7	0.62
16	23	6	1/6	0.80
17	25+	5		
18	32+	4		
19	32+	3		
20	34+	2		
21	35+	1		

It is easy to see that $6+ > 6$.

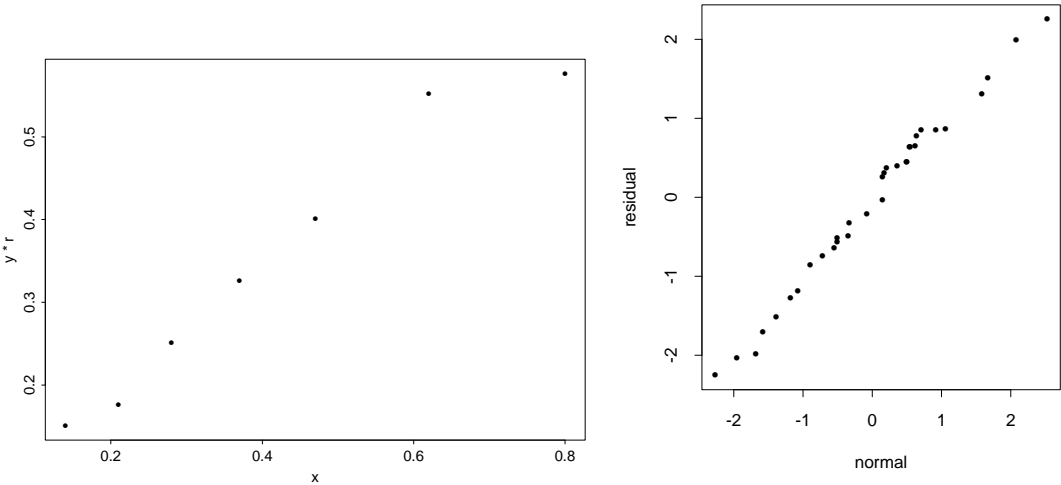


Fig. 2. (Left) : Hazards plot. Fig 3. (Right): Q-Q plot with IC data

Note that the $\hat{H}(6)$ appears three times as there are three exact observations at $x = 6$. However, they are obviously the same. Thus we plot

$$(0.15, 6\rho), (0.21, 7\rho), (0.28, 10\rho), (0.37, 13\rho), (0.47, 16\rho), (0.62, 22\rho), (0.80, 23\rho),$$

$$\text{where } \rho = \sum_{i=1}^n \delta_i / \sum_{i=1}^n Z_i.$$

We explain the reasoning of Hazard plot with RC data.

Let $a_1 < \dots < a_m$ be distinct point among all exact observations. Since $\hat{F}_{pl} \rightarrow F_*$ a.s., we expect sample cumulative hazard $\hat{H} \rightarrow H_o$ a.s. (population cumulative hazard) a.s., thus $(\hat{H}(a_i), H_o(a_i))$'s should be around the line $y = x$ if we know F_o completely.

3 ways for model checking:

- (1) plot the cdf (compare two cdf's curves),
- (2) QQ-plot (check for linearity),
- (3) Hazard plot (check for linearity).

§4.8.3. Diagnostic plot under regression set-up with IC data

In parametric analysis, we often assume that $H_0: F_o(t) = F_*(\frac{t-\mu}{\sigma}) \forall t$, where F_* is known but the location-scale parameters μ and σ may or may not be known.

In general, we can also use the following QQ-plot to check whether the IC data satisfy a certain parametric distribution F_* . We shall consider this in regression set-up. We first explain this method via the following example.

Example 1. If we suspect that conditional \mathbf{Z} , $\ln X \sim N(\beta^t \mathbf{Z}, \sigma^2)$, where β is a $p \times 1$ parameter and X is interval censored, we can use the following procedure to check the assumption, based on our observations (L_i, R_i, \mathbf{z}_i) , $i = 1, \dots, n$.

1. Compute the MLE of β , say $\hat{\beta}$, based on observations (L_i, R_i, \mathbf{z}_i) .
2. Obtain the IIs based on $(\ln L_i - \hat{\beta}^t \mathbf{z}_i, \ln R_i - \hat{\beta}^t \mathbf{z}_i)$, $i = 1, \dots, n$, denoted by A_j , $j = 1, \dots, m$.
3. Obtain the GMLE of the distribution function based on $(\ln L_i - \hat{\beta}^t \mathbf{z}_i, \ln R_i - \hat{\beta}^t \mathbf{z}_i)$, $i = 1, \dots, n$, denoted by \hat{F}_* . The GMLE \hat{F}_* is redefined to be linear on each non-singleton A_j rather than a right continuous step function with jumps only at the right endpoints of the A_j 's, provided the A_j is a finite set.
4. Plot $(m_i, F_*^{-1}(\hat{F}(m_i)))$ for all possible i .

If the assumption is correct, we expect a roughly linear plot.

The justification for the method is as follows.

$\epsilon = \ln X - \beta^t \mathbf{Z}$ conditional on $\mathbf{Z} = \mathbf{z}$ is normal distribution $N(0, \sigma^2)$.

Since X_i is interval censored by the random interval $(L_i, R_i]$,

$\epsilon_i (= \ln X_i - \beta^t \mathbf{z}_i)$ is interval censored by the random interval $(\ln L_i - \beta^t \mathbf{z}_i, \ln R_i - \beta^t \mathbf{z}_i]$,

which can be estimated by $(\ln L_i - \hat{\beta}^t \mathbf{z}_i, \ln R_i - \hat{\beta}^t \mathbf{z}_i]$.

Since the MLE of β and the GMLE of F are consistent under a certain regularity assumptions (see §3 and 4),

it should be approximately normally distributed.

Fig 3. is a result from a simulation study.

We assume that conditional on $\mathbf{Z} = \mathbf{z}$, $\ln X$ has a normal distribution $(N(\beta^t \mathbf{z}, 1))$,

where $\beta = (1, 1, 1)$. $\mathbf{Z} = (Z_1, Z_2, Z_3)$, and

Z_1, Z_2 and Z_3 equal 0 and ± 1 with a certain probabilities.

X is under a mixed case interval censorship model.

The number of follow-up times K is a discrete uniform distribution on $\{1, 2, \dots, 28\}$.

Conditional on $K = k$, the follow-up time Y_i , $i = 1, \dots, k$, satisfy $\ln Y_i = -5 + \sum_{j=1}^i U_j$, where U_i are i.i.d. from uniform distribution $U(0, 1)$.

Here $n = 374$. The MLE of (β, σ) is $(0.92, 1.02, 1.02, 1.06)$.

The resulting QQ-plot supports the normal regression model as we expected.

The method can be viewed as a pivotal method by choosing a function T of (X, \mathbf{Z}, β) so that $T = T(X, \mathbf{Z}, \beta)$ has a distribution function F_T which does not depend on \mathbf{Z} and T is strictly increasing in X . For lognormal,

$$T = \log X - \beta^t \mathbf{Z}.$$

In the above cases, $T \sim N(0, \sigma^2)$.

In general, the procedure is as follows:

1. Find a pivotal function T described as above.
2. Find the MLE $\hat{\beta}$ of parameter β .
3. Obtain the IIs based on $(T(L_i, \mathbf{z}_i, \hat{\beta}), T(R_i, \mathbf{z}_i, \hat{\beta}))$, $i = 1, \dots, n$, denoted by A_j , $j = 1, \dots, m$.

4. Obtain the GMLE of the distribution function based on $(T(L_i, \mathbf{z}_i, \hat{\beta}), T(R_i, \mathbf{z}_i, \hat{\beta}))$, $i = 1, \dots, n$, denoted by \hat{F}_* . The GMLE \hat{F}_* is redefined to be linear on each non-singleton A_j , provided the A_j is a finite set.
 5. Plot $(m_i, F_*^{-1}(\hat{F}_*(m_i)))$, for all possible i .
- If the assumption is correct, we expect a roughly linear plot.

Example 2 Suppose that conditional on $\mathbf{Z} = \mathbf{z}$, X has a Weibull distribution. That is,

$$S_X(x) = e^{-x^\kappa e^{\beta^t \mathbf{z}}}, \quad x > 0.$$

One of such pivoting functions is

$$T = X^\kappa e^{\beta^t \mathbf{z}},$$

where T has an Exponential distribution with survival function $S_T(t) = e^{-t}$, $t > 0$. If X_i is interval censored by $(L_i, R_i]$, then T_i is interval censored by $(T(L_i, \mathbf{z}_i, (\beta, \kappa)), T(R_i, \mathbf{z}_i, (\beta, \kappa))) = (L_i^\kappa e^{\beta^t \mathbf{z}_i}, R_i^\kappa e^{\beta^t \mathbf{z}_i}]$. Since (β, κ) are unknown and neither are $L_{i, \beta, \kappa}$ and $R_{i, \beta, \kappa}$, we replace β and κ by their MLE. Let \hat{F}_T be the GMLE of F_T based on the pivoted data $(T(L_i, \mathbf{z}_i, (\hat{\beta}, \hat{\kappa})), T(R_i, \mathbf{z}_i, (\hat{\beta}, \hat{\kappa})))$, $i = 1, \dots, n$. Finally we plot \hat{F}_* against F_T , the exponential distribution $\text{Exp}(1)$.

§4.8.3.2. Homework:

In the following, make comments on whether the plots suggest $F_o(x) = F_*(\frac{x-\mu}{\sigma})$.

1. The Weibull distribution in Example 2 can also be re-parametrized as a location-scale parameter family. Find the pivot function T and derive the distribution function F_T . What is the revision of the procedure for a diagnostic plot in Example 2 ?
2. Do a QQ-plot using sample 0 in Leukaemia data on page 2 v.s. exponential distribution $\text{Exp}(1)$. Do you think the exponential distribution is appropriate ? If so, what is your guess of ρ according to the slope of fitting straight line ?
3. Q-Q plot:
 - 3.a. Use sample 1 in Leukaemia data on page 2 v.s. exponential distribution.
 - 3.b. Generated a random sample of size 100 from an exponential distribution v.s. exponential distribution, normal distribution.
4. Hazard plot. Generate a random sample of size 100 from a RC model, say $X \sim$ Weibull distribution (nontrivial one), $Y \sim$ Uniform distribution, Draw a hazard plot v.s. Weibull (that you used) and a hazard plot v.s. a normal distribution with the mean and variance equal those of the Weibull distribution you used.
5. Use sample 1 in Leukaemia data on page 2 to check whether Weibull distribution, exponential distribution, lognormal distribution, or log-logistic distribution is appropriate for the data using the idea in the R-program in ch4.r.

Chapter 5. Semi-parametric Analysis

§5.1. Introduction

Suppose

- X is a random survival time,
- \mathbf{Z} is a $p \times 1$ covariate (explanatory) (random) vector
(which sometimes is assumed to be nonrandom);
- X is subject to interval censoring;
- Observable random vector is (L, R, \mathbf{Z}) .

Example 1. Cancer research. In addition to observe the failure time of a patient, we also observe $\mathbf{Z}^t = (Z_1, Z_2, Z_3, Z_4)$, where

- Z_1 — # of relatives who had cancer;
- Z_2 — age of the patient;
- Z_3 — tumor size;

Z_4 — smoking habit.

Example 2. Two-sample problem. There are two independent samples,

$\mathbf{Z} = \mathbf{1}$ (patient is from sample 1).

Two models are considered, among other models.

1. Proportional hazards (PH) model: Conditional on $\mathbf{Z} = \mathbf{z}$, the hazard function

$$h(t|\mathbf{z}) = \psi(\beta, \mathbf{z})h_o(t), \quad t < \tau, \quad (1.1)$$

where $\tau = \sup\{t : S_o(t) > 0\}$, h_o is a (baseline) hazard function, and ψ is a function of (β, \mathbf{z}) . If S_o is a survival function of a continuous random variable, then

$$S_{X|\mathbf{Z}}(t|\mathbf{z}) = (S_o(t))^{\psi(\beta, \mathbf{z})} \quad (H_{X|\mathbf{Z}}(t|\mathbf{z}) = \psi(\beta, \mathbf{z})H_o(t)) \quad (1.2)$$

where S_o is a (baseline) survival function.

The PH model is also called the *Cox regression model*.

For an arbitrary random variable,

Eq. (1.1) defines a PH family (or model), and

Eq. (1.2) defines a Lehmann family or proportional integrated hazards family.

If S is absolutely continuous, then equations (1.1) and (1.2) are the same.

$$(F(x) - F(a)) = \int_a^x F'(t)dt \quad \forall x$$

2. Accelerated lifetime model: Conditional on $\mathbf{Z} = \mathbf{z}$,

$$X = X_o/\psi(\beta, \mathbf{z}) \quad (\ln X = \ln X_o - \ln \psi(\beta, \mathbf{z})). \quad (1.3)$$

$$S_{X|\mathbf{Z}}(t|\mathbf{z}) = S_o(\psi(\beta, \mathbf{z})t) \quad (S_o = S_{X_o}).$$

Under either model, it is a parametric problem, if S_o is known, and is of a parametric form; otherwise, it is a semi-parametric one (*i.e.*, S_o is an arbitrary survival function).

Most usual forms of ψ : ($\psi \geq 0$)

$e^{\beta^t \mathbf{z}}$ ——— log linear;

$\log(1 + e^{\beta^t \mathbf{z}})$ ——— logistic.

§5.2. PH model with RC data.

§5.2.1. Continuous RC data.

Assume:

Conditional on $\mathbf{Z} = \mathbf{z}$, $X \sim F(\cdot|\mathbf{z})$ with its hazard satisfies (1.1);

$Y \sim G$;

(X, \mathbf{Z}) and Y are independent;

F and G are absolutely continuous;

Observe $(M, \delta, \mathbf{Z}) = (\min\{X, Y\}, \mathbf{1}_{(X \leq Y)}, \mathbf{Z})$.

Let $(M_i, \delta_i, \mathbf{z}_i)$, $i = 1, \dots, n$ be i.i.d. copies of (M, δ, \mathbf{Z}) .

The log likelihood function is

$$\begin{aligned} \mathcal{L}(\beta) &= \ln \prod_{i=1}^n (f(M_i|\mathbf{z}_i))^{\mathbf{1}_{e,i}} (S(M_i|\mathbf{z}_i))^{\mathbf{1}_{r,i}} \\ &= \ln \left[\prod_{i: ex} h(M_i|\mathbf{z}_i) \prod_{i=1}^n S(M_i|\mathbf{z}_i) \right] \\ &= \sum_{i: ex, M_i < \tau} \ln \psi(\beta, \mathbf{z}_i) + \sum_{i: ex, M_i < \tau} \ln h_o(M_i) + \sum_{i=1}^n \psi(\beta, \mathbf{z}_i) \ln S_o(M_i). \end{aligned}$$

This approach needs to estimate β and S_o in the same time. Cox (1972) uses a conditional probability approach and a partial likelihood approach with some assumptions to define a new likelihood function, which only involves β . We first give the likelihood functions and then introduce the derivation.

Notation:

$a_1 < \dots < a_g$ — all the distinct exact observations.

By rearranging the index, assume $X_i = a_i$, $i = 1, \dots, g$.

$\mathcal{R}_j = \mathcal{R}(a_j) = \{i : M_i \geq a_j\}$,

$\phi(i) = \psi(\beta, \mathbf{z}_i)$,

$\mathcal{D} = \{i : \delta_i = 1, i = 1, \dots, n\}$ (note that all exact observations are distinct),

Define a modified likelihood function

$$lik = \prod_{i \in \mathcal{D}} \frac{\phi(i)}{\sum_{k \in \mathcal{R}_i} \phi(k)}. \quad (2.1)$$

The log likelihood

$$l(\beta) = \ln lik = \sum_{i \in \mathcal{D}} [\ln \phi(i) - \ln \sum_{k \in \mathcal{R}_i} \phi(k)]. \quad (2.2)$$

The Maximum Partial likelihood estimator MPLE $\hat{\beta}$ of β is a value of b that maximizes $l(b)$.

Remark. Under certain assumptions, the MPLE $\hat{\beta}$ is consistent and asymptotically normally distributed. Its covariance matrix can be estimated by $-\left(\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^t} \Big|_{\beta=\hat{\beta}}\right)^{-1}$.

Hereafter to the end of §5.2, let

$$\psi(\beta, \mathbf{z}) = e^{\beta^t \mathbf{z}}. \quad (2.3)$$

Then (2.2) becomes

$$l(\beta) = \ln lik = \sum_{i \in \mathcal{D}} [\beta^t \mathbf{z}_i - \ln \sum_{k \in \mathcal{R}_i} e^{\beta^t \mathbf{z}_k}].$$

Example 1. 5 observations $(M_i, \delta_i, \mathbf{z}_i)$'s: $(2.5, 1, 2)$, $(2, 0, 5)$, $(4, 0, 1)$, $(1, 1, 1)$, $(7, 1, 2)$.
lik= ?

Sol: Reorder $(M_i, \delta_i, \mathbf{z}_i)$'s as

$(1, 1, 1), (2, 0, 5), (2.5, 1, 2), (4, 0, 1), (7, 1, 2)$.

Rearrange (according to exact obs.) $(M_i, \delta_i, \mathbf{z}_i)$'s as

$(1, 1, 1), (2.5, 1, 2), (7, 1, 2), (2, 0, 5), (4, 0, 1)$.

We use this order from now on.

3 exact observations: $(a_1, a_2, a_3) = (1, 2.5, 7)$.

$\mathcal{D} = \{1, 2, 3\}$,

$\mathcal{R}_1 = \mathcal{R}(1) = \{1, 2, 3, 4, 5\}$, $\mathcal{R}_2 = \mathcal{R}(2.5) = \{2, 3, 5\}$, $\mathcal{R}_3 = \mathcal{R}(7) = \{3\}$;

$\phi(i)$'s are e^β , $e^{2\beta}$, $e^{2\beta}$, $e^{5\beta}$, e^β ,

$$lik = \frac{e^\beta}{2e^{2\beta} + e^{5\beta} + 2e^\beta} \cdot \frac{e^{2\beta}}{2e^{2\beta} + e^\beta} \cdot \frac{e^{2\beta}}{e^{2\beta}} = \frac{1}{(2 + e^{-\beta})(2e^\beta + e^{4\beta} + 2)}.$$

§5.2.1.2. Homework

1. Derive the estimate of the covariance matrix of the MPLE $\hat{\beta}$ under the assumptions in this section and assuming (2.3).
2. Derive the MPLE of β in Example 1.
3. Construct a level-0.10 two-sided test for $H_0 : \beta_2 = 0$, where $\beta^t = (\beta_1, \beta_2)$. Give the expression as explicitly as possible.

§5.2.2. Discrete RC data.

If F_o and G are continuous, the order statistics of the observations satisfy

$$M_{(1)} < M_{(2)} < \dots < M_{(n)}.$$

In this section, we consider the case that there are ties in the observations. Using the idea of conditional probability, Cox suggests a likelihood function for the discrete RC data as follows.

Notations:

$a_1 < \dots < a_g$ are all distinct exact observations;
 d_j is the # of deaths at a_j ;
 \mathcal{S}_j is the collection of all the combinations of selecting d_j elements out of those in $\mathcal{R}(a_j)$;
 $\mathcal{D}_j = \{i : \delta_i = 1, M_i = a_j\}$;
 $r_j = |\mathcal{R}(a_j)|$;
 $\mathcal{R}(a_j)$ and ϕ are the same as in § 5.2.1.
 A log likelihood is defined as

$$l(\beta) = \ln \prod_{j=1}^g \frac{\prod_{i \in \mathcal{D}_j} \phi(i)}{\sum_{(i_1, \dots, i_{d_j}) \in \mathcal{S}_j} \phi(i_1) \cdots \phi(i_{d_j})}.$$

The MPLE $\hat{\beta}$ of β maximizes $l(\beta)$.

Remark. The likelihood $l(\beta)$ is actually equivalent to

$$l(\beta) = \begin{cases} \ln \prod_{j=1}^g \frac{\prod_{i \in \mathcal{D}_j} \phi(i)}{\sum_{(i_1, \dots, i_{d_j}) \in \mathcal{S}_j} \phi(i_1) \cdots \phi(i_{d_j})} & \text{if } \delta_{(n)} = 0; \\ \ln \prod_{j=1}^{g-1} \frac{\prod_{i \in \mathcal{D}_j} \phi(i)}{\sum_{(i_1, \dots, i_{d_j}) \in \mathcal{S}_j} \phi(i_1) \cdots \phi(i_{d_j})} & \text{if } \delta_{(n)} = 1. \end{cases}$$

as $h(x) = \psi h_o(x)$, $x < \tau$, where $\tau = \sup\{t : S_o(t) > 0\}$,
 as well as the last factor is 1 if $\delta_{(n)} = 1$.

Remark. For discrete r.v., the form of log likelihood function

$$\mathcal{L}(\beta) = \sum_{i: ex} \ln \psi(\beta, \mathbf{z}_i) + \sum_{i: ex} \ln h_o(M_i) + \sum_{i=1}^n \psi(\beta, \mathbf{z}_i) \ln S_o(M_i)$$

is not applicable, as $h(t) = \frac{f(t)}{S(t-)} \neq \frac{f(t)}{S(t)}$ and $S(t|\mathbf{z}) \neq (S_o(t))^{\exp(\beta \mathbf{z})}$.

Remark. Under certain assumptions, the MPLE $\hat{\beta}$ is consistent and asymptotically normally distributed. Its covariance matrix can be estimated by $-\left(\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^t} \Big|_{\beta=\hat{\beta}}\right)^{-1}$.

§5.2.2.2. Homework

1. Suppose that ψ is log linear and there are 6 data: $(M_i, \delta_i, (z_1, z_2))$'s are

$$(4, 1, (1, 0)), (1, 0, (1, 1)), (5, 1, (0, 1)), (4, 1, (1, 0)), (1, 1, (0, 0)), (4, 0, (0, 1)).$$

Compute the MPLE of $\hat{\beta}$ if it exists.

§5.2.3. Conditional probability approach in continuous case.

By continuous assumption, there are no ties in exact observations.

Since we shall make use of independence, we would not reorder the M_i 's.

For each uncensored time $a_j = M_i$,

denote $\mathcal{R}_{i_j} = \mathcal{R}(a_j)$, $j = 1, \dots, g$,

where $g = \sum_{i=1}^n \delta_i$.

Cox made the following assumptions:

1. there is only one death at a_j and thus we can pretend that there is only one death in $[a_j, a_j + \epsilon)$ for a small ϵ ;
2. death only occurs in $[a_j, a_j + \epsilon)$, $j = 1, \dots, g$.

Then

$$\begin{aligned}
 & P\{\text{the } i\text{-th person died in } [a_j, a_j + \epsilon) | \text{the } i\text{-th person survived time } a_j\} \\
 &= P(X_i \in [a_j, a_j + \epsilon) | X_i \geq a_j) \\
 &= P(X_i \in [a_j, a_j + \epsilon)) / P(X_i \geq a_j) \\
 &\approx h(a_j | \mathbf{z}_i) \epsilon,
 \end{aligned} \tag{0}$$

$$\begin{aligned}
&= P(A_{i_1}|B_{i_1})P(A_{i_2}|A_{i_1}B_{i_1}B_{i_2}) \cdots P(A_{i_g}|A_{i_1} \cdots A_{i_{g-1}}B_{i_1} \cdots B_{i_g}) \\
&\quad \times P(B_{i_1})P(B_{i_2}|A_{i_1}B_{i_1}) \cdots P(B_{i_g}|A_{i_1} \cdots A_{i_{g-1}}B_{i_1} \cdots B_{i_{g-1}}) \\
&= P(A_{i_1}|B_{i_1})P(A_{i_2}|B_{i_2}) \cdots P(A_{i_g}|B_{i_g}) \quad (=lik) \\
&\quad \times P(B_{i_1})P(B_{i_2}|A_{i_1}B_{i_1}) \cdots P(B_{i_g}|A_{i_1} \cdots A_{i_{g-1}}B_{i_1} \cdots B_{i_{g-1}})
\end{aligned}$$

by independence.

$$lik = P(A_{i_1}|B_{i_1})P(A_{i_2}|B_{i_2}) \cdots P(A_{i_g}|B_{i_g}). \quad ((\text{see (3) in } \S 5.2.3))$$

Thus *lik* is also called the *partial likelihood* by Cox.

§5.2.5. Nonparametric estimation of S_o .

Under the PH model: $h(t|\mathbf{z}) = \psi(\beta, \mathbf{z})h_o(t)$ and

if X is continuous then $S(t|\mathbf{z}) = (S_o(t))^{\psi(\beta, \mathbf{z})}$.

Baseline integrated hazard:

$$H_o(t) = \int_0^t h_o(u)du \quad (\text{cts})$$

$$\text{or } \sum_{u \leq t} \ln(1 - h_o(u)) \approx \sum_{u \leq t} h_o(u) \quad (\text{discrete}).$$

$$S_o(t) = \exp(-H_o(t)).$$

There is no MPLE of S_o under Cox's assumption, though

Cox's maximum partial likelihood estimator (MPLE) of β is a semi-parametric estimator.

Cox proposes an estimator

$$\hat{S}_o(t) = \exp\left(-\sum_{a_j \leq t} \frac{d_j}{\sum_{l \in \mathcal{R}(a_j)} \hat{\phi}(l)}\right).$$

Note that if $\phi(i) = 1$,

$$\hat{S}_o(t) = \exp\left(-\sum_{a_j \leq t} \frac{d_j}{r_j}\right) = \exp\left(-\sum_{a_j \leq t} \hat{h}(a_j)\right).$$

Breslow (1972, JRSS,B) also proposes another estimator. They can be computed by R codes:

```

library(MASS)
library(survival)
u=coxph(Surv(m,d)~ x)
y=survfit(u)
plot(y)

```

The third estimator is the SMLE which maximizes the likelihood function directly, where

$$L_o(F) = \prod_{i=1}^n (S(M_i - |z_i) - S(M_i|z_i))^{\delta_i} (S(M_i|z_i))^{1-\delta_i}.$$

which corresponding to the PH model, and

$$L(F) = \prod_{i: rc} (S_o(M_i))^{\psi(\beta, \mathbf{Z}_i)} \prod_{i: ex} [(S_o(M_i-))^{\psi(\beta, \mathbf{Z}_i)} - (S_o(M_i))^{\psi(\beta, \mathbf{Z}_i)}].$$

which corresponding to the Lehmann model. They are the same if X is continuous.

§5.3. Extension of PH model with IC data

The conditional probability approach does not work for IC data under the PH model. Finkelstein (1986, Biometrics) first considers the following extension of the PH model to IC data. She defines the SMLE of β to be the one that maximizes the likelihood function $L(\beta, F_o) =$

$$\prod_{i: L_i < R_i} [(S_o(L_i))^{\psi(\beta, \mathbf{Z}_i)} - (S_o(R_i))^{\psi(\beta, \mathbf{Z}_i)}] \prod_{i: L_i = R_i} [(S_o(L_i-))^{\psi(\beta, \mathbf{Z}_i)} - (S_o(R_i))^{\psi(\beta, \mathbf{Z}_i)}]. \quad (5.3.1)$$

β and S_o need to be estimated simultaneously.

Remark. The likelihood corresponds to proportional integrated hazards model, not really the PH model, unless one assumes that X is continuous. If X is discrete, then the likelihood under the PH model is different.

Since $h(t) = f(t)/S(t-)$ and $S(t) = \prod_{x_i \leq t, x_i \in D_f} (1 - h(x_i))$,

$f(t|\mathbf{z}) = h(t|\mathbf{z})S(t-|\mathbf{z})$, where D_f is the set of points at which $f > 0$), the correct generalized likelihood of the PH model with discrete IC data is

$$\mathbf{L}(\beta, h_o) = \prod_{i: L_i < R_i} [S(L_i|\mathbf{z}_i) - S(R_i|\mathbf{z}_i)] \times \prod_{i: L_i = R_i} [h(L_i|\mathbf{z}_i)S(L_i - |\mathbf{z}_i)] \quad (5.3.2)$$

where $h(t|\mathbf{z}) = \begin{cases} \exp(\beta\mathbf{z})h_o(t) & \text{if } t < \tau \text{ or } t < M_{(n)} \\ 1 & \text{if } t = \tau \text{ or } t = M_{(n)} \end{cases}$,

$S(t|\mathbf{z}) = \prod_{x_i \leq t, x_i \in D_{f_o}} (1 - h(x_i|\mathbf{z}))$, and

$S(t-|\mathbf{z}) = \prod_{x_i < t, x_i \in D_{f_o}} (1 - h(x_i|\mathbf{z}))$.

Notice that $h_o(t) = f_o(t)/S_o(t-)$ and $S_o(t-) = \prod_{x_i < t, x_i \in D_{f_o}} (1 - h_o(x_i))$.

Remark. The likelihood in (5.3.2) is actually applicable for both continuous and discrete $X|\mathbf{z}$, though (5.3.2) and (5.3.1) will result in different estimates (due to $S(L_i-)$ in 5.3.2)).

In order to compute the SMLE of S_o in both the PIH model (5.3.1) or the PH model (5.3.2), let A_1, \dots, A_m be the Π 's induced by I_i 's, the observed intervals, and $p_j = \mu_F(A_j)$. Consider S_o of form

$$S_o(t) = \sum_{j: A_j \cap (t, \infty] \neq \emptyset} p_j.$$

The variance of the SMLE $\hat{\beta}$ can be estimated by a $p \times p$ matrix V_{11} , where

$$\begin{pmatrix} V_{11} & V_{12} \\ V_{12}^t & V_{22} \end{pmatrix} = J^{-1},$$

where J is the empirical Fisher information matrix, i.e., minus the second partial derivatives matrix of the log likelihood function with respect to $(\beta^t, p_1, \dots, p_{m-1})$, with all parameters replaced by their SMLE.

Finkelstein (1986, Biometrics) suggests to use the Newton-Raphson method to compute the SMLE of (β, S_o) . It turns out most of the time, the approach does not work. The reason is that it is often that the SMLE of some p_j is zero and the algorithm will lead to a point which is not an SMLE but $p_j = 0$ for some j . The derivative $\frac{\partial \mathbf{L}}{\partial p_j} < 0$. Then the algorithm has to stop. The following is such a counterexample.

Example 1. Consider fitting the PIH model with 5 observations

$(L_i, R_i, Z_i): (2, 5, 0), (3, 4, 0), (5, 9, 1), (1, 6, 1), (7, 8, 0)$.

It can be viewed as data from two groups, corresponding to $Z_i = 0$ or 1.

Then, the innermost intervals are $(3, 4]$, $(5, 6]$ and $(7, 8]$.

Let the weights on these innermost intervals be

p_1, p_2 and p_3 , with $p_1 + p_2 + p_3 = 1$ and $p_i \geq 0$.

Note that the baseline survival function S satisfies

$$\begin{aligned} S(4-) &= 1, \\ S(4) &= S(6-) = p_2 + p_3, \\ S(6) &= S(8-) = p_3 \text{ and} \\ S(8) &= 0. \end{aligned}$$

For this example, it is more convenient to express the likelihood as a function of p_i 's rather than S .

The likelihood is

$$\mathbf{L} = p_1 p_1 (p_2 + p_3)^{e^\beta} (1 - p_3^{e^\beta}) p_3.$$

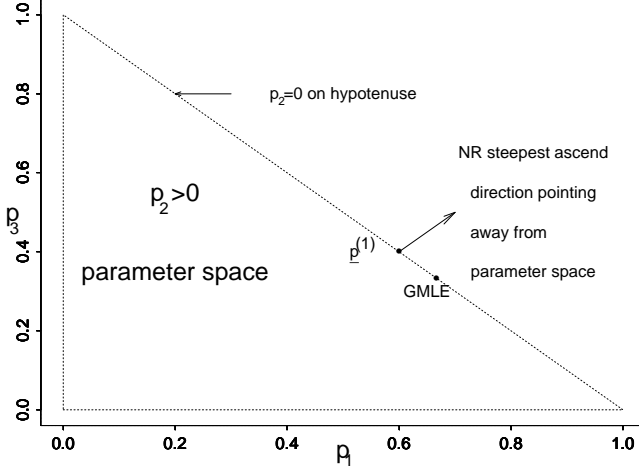
Since $p_1 + p_2 + p_3 = 1$, in view of L , it is simpler to write the log likelihood as

$$\mathcal{L} = \log[p_1^2 p_3 (1 - p_1)^{e^\beta} (1 - p_3^{e^\beta})].$$

The parameter space is $\Omega = \{(\beta, p_1, p_3) : \beta \in (-\infty, \infty), p_1 \geq 0, p_3 \geq 0, p_1 + p_3 \leq 1\}$ with $p_2 = 1 - p_1 - p_3$. For convenience, we write $\alpha = e^\beta$ hereafter. Thus,

$$\mathcal{L} = 2 \log p_1 + \log p_3 + \alpha \log(1 - p_1) + \log(1 - p_3^\alpha).$$

Since the likelihood function has only three variables, it can be shown by direct derivation that the GMLE of (β, p_1, p_2, p_3) is approximately $(-0.461, 2/3, 0, 1/3)$.



A feasible algorithm for the GMLE. Abusing notations, we identify S with a vector (S_1, \dots, S_m) . Similarly, we identify $S^{(i)}$ with $(S_1^{(i)}, \dots, S_m^{(i)})$.

Step 0. Let $b^{(0)} = 0$ be the initial estimate of β and the GMLE of a survival function with observations (L_j, R_j) , $j = 1, \dots, n$ be the initial estimate of $S^{(0)}$.

Step $i + 1$ ($i \geq 0$). Let $b^{(i)}$ and $S^{(i)}$ be the updated values of b and S at Step i .

Do b -step and S -step as follows.

* (**b -step**) With $S = S^{(i)}$ fixed, find a b so that the log likelihood function $\mathcal{L}(S^{(i)}, \cdot)$ increases. Denote the up-dated estimate b by $b^{(i+1)}$. In particular, one can use the NR method to obtain the maximum point b of the log likelihood function with the given $S = S^{(i)}$.

* (**S -step**) With $b = b^{(i+1)}$ fixed, search a non-increasing S so that the log likelihood function $\mathcal{L}(\cdot, b^{(i+1)})$ is maximized (or increases). Since $S_i = p_{i+1} + \dots + p_m$ for some i , let $p^{(i+1),0} = p^{(i)}$. At Sub-step j ($j = 1, \dots, m$), update (p_1, \dots, p_m) by $(p_1^{(i+1),j}, \dots, p_m^{(i+1),j})$, where $p_h^{(i+1),j} = p_{h,u_o}$ and

$$p_{h,u} = \begin{cases} \frac{p_h^{(i+1),j-1} + u}{1+u} & \text{if } h = j, \\ \frac{p_h^{(i+1),j-1}}{1+u} & \text{if } h \neq j, \end{cases}$$

$h = 1, \dots, m$, $u_o > 0$ is a number maximizing $L(b^{(i+1)}, S_{\cdot,u})$ where $S_{\cdot,u} = (S_{1,u}, \dots, S_{m,u})$ and $S_{i,u} = p_{i+1,u} + \dots + p_{m,u}$.

Note: If such u_o is difficult to choose, one may choose a u_o satisfying

$$L(b^{(i+1)}, S^{(i+1),j}) > L(b^{(i+1)}, S^{(i+1),j-1}). \quad (3.1)$$

In particular, if $\frac{\partial}{\partial u} \ln L(b^{(i+1)}, S_{\cdot,u})|_{u=0} > 0$, $u_o = c^k \frac{\partial}{\partial u} \ln L(b^{(i+1)}, S_{\cdot,u})|_{u=0}$, where $S_{\cdot,u} = (S_{1,u}, \dots, S_{m,u})$, $c \in (0,1)$, and k is the smallest non-negative integer such that Inequality (3.1) holds.

Stop at convergence.

Remark. The restriction $u > 0$ can be replaced by $u > -p_h^{(i+1),j-1}$.

If $X|\mathbf{z}$ is not continuous, (5.3.1) is not the likelihood function of the PH model.

Now consider fitting Cox's regression model (5.3.2).

First compute $S(L_i|\mathbf{z}_i) - S(R_i|\mathbf{z}_i)$ in the following table.

(L, R, z)	$S(L), S(R)$	$S(L z) - S(R z)$	<i>simplification</i>
(2, 5, 0)	1, $S(4 0)$	$1 - (p_2 + p_3)$	p_1
(3, 4, 0)	1, $S(4 0)$	$1 - (p_2 + p_3)$	p_1
(5, 9, 1)	$S(4 1), 0$	$(1 - e^\beta p_1) - 0$	$1 - e^\beta p_1$
(1, 6, 1)	1, $S(6 1)$	$1 - (1 - e^\beta p_1)(1 - e^\beta \frac{p_2}{p_2+p_3})$	$1 - (1 - e^\beta p_1)(1 - e^\beta \frac{p_2}{1-p_1})$
(7, 8, 0)	$S(6 0), 0$	$p_3 - 0$	$1 - p_1 - p_2$

The likelihood is $\mathcal{L} = \ln \prod_i (S(L_i|\mathbf{z}_i) - S(R_i|\mathbf{z}_i)) =$

$$2\ln p_1 + \ln(1 - p_1 - p_2) + \ln[1 - e^\beta p_1] + \ln[1 - (1 - e^\beta p_1)(1 - e^\beta \frac{p_2}{1-p_1})]$$

§5.3.2. Homework.

1. Verify the GMLE of (β, p_1, p_2, p_3) for the data related to the figure is approximately $(-0.461, 2/3, 0, 1/3)$.
You do not need to use the algorithm mentioned above.
2. Show that the GMLE of (β, p_1, p_2, p_3) under likelihood (5.3.2) is $(-0.288, 2/3, 0, 1/3)$ ($e^\beta = 3/4$).

§5.4. Accelerated lifetime model and regression analysis with IC data

Assume that conditional on $\mathbf{Z} = \mathbf{z}$,

$$X = T/e^{\beta t \mathbf{z}}.$$

For simplicity, we consider $p = 1$. Then conditional on $Z = z$, $\ln T = \ln X + \beta z$ has a distribution which does not depend on β and Z and has the mean α and variance σ^2 . For simplicity, we shall replace β , $\ln X$ and $\ln T$ by $-\beta$, X and ϵ , respectively. That is, the ordinary linear regression set-up:

$$X = \beta z + \epsilon, \text{ where } E(\epsilon) = \alpha \text{ and } Var(\epsilon) = \sigma^2.$$

Here $\beta \mathbf{z}$ can be interpreted as $\beta' \mathbf{z}$. For simplicity, we only consider the univariate case.

Question: $(\alpha, \beta) = ?$

If one has complete data, (X_i, z_i) , $i = 1, \dots, n$, the least squares estimate of β is the one that minimizes

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

The solution is

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = (U^t U)^{-1} U^t \mathbf{X},$$

where $U = \begin{pmatrix} 1 & z_1 \\ \vdots & \vdots \\ 1 & z_n \end{pmatrix}$ and $\mathbf{X}^t = (X_1, \dots, X_n)$.

The least squares estimate can be viewed as the solution of

$$(\alpha, \beta) = \operatorname{argmin} r(a, b), \text{ where } r(\alpha, \beta) = n \int t^2 d\hat{F}_{\alpha, \beta}(t),$$

where $\hat{F}_{\alpha, \beta}$ is the e.d.f. based on observations $T_i(\alpha, \beta) = X_i - \alpha - \beta z_i$, $i = 1, \dots, n$.

Moreover, the LSE is also the MLE under the normal assumption.

Thus it is the solution to the equation $\frac{\partial \ln \mathbf{L}}{\partial \beta} = 0$ and $\frac{\partial \ln \mathbf{L}}{\partial \alpha} = 0$.

§5.4.1. **Miller Estimator** (1976, Biometrika). The estimator of (α, β) with RC data is $(\hat{\alpha}, \hat{\beta}) = \operatorname{argmin}_{a, \mathbf{b}} r(a, \mathbf{b})$ (that minimizes $r(a, \mathbf{b})$ over all (a, \mathbf{b})), where

$$r(a, \mathbf{b}) = n \int t^2 d\hat{F}_{a, \mathbf{b}}(t), \quad \hat{F} = \hat{F}_{a, \mathbf{b}} \text{ is the GMLE of } F_\epsilon \text{ based on } (T_i(a, \mathbf{b}), \delta_i)\text{'s,}$$

and $T_i(a, \mathbf{b}) = M_i - a - \mathbf{b}'\mathbf{z}_i$. There are three issues:

- (1) How to justify the estimator ?
- (2) How to derive the estimate ?
- (3) Is it consistent or efficient ?

Actually, if $F_X(\tau_Y) < 1$, then Miller's estimator is not consistent, as $\alpha = E(\epsilon)$ is not identifiable.

Let s_i, \dots, s_m be the weight assigned by \hat{F} to the II's. If the largest observation is right censored, he suggested to pretend that it is exact to avoid assigning weight to $+\infty$ and thus $r(\alpha, \beta) = +\infty$. Note that s_i 's are only functions of β not of α , as changing α only shifts the II's and the corresponding intervals, but not the weights. If we let \hat{F}_0 be the PLE based on $(M_i - \beta z_i, \delta_i)$'s and \hat{F}_α the one based on $(M_i - \alpha - \beta z_i, \delta_i)$, then $\hat{F}_\alpha(t) = \hat{F}_0(t + \alpha)$. Let η_1, \dots, η_m be all the distinct exact observations based on $M_i - \beta z_i$'s.

$$\begin{pmatrix} \alpha : & 0 & c \\ II's : & \eta_j & \eta_j - c \\ weights : & s_j & s_j \end{pmatrix}$$

So we write $s_j = s_j(\beta)$ and $\eta_j = \eta_j(\beta)$. Denote $w_i(\beta)$ the weight assigned by \hat{F} to each observation $(M_i - \beta z_i, \delta_i)$ (treating each observation as one unit, even if there are ties).

Note that we may have ties at η_j , say, there are h exact observations such that $X_{i_k} - \beta z_{i_k} = \eta_j$ for $k = 1, \dots, h$. Then $w_{i_k} = s_j/h$ for each k . Then

$$r(\alpha, \beta) = n \sum_{i=1}^n (M_i - \alpha - \beta z_i)^2 w_i(\beta).$$

Taking derivative of $r(\alpha, \beta)$ w.r.t. α and setting the derivative to be zero yield

$$\alpha = \hat{\alpha}(\beta) = \sum_{i=1}^n w_i(\beta) (M_i - \beta z_i). \quad (1.0)$$

Thus it suffices to search $\hat{\beta}$ that minimizes

$$H(\beta) = \sum_{j=1}^n w_j(\beta) [M_j - \hat{\alpha}(\beta) - \beta z_j]^2. \quad (1.1)$$

Miller suggested the following iterative procedure:

1. Assign an initial value $\beta = \frac{\sum_{i: \epsilon x} X_i (z_i - \bar{z}_u)}{\sum_{j: \epsilon x} (z_j - \bar{z}_u)^2}$, where \bar{z}_u (\bar{X}_u) is the average of z_i 's (X_i 's) corresponding to exact observations of X_i .
2. Obtain the II's, $\eta_j(\beta)$ and the GMLE of $s_j(\beta)$'s based on $(M_i - \beta z_i, \delta_i)$ with the given β , and compute $w_i(\beta)$.
3. Update β by $\frac{\sum_{i=1}^n w_i(\beta) X_i (z_i - \bar{z}_w)}{\sum_{j=1}^n w_j(\beta) (z_j - \bar{z}_w)^2}$, where $\bar{z}_w = \sum_{j=1}^n w_j z_j$.
4. Repeat steps 2 and 3 until β converges or oscillates between two values. In the latter case, take the midpoint as an estimate of β .

The variance of β can be estimated by

$$\hat{\sigma}_{\hat{\beta}}^2 = \frac{\sum_{j=1}^n w_j(\hat{\beta}) (M_j - \hat{\alpha} - \hat{\beta} z_j)^2}{n \sum_{j=1}^n w_j(\hat{\beta}) (z_j - \bar{z}_w)^2}, \text{ comparing to } \sigma^2 / \sum_{i=1}^n (z_i - \bar{z})^2.$$

The consistency and asymptotic normality were considered under the assumption that

$$P\{X \text{ is not censored} | X = t\} > 0 \text{ for all possible } t \quad (1.2)$$

and the censoring distribution is of form

$$G(y|z) = G_o(y - \beta z), \text{ where } G_o \text{ is a cdf.}$$

However, the estimator has not been proved to be asymptotically efficient even under the normal assumption.

Remark For IC data, if we replace the PLE by GMLE, the above procedure can be adopted, provided we define that the GMLE only has jumps at midpoints of the II's. However, the consistency and asymptotic normality have not been verified.

§5.4.1.2. Homework. There are 4 observations $(M_i, \delta_i, \mathbf{z}_i)$'s: $(3, 1, 2)$, $(4, 0, 1)$, $(1, 1, 1)$, $(7, 1, 2)$. Find the Miller estimator of β under the linear regression model. $X_i = \beta z_i + \epsilon_i$. You are able to find the solution explicitly, because there are at most 6 distinct values of $\hat{S}_b(T_i(b))$ (as a function of b) for each fixed i , where $T_i = M_i - bz_i$.

§5.4.2. Buckley-James Estimator with RC data (1979, Biometrika).

Consider $X_i = \alpha + \beta'z_i + \epsilon_i$, $i = 1, \dots, n$, $\beta \in \mathcal{R}^p$, with complete data. Recall that the LSE is the solution to

$$\operatorname{argmin}_{a, \mathbf{b}} \sum_{i=1}^n (X_i - a - \mathbf{b}'z_i)^2 \quad (= \operatorname{argmin}_{a, \mathbf{b}} n \int t^2 d\hat{F}_{a, \mathbf{b}}(t)),$$

where $\hat{F} = \hat{F}_{a, \mathbf{b}}$ is the GMLE based on $X_i - a - \mathbf{b}'z_i$'s. Or the root of

$$\begin{aligned} H(b) &= \sum_{i=1}^n (X_i - \bar{X} - \mathbf{b}'(\mathbf{z}_i - \bar{\mathbf{z}}))(\mathbf{z}_i - \bar{\mathbf{z}}) = 0 \text{ with } a = \bar{X} - \mathbf{b}'\bar{\mathbf{z}} \\ H(b) &= \sum_{i=1}^n (X_i - \mathbf{b}'z_i)(z_i - \bar{\mathbf{z}}) \quad \left(\sum_{i=1}^n (\bar{X} - \mathbf{b}'\bar{\mathbf{z}})(z_i - \bar{\mathbf{z}}) = 0 \right). \end{aligned}$$

Given RC data, $(M_i, \delta_i, \mathbf{z}_i)$ s, denote $T_i = T_i(\mathbf{b}) = M_i - \mathbf{b}'z_i$, where $\mathbf{b} \in \mathcal{R}^p$.

Denote $\hat{S}_{\mathbf{b}}$ the PLE of the survival function based on (T_i, δ_i) , $i = 1, \dots, n$.

It is worth mentioning that $\hat{S}_{\mathbf{b}}$ depends on \mathbf{b} as T_i 's depend on \mathbf{b} .

Denote $\hat{f}_{\mathbf{b}}$ and $\hat{F}_{\mathbf{b}}$ the PLE's of the density function and the cdf, respectively.

Let $\hat{\mathbf{X}}^* = \hat{\mathbf{X}}^*(\mathbf{b}) = (\hat{X}_1^*, \dots, \hat{X}_n^*)'$, where

$$\hat{X}_i^* = M_i \delta_i + (1 - \delta_i) \left[\mathbf{b}'z_i + \frac{\sum_{t \in \mathcal{A}_i} t \hat{f}_{\mathbf{b}}(t)}{\hat{S}_{\mathbf{b}}(T_i)} \right], \quad (2.1)$$

$$\mathcal{A}_i = \{t : t > T_i, \hat{f}_{\mathbf{b}}(t) > 0\}.$$

Notice. If the largest observation $T_{(n)}$ among T_i 's is right censored, then Buckley and James suggest to treat $T_{(n)}$ as uncensored and

define the Buckley and James estimator (BJE) $(\hat{\alpha}, \hat{\beta})$ by

$$\hat{\alpha} = \overline{\hat{X}^*} - \hat{\beta}'\bar{\mathbf{z}} \text{ and}$$

$$b = \hat{\beta} \text{ being a solution to } H(b) = 0,$$

where $\bar{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i$ and

$$H(\mathbf{b}) = \sum_{i=1}^n (\hat{X}_i^*(\mathbf{b}) - \mathbf{b}'z_i)(z_i - \bar{\mathbf{z}}).$$

Buckley and James point out that the solution may not exist. In the latter case, if $p = 1$, the BJE of β is modified as the zero-crossing point of $H(b)$ (see Lai and Ying (1991)). One says that \hat{b} is a zero-crossing point of a function H if $H(b)$ changes its sign at $b = \hat{b}$.

Remark. With Buckley and James' modification, (2.1) defines X_i^* and is used in this note. One can also define $\hat{S}_{\mathbf{b}}(T_{(n)} + 1) = 0$ if $\delta_{(n)} = 0$. Then no need to modify $\delta_{(n)}$ if $\delta_{(n)} = 0$.

Two ways of justifying the BJE:

(1) Conditional expectation.

Note that

$$E(X|Z) = \alpha + \beta Z.$$

We only observe $M = X \wedge Y$ and

$$E(M|Z) \neq \alpha + \beta Z.$$

Define

$$X^* = X\delta + E(X|X > Y)(1 - \delta) \text{ (not the same as (2.1))} \quad (2.2)$$

where $E(X|X > Y)$ is defined as $g(\mathbf{1}(X > Y))$, and $g(y) = E(X|\mathbf{1}(X > Y) = y)$. Then

$$E(X^*|Z) = \alpha + \beta Z.$$

Reason: abusing notation, write $E_z(W(X, Y)) = E(W(X, Y)|Z)$, where $W(\cdot, \cdot)$ is a function, then

$$\begin{aligned} & E_z(X^*) \\ &= E_z(E_z(X^*|\delta)) \\ &= P(\delta = 1)E_z(X^*|\delta = 1) + P(\delta = 0)E_z(X^*|\delta = 0) \\ &= P(\delta = 1)E_z(X|\delta = 1) + P(\delta = 0)E_z(E_z(X|X > Y)|\delta = 0) \\ &= P(\delta = 1)E_z(X|\delta = 1) + P(\delta = 0)E_z(E_z(X|\delta = 0)|\delta = 0) \\ &= P(\delta = 1)E_z(X|\delta = 1) + P(\delta = 0)E_z(X|\delta = 0) \\ &= E_z\{E_z(X|\delta)\} \\ &= E_z(X) \\ &= \alpha + \beta Z. \end{aligned} \quad (2.3)$$

If we could observe X_1^*, \dots, X_n^* , we can use

$$\hat{\alpha} = \overline{X^*} - \hat{\beta}\bar{z} \text{ and } \hat{\beta} = \frac{\sum_{i=1}^n X_i^*(z_i - \bar{z})}{\sum_{j=1}^n (z_j - \bar{z})^2}. \quad (2.4)$$

Since we cannot observe all X_i^* 's, in view of (2.3) we replace X_i^* 's by their predictors

$$\hat{X}_i^* = X_i\delta_i + \hat{E}(X_i|X_i > Y_i)(1 - \delta_i)$$

(which is (2.1)), where

$$\hat{E}(X_i|X_i > Y_i) = \hat{\beta}z_i + \frac{\sum_{t > M_i - \hat{\beta}z_i} t \hat{f}_{\hat{\beta}}(t)}{\hat{S}(M_i - \hat{\beta}z_i)}, \quad (2.5)$$

and \hat{S} is the PLE of S based on observations $(M_i - \hat{\beta}z_i, \delta_i)$'s.

(2) An M-estimator based on $N(\mu, \sigma)$.

An M-estimator is the solution to $\frac{\partial \ln \mathbf{L}}{\partial \beta} = 0$. The likelihood function for given RC data is

$$\mathbf{L} = \prod_{i=1}^n (f(T_i))^{\delta_i} (S(T_i))^{1-\delta_i}.$$

$$\frac{\partial \ln L}{\partial \beta} = \sum_{i=1}^n \left\{ \delta_i \frac{f'}{f}(T_i) + (1 - \delta_i) \frac{-f}{S}(T_i) \right\} (-\mathbf{z}_i).$$

Under $N(\mu, \sigma^2)$ assumption, $f(t) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(t-\mu)^2}{2\sigma^2})$,

$$\frac{f'}{f}(t) = -\frac{t-\mu}{\sigma^2},$$

$f(t) = -\int_t^\infty f'(x)dx = -\int_t^\infty \frac{f'}{f}(x)f(x)dx = \int_t^\infty \frac{x-\mu}{\sigma^2} f(x)dx$. Then

$$\frac{\partial \ln L}{\partial \beta} = \sum_{i=1}^n \left\{ \delta_i (T_i - \mu) + (1 - \delta_i) \frac{\int_{T_i}^\infty (x - \mu) f(x) dx}{S(T_i)} \right\} (\mathbf{z}_i / \sigma^2). \quad (2.6)$$

The normal equation $H(\mathbf{b})$ is obtained by

replacing f and S by their PLE $\hat{f}_{\mathbf{b}}$ and $\hat{S}_{\mathbf{b}}$,

replacing μ by \bar{T}^* and

multiplying σ^2 .

In other words, (2.1) is the same as

$$H(\mathbf{b}) = \sum_{i=1}^n (T_i^* - \bar{T}^*) \mathbf{z}_i, \text{ where } T_i^* = \hat{X}_i^* - \mathbf{b} \mathbf{z}_i.$$

Notice

$$\begin{aligned} H(b) &= \sum_{i=1}^n (X_i - \bar{X} - \mathbf{b}'(\mathbf{z}_i - \bar{\mathbf{z}}))(\mathbf{z}_i - \bar{\mathbf{z}}) \\ &= \sum_{i=1}^n (X_i - \mathbf{b}'\mathbf{z}_i)(\mathbf{z}_i - \bar{\mathbf{z}}) \\ &= \sum_{i=1}^n (X_i - \bar{X} - \mathbf{b}'(\mathbf{z}_i - \bar{\mathbf{z}}))\mathbf{z}_i \quad \left(\sum_{i=1}^n (X_i - \bar{X} - \mathbf{b}'(\mathbf{z}_i - \bar{\mathbf{z}}))\bar{\mathbf{z}} = 0 \right) \end{aligned}$$

The consistency and the asymptotic properties of the BJE have been established

under continuous assumptions by Lai and Ying (1991) and

under discrete assumptions by Kong and Yu (2006).

In particular, if $\epsilon \sim N(\mu, \sigma^2)$, then the BJE is asymptotically efficient, just like the LSE.

Otherwise, it is not efficient.

Moreover, under certain discontinuous assumptions,

the BJE may not have asymptotic normality (see Kong and Yu (2006)).

Estimation of covariance matrix of the BJE, say $\Sigma_{\hat{\beta}}$.

With complete data the BJE becomes the LSE. Under the assumption that $\epsilon_1, \dots, \epsilon_n$ are i.i.d. with variance σ^2 , given U (or \mathbf{z}_i s), the covariance matrix of the LSE is

$$\Sigma_{\hat{\beta}} = \sigma^2 (U'U)^{-1}. \quad (2.7)$$

Reason:

$$\begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix} = \begin{pmatrix} \mathbf{z}'_1 \\ \vdots \\ \mathbf{z}'_n \end{pmatrix} \beta + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} \text{ or } \mathbf{X} = U\beta + \underline{\epsilon}.$$

$$\hat{\beta} = (U'U)^{-1}U'\mathbf{X} = (U'U)^{-1}U'(U\beta + \underline{\epsilon}) = (U'U)^{-1}U'U\beta + (U'U)^{-1}U'\underline{\epsilon}$$

$$Cov(\hat{\beta}) = (U'U)^{-1}U'Cov(\underline{\epsilon})(U'U)^{-1}U' = \sigma^2 (U'U)^{-1}$$

Formula (2.7) is valid no matter whether $X|\mathbf{z}$ is continuous or discrete.

Lai and Ying (1991, p.1389) present an estimator of $\Sigma_{\hat{\beta}}$ under smoothness assumptions with RC data. Kong and Yu present (2006) another estimator of $\Sigma_{\hat{\beta}}$ under discrete assumptions with RC data. Both expressions are complicated and are not given here.

An estimator under the normal assumption is their empirical Fisher information matrix $(I(\hat{\beta}))^{-1}$, where

$$I(\hat{\beta}) = \sum_{i=1}^n (T_i^* - \bar{T}^*)(\mathbf{z}_i - \bar{\mathbf{z}}) \{(T_i^* - \bar{T}^*)(\mathbf{z}_i - \bar{\mathbf{z}})\}' / \hat{\sigma}^4$$

and $\hat{\sigma}^2$ is an estimator of σ^2 .

For each i , denote

$$m = \sum_i \delta_i.$$

The parameter σ can be estimated in two ways.

If the largest M_i is not censored, then one can estimate it by

$$\hat{\sigma}^2 = \sum_i T_i^2 \hat{f}_{\hat{\beta}}(T_i) - (\sum_i T_i \hat{f}_{\hat{\beta}}(T_i))^2.$$

Otherwise, we can use the least squares method as follows.

We can find the quantiles of $\hat{F}(T_i)$ under $N(0, 1)$, say q_i 's.

Then we find the least squares estimate of (μ, σ) that minimizes $\sum_{i=1}^n \delta_i (\frac{T_i - \mu}{\sigma} - q_i)^2$.

It can be shown that the LSE is

$$\hat{\sigma} = \frac{T\bar{q} - \bar{T}\bar{q}}{q^2 - (\bar{q})^2},$$

where $\bar{T} = \frac{1}{m} \sum_{i=1}^n \delta_i T_i$,

$$\bar{q} = \frac{1}{m} \sum_{i=1}^n \delta_i q_i,$$

$$\bar{q}^2 = \frac{1}{m} \sum_{i=1}^n \delta_i q_i^2 \text{ and}$$

$$\bar{T}q = \frac{1}{m} \sum_{i=1}^n \delta_i T_i q_i.$$

An estimator of the variance of $\hat{\beta}$ given by Buckley and James is

$$\hat{\sigma}_{\hat{\beta}}^2 = \frac{\hat{\sigma}_u^2}{\sum_{i: ex} (z_i - \bar{z}_u)^2},$$

where $\hat{\sigma}_u^2 = \frac{1}{\sum_{i=1}^n \delta_i - 2} \sum_{i: ex} [X_i - \bar{X}_u - \hat{\beta}(z_i - \bar{z}_u)]^2$,

$$\bar{X}_u = \frac{\sum_{i: ex} X_i}{\sum_{i=1}^n \delta_i} \text{ and } \bar{z}_u = \frac{\sum_{i: ex} z_i}{\sum_{i=1}^n \delta_i}.$$

An alternative estimator of the variance of $\hat{\beta}$ is

$$\tilde{\sigma}_{\hat{\beta}}^2 = \frac{1}{\sum_{i=1}^n \delta_i - 2} \Sigma_z^{-1} \tilde{\sigma}_\epsilon^2,$$

where $\tilde{\sigma}_\epsilon = \int (t - \mu_{\epsilon, u})^2 d\hat{F}_{\hat{\beta}}(t)$,

$$\mu_{\epsilon, u} = \int t d\hat{F}_{\hat{\beta}}(t),$$

$$\tilde{\Sigma}_z^2 = \int_{t < \infty} (z(t) - \mu_{z, u})(z(t) - \mu_{z, u})' d\hat{F}_{\hat{\beta}}(t),$$

$$\mu_{z, u} = \int z(t) d\hat{F}_{\hat{\beta}}(t),$$

$\hat{F}_{\hat{\beta}}$ is the modified PLE that moves the weight from $+\infty$ to the largest observation, and

$$z(t) = \text{average of } z \text{ in } \{z_i : T_i(\hat{\beta}) = t, i = 1, \dots, n\}.$$

The last two estimators try to mimic the expression $\sigma^2(U'U)^{-1}$, but they are not consistent estimators.

The extension of Buckley-James estimator to the IC data are considered by Li and Pu (1999) and Rabinowitz, Tsiatis, and Aragon (1995).

An iteration algorithm for the BJE (Buckley and James (1979)).

1. Give initial values to β .
2. Obtain \hat{X}_i^* 's using (2.1) with the given β .
3. Update β using (2.4) with the given \hat{X}_i^* 's: $\hat{\alpha} = \overline{\hat{X}^*} - \hat{\beta}\bar{z}$ and $\hat{\beta} = \frac{\sum_{i=1}^n X_i^*(z_i - \bar{z})}{\sum_{j=1}^n (z_j - \bar{z})^2}$.
4. Repeat steps 2 and 3 until β converges or oscillates between two values. In the latter case, take the midpoint as an estimate of β .

Remark.

1. The algorithm may not converges to a solution of the BJE even if the zero point of H exits (see Yu and Wong (2002)).
2. The BJE of β may not be unique. If there are both root and non-root zero-crossing point to $H(\mathbf{b})$, the iterative algorithm may present non-root zero-crossing-point of $H(\mathbf{b})$.

R codes:

```
library(rms)
library(MASS)
library(splines)
library(survival)
bj(Surv(m, d) ~ x, link="identity",control=list(iter.max=50))
```

A non-iterative algorithm for obtaining all BJE's (for $p = 1$) (Yu and Wong (2002a)):

1. Let b_{ij} be the solution to an equation $T_i(b) = T_j(b)$, where $z_i \neq z_j$ and $\delta_i \neq \delta_j$. Let $q_1 < \dots < q_m$ be all the distinct solutions b_{ij} 's. Let $q_0 = -\infty$ and $q_{m+1} = \infty$.
2. (Case (1)). For each $h = 0, 1, \dots, m$, first compute the PLE \hat{S}_b for a $b \in (q_h, q_{h+1})$. For example, let b be the midpoint of the interval if $0 < i < m$, $b = q_1 - 1$ if $i = 0$, and $b = q_m + 1$ if $i = m$. Then compute $(M_i^*(b), z_i^*(b))$'s and

$$\hat{b}_h = \frac{\sum_{j=1}^n (z_j - \bar{z}) M_j^*}{\sum_{k=1}^n (z_k - \bar{z}) z_k^*}, \quad \text{where}$$

$$M_i^* = M_i \delta_i + (1 - \delta_i) \frac{\sum_{t > T_i(\mathbf{b})} \hat{f}_{\mathbf{b}}(t) \frac{\sum_{j=1}^n M_j \mathbf{1}_{(T_j(\mathbf{b})=t, \delta_j=1)}}{\sum_{k=1}^n \mathbf{1}_{(T_k(\mathbf{b})=t, \delta_k=1)}}}{\hat{S}_{\mathbf{b}}(T_i(\mathbf{b}))}$$

$$z_i^* = \delta_i z_i + (1 - \delta_i) \left[\frac{\sum_{t > T_i(\mathbf{b})} \hat{f}_{\mathbf{b}}(t) \frac{\sum_{j=1}^n z_j \mathbf{1}_{(T_j(\mathbf{b})=t, \delta_j=1)}}{\sum_{k=1}^n \mathbf{1}_{(T_k(\mathbf{b})=t, \delta_k=1)}}}{\hat{S}_{\mathbf{b}}(T_i(\mathbf{b}))} \right]$$

If $\hat{b}_h \in (q_h, q_{h+1})$ then \hat{b}_h is a root of $H(b)$ (see (2.8)) and thus is a BJE of β .

3. (Case (2)). Compute $H(q_i-)$, $H(q_i)$ and $H(q_i+)$, $i = 1, \dots, m$. If $H(q_i-)H(q_i+) \leq 0$, or $H(q_i-)H(q_i) \leq 0$, or $H(q_i)H(q_i+) \leq 0$, then q_i is a zero-crossing point of H and thus is a BJE of β .

Remark. In computing $H(b)$ is better off to use the following equivalent expression

$$H(b) = \sum_{i=1}^n (M_i^* - \mathbf{b}' \mathbf{z}_i^*) (\mathbf{z}_i - \bar{\mathbf{z}}), \quad (2.8)$$

rather than

$$H(\mathbf{b}) = \sum_{i=1}^n (\hat{X}_i^*(\mathbf{b}) - \mathbf{b}' \mathbf{z}_i) (\mathbf{z}_i - \bar{\mathbf{z}}),$$

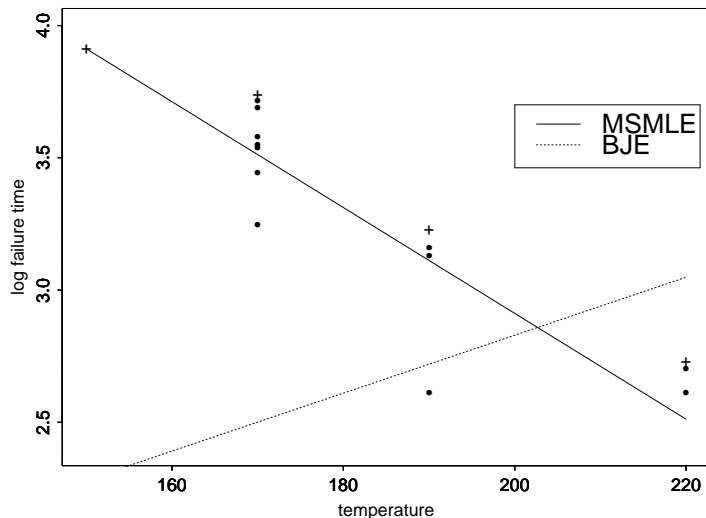
because even though (M_i^*, z_i^*) depends on b , it is constant in b in each interval (q_j, q_{j+1}) .

Remark. If data are complete, both the BJE $\hat{\beta}_{BJ}$ and the Miller estimator $\hat{\beta}_M$ reduce to the LSE. In fact, $\hat{\beta}_M = \operatorname{argmin}_{\mathbf{b}} \int t^2 d\hat{F}_{\mathbf{b}}(t)$, where $\hat{F}_{\mathbf{b}}$ is the PLE of F_o based on $X_i - \bar{X} - \mathbf{b}'(z_i - \bar{z})$'s is an extension of $\operatorname{argmin}_{\mathbf{a}, \mathbf{b}} \sum_{i=1}^n (X_i - a - \mathbf{b}'z_i)^2/n$, and $\mathbf{b} = \hat{\beta}_{BJ}$ is the zero-crossing point of $H(\mathbf{b})$ and $H(b)$ is the extension of $\frac{d}{d\mathbf{b}} \sum_{i=1}^n (X_i - a - \mathbf{b}'z_i)^2$.

Example 2.1. (Insulation data (Nelson 1973)). To evaluate a new Class-B insulation for electric motors, temperature-accelerated life testing was conducted on 40 motorettes. The main purpose was to estimate the distribution of insulation at the design temperature of 130°C. Ten motorettes were put on test at each of four temperatures (150°C, 170°C, 190°C, and 220°C). Let z be the temperature (in °C) and X (or M) the logarithm of hours to failure of an insulation at temperature x . The data are plotted in Figure 1, “+” stands for right-censored observations and “.” stands for exact observations.

If one sets $\hat{f}_b(M_{(n)}(b)) = 0$ when $\delta_{(n)} = 0$, rather than $\hat{f}_b(M_{(n)}(\mathbf{b})) = \hat{S}_b(M_{(n)}(\mathbf{b}))$, no matter what initial point is used, the existing algorithms always result in an estimate $\hat{\beta}_1 = 0.0109$, which is the unique solution to equation (2.3). The fitted line is plotted in Figure 1 in broken line. $\hat{\beta}_1$ does not make sense, as it should be negative according to the data (see Figure 1.). Yu and Wong (2002a) present an algorithm that can find all possible solutions for the BJE. Using this algorithm, we found that there are exactly 3 zero-crossing points: -0.0207 , -0.0205 and -0.0193 . They are approximately -0.02 . We plot the BJE fitted line corresponding to $\hat{\beta} = -0.02$ in Figure 1 (in solid line). It appears to be a reasonable estimate and it is actually the semi-parametric MLE (SMLE).

Fig. 5.1. MSMLE vs. BJE For Insulation Data



If one sets $\hat{f}_b(M_{(n)}(\mathbf{b})) = \hat{S}_b(M_{(n)}(\mathbf{b}))$ (when $\delta_{(n)} = 0$), then there is just unique BJE and it is a root of H .

Table 4.1 presents two data sets, of sample size 30 each, generated from simulation. Both set $\hat{f}_b(M_{(n)}(\mathbf{b})) = \hat{S}_b(M_{(n)}(\mathbf{b}))$ if $\delta_{(n)} = 0$.

The first data set has no zero-crossing point of the sum of least squares

and has exactly two solutions to equation (2.3), which are $(a, b) = (0.873114, 0.939279)$ and $(a, b) = (0.856106, 0.944563)$, respectively.

The second one does not have a solution to Eq. (2.3),

but has a unique zero-crossing point at $b = -0.199549$.

The current iterative algorithm oscillates between -0.642338 and 0.805032 .

The midpoint approach results in the midpoint BJE of β as 0.081347 and the BJE of β in the data set is -0.199549 .

They are not the same even at the first non-zero digit.

Thus the current midpoint approach is not close in approximation and our method gives the precise value.

§5.4.2.2. Homework

1. There are 4 observations $(M_i, \delta_i, \mathbf{z}_i)$'s: $(3, 1, 2)$, $(4, 0, 1)$, $(1, 1, 1)$, $(7, 1, 2)$. Show that there is only one BJE of β under the linear regression model and it is 2.
2. Under interval censoring, (2.2) can be rewritten as

$$X_i^* = E(X_i | X_i \in I_i), \text{ where } I_i \text{ is the } i\text{-th observed interval.}$$

2.a. Verify (2.3) under the mixed case IC model with continuous random vectors.

2.b. Give an estimator of X_i^* corresponding expressions for (2.5) and give a corresponding expressions for (2.6).

Data1 :			Data2 :		
5.182270	0	8.533530	14.822689	0	7.212610
13.111501	0	2.663143	6.457941	0	6.103648
8.424795	0	1.969974	8.239654	1	7.739654
10.160930	1	9.660930	0.805357	0	6.010412
8.389024	0	8.689664	11.343863	0	1.808932
0.680877	1	0.180877	6.929858	1	6.429858
12.061478	0	7.999715	7.445796	0	6.180522
9.582826	1	9.082826	15.379855	0	1.814571
6.589164	0	9.442880	0.419907	0	0.091820
8.013080	0	1.630284	17.549400	1	2.049400
0.811283	1	0.311283	16.841743	0	5.595673
14.260701	0	1.824968	7.904303	0	4.497519
16.177659	0	3.087899	2.657995	1	2.157995
6.984444	0	9.305123	3.772100	0	5.745633
6.078930	0	5.328123	7.480979	0	6.541589
17.478192	0	7.680852	4.419456	1	3.919456
1.635554	0	8.539667	1.416040	1	0.916040
0.355867	0	4.190535	9.712842	1	9.212842
4.687164	1	4.187164	15.500627	0	0.147172
3.497800	1	2.997800	15.202233	0	1.296472
0.871094	1	0.371094	0.235298	0	3.911587
5.445252	1	4.945252	0.271033	0	9.242885
13.928942	0	0.590025	11.917236	0	6.549110
3.919157	1	3.419157	0.532521	1	0.032521
6.351499	1	5.851499	14.215613	0	6.540187
4.531607	1	4.031607	14.370996	0	7.267528
3.440897	0	8.922772	1.883544	1	1.383544
1.657479	1	1.157479	1.508759	1	1.008759
7.315708	1	6.815708	15.451111	0	0.100477
11.559872	0	3.293425	15.694594	1	0.194594

Table 4.1. Simulation Examples

§5.4.3. An M-estimation approach

Recall that Huber (1964) proposed an M-estimator

which is a zero point of a score function $\sum_{i=1}^n \psi(\theta, T_i)$ ($= \frac{\partial}{\partial \theta} \ln L$)

where θ is the parameter of interest and T_i 's are observations.

Modifying Huber's M-estimation (Huber, 1964),

Zhang and Li (1996) consider another M-estimation approach with interval-censored data.

The idea is to find a zero point of an *estimate* of the score function $\frac{\partial}{\partial \mathbf{b}} \ln L$ in \mathbf{b} .

We shall first illustrate via RC data.

Note that the likelihood function can be written as

$$L(b, S, f) = \prod_{i=1}^n (f(M_i - b'z_i))^{\delta_i} (S(M_i - b'z_i))^{1-\delta_i}.$$

Assuming that f and S are differentiable, the "MLE" of β is a critical point of L , where

a critical point is a point that either L is not differentiable or $\frac{\partial L}{\partial \mathbf{b}} = 0$.

Moreover, to eliminate the effect of α , one needs to centralize z_i in L .

Thus the derivative of $-\ln L$ is

$$\Phi = \sum_{i=1}^n (z_i - \bar{z}) \left(\delta_i \left(\frac{f'}{f} \right) (M_i - bz_i) - (1 - \delta_i) \left(\frac{f}{S} \right) (M_i - bz_i) \right). \quad (3.1)$$

Since (under certain assumptions)

$$\int_{x>t} \frac{f'(x)}{f(x)} f(x) dx = \int_{x>t} df(x) = f(x) \Big|_t^{\infty} = -f(t),$$

as $f(\infty) = 0$, we have

$$\frac{f}{S}(t) = \frac{-\int_{x>t} \frac{f'(x)}{f(x)} f(x) dx}{S(t)} = \frac{-\int_{x>t} \frac{f'(x)}{f(x)} dF(x)}{S(t)} = \frac{\int_{x>t} \frac{f'(x)}{f(x)} dS(x)}{S(t)}.$$

$$\text{Thus } \Phi = \Phi(b, S, \frac{f'}{f}) = \sum_{i=1}^n (z_i - \bar{z}) \left(\delta_i \left(\frac{f'}{f} \right) (T_i(b)) - (1 - \delta_i) \frac{\int_{x>T_i(b)} \frac{f'(x)}{f(x)} dS(x)}{S(T_i(b))} \right),$$

where $T_i(b) = M_i - bz_i$. Note that S , f and f' are all unknown. Thus they need to be estimated. If one replaces S by its PLE and chooses $f = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$ (thus $(f'/f)(x) = -x$), then

$$\Phi = - \sum_{i=1}^n (z_i - \bar{z}) \left(\delta_i(T_i(b)) - (1 - \delta_i) \frac{\int_{x>T_i(b)} x d\hat{S}_b(x)}{\hat{S}_b(T_i(b))} \right)$$

and the M-estimator reduces to the BJE. Thus the BJE is an M-estimator.

Zhang and Li suggest to look for a root of an estimate of $\Phi(b, S, \frac{f'}{f})$, say

$$\Phi(b, \hat{S}_b, \frac{\hat{f}'_b}{\hat{f}_b})$$

where \hat{S}_b is the PLE of S_o based on $(M_i - bz_i, \delta_i)$'s, and \hat{f}_b is a kernel estimator

$$\hat{f}_b(t) = \frac{1}{h} \int K\left(\frac{x-t}{h}\right) d\hat{F}_b(x), \text{ with } K \geq 0, \int K(x) dx = 1,$$

the set $\{x : K'(x) \neq 0\}$ is not a null set, $n^{-1/2}/h \rightarrow 0$

e.g., $h = cn^{-1/5}$ and c is a predetermined constant. (3.2)

Examples of such kernels are

$$K(x) = \frac{3}{4}(1-x^2)\mathbf{1}_{(|x|\leq 1)},$$

and

$$K(x) = (1-|x|)\mathbf{1}_{(|x|\leq 1)}.$$

Other examples can be found in Härdle W. (1990, p.45). It is obvious that an M-estimate can only be obtained by iterative algorithms. Zhang and Li show that under certain regularity conditions, the M-estimator is consistent and is asymptotically efficient if the initial point in the iterative algorithm is close enough to β . The asymptotic variance of the M-estimator is expected to be $(E(\Phi(\beta)\Phi'(\beta)))^{-1}$ and can be estimated by $\hat{\Sigma} =$

$$\left(\sum_{i=1}^n \left(\delta_i \left[\left(\frac{\hat{f}'}{\hat{f}} \right) (M_i - bz_i) \right]^2 - (1 - \delta_i) \left(\frac{\hat{f}}{\hat{S}} \right)^2 (M_i - bz_i) \right) (z_i - \bar{z})(z_i - \bar{z})' \right)^{-1},$$

where \hat{S} , \hat{f} and \hat{f}' are estimates of S_o , f_o and f'_o , respectively.

However, in practice, there are some outstanding computational issues with this approach:

1. it is not clear how to select an initial point that is really "close" to β .
2. It is not clear how to select a constant c in (3.2).
3. It is not clear what is an optimal choice of the kernel K .
4. A similar phenomenon like the BJE, which is also a special case of M-estimators, may also occur. That is, there does not exist a zero point of Φ .

5. Even if one may obtain a zero point of Φ , it may not be the solution that is closed to the maximum point of the likelihood (*i.e.*, as point that is near a local maximum or even a local minimum of the likelihood). Consequently, the estimate is not good.

An algorithm for the case $p = 1$ maybe as follows, assuming $\Phi(-\infty)\Phi(\infty) < 0$.

1. Choose a tolerance number $\eta > 0$, *e.g.*, $\eta = 0.00001$.
2. Choose a $b_1 > 0$ and let $b_2 = -b_1$ such that $\Phi(b_1)\Phi(b_2) < 0$. If $\Phi(b_i) \in (-\eta, \eta)$ for an $i \in \{1, 2\}$, stop and let b_i be the “estimate” (**treat as a zero point of Φ**). Otherwise, go to next step.
3. Let $b_3 = (b_2 + b_1)/2$. If $\Phi(b_3) \in (-\eta, \eta)$, stop and let b_3 be the “estimate”. Otherwise, go to next step.
4. Set $b_{i+1} = \begin{cases} (b_{i-1} + b_i)/2 & \text{if } \Phi(b_{i-1})\Phi(b_i) < 0 \\ (b_{i-2} + b_i)/2 & \text{if } \Phi(b_{i-2})\Phi(b_i) < 0 \end{cases}$ for $i \geq 2$ iteratively until either $|b_i - b_{i+1}| < \eta$ (**zero point**), or $\Phi(b_{i+1}) \in (-\eta, \eta)$ (**zero-crossing**).

Remark. Verify that in the case of complete data,

$$\Phi = -2\left(\sum_{i=1}^n (z_i - \bar{z})(X_i - \bar{X}) - b \sum_{i=1}^n (z_i - \bar{z})^2\right).$$

Thus $\Phi(-\infty) > 0$ and $\Phi(\infty) < 0$.

One expects that the assumption $\Phi(-\infty)\Phi(\infty) < 0$ holds in general.

Another algorithm is to find a minimum point of $|\Phi(\mathbf{b})|$ by Monte Carlo method.

That is, randomly select a sequence of values of \mathbf{b} and find the up-to-date minimum point until it is stable.

§5.4.3.2. Homework.

1. Derive the expressions of Φ when f is the density of $U(0, \theta)$ and when $f(x) = e^{-x}$, $x > 0$.
2. Give the expression of Φ when data are mixed IC type.

§5.4.4. An SMLE approach

Recall we consider regression model $X_i = \beta' \mathbf{z}_i + \epsilon_i$, $i = 1, \dots, n$. Yu and Wong (2003a,b,c) proposed the SMLE of (β, S_o) , based on complete data, RC data and IC data, respectively. The semi-parametric likelihood function is

$$\mathbf{L} = \prod_{i=1}^n \mu_F(I_i - bz_i),$$

where F is a cdf, I_i 's are the observed intervals and $I_i - c$ is a shift of the interval I_i by c units. Thus the SMLE of (β, F_o) is the value of (b, F) that maximizes $\mathbf{L}(b, F)$.

Without loss of generality, we assume that the dimensions of b and \mathbf{z}_i are 1.

For fixed b , the likelihood function is maximized by the GMLE of F_o based on $I_i - bz_i$'s, denoted by \hat{F}_b .

Thus it suffices to maximizes $\mathbf{L}(b, \hat{F}_b)$.

Since \hat{F}_b , as a function of b , only takes finitely many values, so is \mathbf{L} .

Let $a_1 < \dots < a_m$ be all the solutions to equations of form

$$M_k - bz_k = M_j - bz_j, \quad z_k \neq z_j, \quad \text{where } M_i = L_i \text{ or } R_i, \quad a_0 = -\infty \text{ and } a_{m+1} = \infty.$$

Note that the

GMLE $\hat{F}_b(L_j - bz_j)$ and $\hat{F}_b(R_j - bz_j)$'s only depend on the ranks of $L_i - bz_i$'s and $R_i - bz_i$'s, and the ranks of $L_i - bz_i$'s and $R_i - bz_i$'s are constant for $b \in (a_k, a_{k+1})$.

Thus for each fixed i , $\mathbf{L}(b, \hat{F}_b)$ is constant in b on the interval (a_i, a_{i+1}) .

Example 1. Consider a simple example of RC data, say,

$$3 \text{ } (M_i, \delta_i, z_i) \text{ s are } (1, 1, 1), (2, 0, 1), (3, 1, 0).$$

Let $T_i(b) = M_i - bz_i$. $M_k - bz_k = M_j - bz_j$ and $z_k \neq z_j$ yield

$$1 - b = 3 \text{ } (T_1(b) = T_3(b)) \text{ and } 2 - b = 3 \text{ } (T_2(b) = T_3(b)).$$

Thus $a_1 = -2$ and $a_2 = -1$.

Let (r_1, r_2, r_3) be the ranks of $(T_1, T_2, T_3)(b)$ s or

$$(1 - b, (2 - b)^+, 3).$$

$I_i(b)$ s are $([T_1, T_1], (T_2, \infty), [T_3, T_3])$.

$b \in$	$(-\infty, -2)$	$\{-2\}$	$(-2, -1)$	$\{-1\}$	$(-1, \infty)$
$r_i s$	$(2, 3, 1)$	$(1.5, 3, 1.5)$	$(1, 3, 2)$	$(1, 3, 2)$	$(1, 2, 3)$
$\mu_{\hat{F}_b}(I_i(b))s$	$(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$	$(\frac{2}{3}, \frac{1}{3}, \frac{2}{3})$	$(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$	$(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$	$(\frac{1}{3}, \frac{2}{3}, \frac{2}{3})$
$maxL$	no	yes	no	no	yes

Each point in $\{-2\} \cup (-1, \infty)$ is an SMLE of β .

In the case of RC data with $p = 1$, there are finitely many such disjoint intervals and the PLE and L have explicit forms.

Thus, the SMLE can be obtained by a non-iterative algorithm.

Let \mathcal{A} be the union of all a_i 's, all midpoints $(a_i + a_{i+1})/2$, and $a_1 - 1$ and $a_m + 1$.

A point in \mathcal{A} that maximizes $L(b, \hat{F}_b)$ over all $b \in \mathcal{A}$ is an SMLE of β .

We summarize as an algorithm:

1. Derive \mathcal{B} and \mathcal{A} , where $\mathcal{B} = \{b : b = \frac{M_i - M_j}{z_i - z_j}, z_i \neq z_j, (\delta_i, \delta_j) \neq (0, 0), 1 \leq i < j \leq n\}$.
Let $a_1 < \dots < a_m$ be the distinct elements of \mathcal{B} .
2. Derive $\mathcal{L}(b) = L(b, \hat{F}_b)$ for each $b \in \mathcal{A}$.
3. The maximizer of $\mathcal{L}(b)$, $b \in \mathcal{A}$ is an SMLE of β , denoted by $\hat{\beta}$. Moreover, if $a, \hat{\beta} \in (a_i, a_{i+1})$, then a is also an SMLE of β . The SMLE of S_o is $\hat{S}_{\hat{\beta}}$.

In the case of IC data, the GMLE may not have a closed form expression, one has to use iterative algorithm to obtain the GMLE and thus the SMLE may be obtained by an iterative algorithm.

The SMLE may not be unique, thus one can choose one that is the closest to the median of all SMLE's. Note that the median may not be an SMLE.

Question: If we have complete data and $\mathbf{z}_i = 1 \forall i$, SMLE of $\beta = ?$

Each real number is an SMLE of β , as the edf gives equal weight to each $X_i - b$.

Properties of the SMLE with RC data:

1. If F_o is continuous and $P\{\delta = 1\} = 1$, there are inconsistent SMLE's and consistent SMLE's (see Yu and Wong (2003b)).
2. If F_o is discontinuous, then $P\{\hat{\beta}_n \neq \beta \text{ infinitely often}\} = 0$
or $P\{\hat{\beta}_n = \beta \text{ for large enough } n\} = 1$.
3. The SMLE and the BJE cannot dominate each other. In particular, if $\epsilon \sim N(\mu, \sigma^2)$, then the BJE is efficient, but not the SMLE.
It is conjectured that the SMLE $\hat{\beta}$ (or $\hat{\beta}_n$) has the following properties:
4. If F_o is continuous and $P\{\delta = 1\} \in (0, 1)$, then the SMLE is consistent.
5. If F_o is continuous and $P\{\delta = 1\} \in (0, 1)$, then the SMLE is asymptotically normal, with estimated asymptotic covariance matrix

$$\hat{\Sigma} = \left(\sum_{i=1}^n (1 - \delta_i) \left(\frac{\tilde{f}(T_i(\hat{\beta}))}{\hat{S}_{\hat{\beta}}(T_i(\hat{\beta}))} \right)^2 z_i z_i' \right)^{-1}, \quad (4.1)$$

where \tilde{f} is a kernel estimate of the df of F_o . Note that

$$\begin{aligned} \ln(f^\delta(T(b))S^{1-\delta}(T(b))) &= \delta \ln f(T(b)) + (1 - \delta) \ln S(T(b)) \\ \frac{\partial \ln(f^\delta(T(b))S^{1-\delta}(T(b)))}{\partial b} &= \left[-\delta \frac{f'(T(b))}{f_o(T(b))} + (1 - \delta) \frac{f_o'(T(b))}{S_o(T(b))} \right] Z \end{aligned}$$

Thus the conjecture says that the first term is missing in $\hat{\Sigma}$.

6. If F_o is continuous, then the SMLE is not efficient. In fact, if $\hat{\Sigma}$ in (4.1) is true, this is obvious as the efficient covariance matrix is

$$E\left[\left(\delta \frac{f'(T(b))}{f(T(b))} + (1 - \delta) \frac{f(T(b))}{S(T(b))} \right)^2 Z Z' \right]$$

We now verify property 1.

Assume that $p = 1$, ϵ and \mathbf{z} are both continuous independent random variables.

Let (X_i, z_i) , $i = 1, \dots, n$ are observations.

Then $\hat{\beta} = \frac{X_1 - X_2}{z_1 - z_2}$ **Why ??**

$X_i, z_i, \hat{\beta}$ are all continuous random variables.

$$\hat{\beta} = \beta + \frac{\epsilon_1 - \epsilon_2}{z_1 - z_2}.$$

Let $\hat{\beta}_n = \max\{b : b \in \mathcal{B}\}$, then $\hat{\beta}_n \rightarrow \infty$ a.s., as $\min\{|Z_1 - Z_2|\} \rightarrow 0$ a.s..

That is, $\hat{\beta}_n$ is an inconsistent SMLE.

On the other hand, the SMLE that is closest to the LSE is consistent.

We shall present some of the results in simulation studies. The main purpose is to study the properties of the SMLE when F_o is arbitrary, *i.e.*, continuous, or discontinuous but not necessarily discrete. We assume that Z, ϵ and Y are independent. We consider several cases in our simulation studies:

- (1) F_o is continuous (Examples 4.3, 4.4 and 4.5), or discrete (Example 4.1), or discontinuous but not discrete (Examples 4.2).
- (2) All the underlying distributions belong to the exponential family (Examples 4.3 and 4.5) or F_o does not belong to the exponential family (other examples).

In the following examples, let $X = \beta Z + \epsilon$ and $E(\epsilon) = \alpha$.

Example 4.1. Suppose ϵ equals 13.5 and 38.5, with probabilities 0.5 and 0.5, respectively, $Z \sim U(2, 3)$, $Y \sim U(24, 24.2)$, and $(\alpha, \beta) = (26, 1)$.

Example 4.2. Suppose ϵ is a mixture of $U(0, 0.5)$ and a unit point mass concentrated at 0.25, with probabilities 0.5 and 0.5, respectively, $Z \sim U(1, 2)$, $Y \sim U(0, 4)$, and $(\alpha, \beta) = (0.25, 1)$.

Table 1. Simulation Results on estimating (α, β)				
		(α, β)	SMLE $(\hat{\alpha}, \hat{\beta})$	BJE
Example 4.1 (discrete F_o).				
n=32	average	(13, 1)	(0.941, 1.000)	(-5.388, 1.749)
	SE		(0.044, 0.000)	(22.368, 6.702)
n=200	average	(13, 1)	(0.360, 1.000)	(1.388, 0.301)
	SE		(0.000, 0.000)	(1.075, 0.255)
Example 4.2 (discontinuous F_o).				
n=32	average	(0.25, 1)	(0.248, 1.001)	(0.258, 0.997)
	SE		(0.038, 0.021)	(0.142, 0.961)
n=200	average	(0.25, 1)	(0.250, 1.000)	(0.247, 1.002)
	SE		(0.010, 0.000)	(0.050, 0.033)
Example 4.3 (continuous F_o).				
n=32	average	(5, 1)	(3.693, 1.415)	(4.035, 0.884)
	SE		(3.784, 3.227)	(2.106, 1.708)
n=200	average	(5, 1)	(4.701, 0.957)	(4.701, 0.959)
	SE		(1.294, 0.693)	(0.798, 0.753)
Example 4.4 (continuous F_o).				
n=32	average	(0, 2)	(1.446, 2.129)	(0.008, 2.001)
	SE		(5.411, 0.464)	(0.291, 0.022)
n=200	average	(0, 2)	(0.821, 2.077)	(-0.002, 2.000)
	SE		(0.876, 0.080)	(0.093, 0.007)
Example 4.5 (continuous F_o).				
n=32	average	(0, 1)	(-0.2093, 1.1651)	(-0.0110, 1.0145)
	SE		(1.5933, 0.8971)	(0.7867, 0.3814)
n=200	average	(0, 1)	(-0.2598, 1.1516)	(0.0048, 0.9942)
	SE		(0.7068, 0.3831)	(0.2350, 0.1116)

Hereafter denote $Exp(\mu, \sigma)$ a distribution with the df

$$f(x) = \frac{1}{\sigma} e^{-[\frac{x-\mu}{\sigma}+1]} \mathbf{1}_{(x>\mu-\sigma)}.$$

Q: Does it belong to the exponential family ?

Example 4.3. Suppose ϵ , Y and \mathbf{Z} have distributions $Exp(5, 2)$, $Exp(3, 4)$ and $Exp(2, 2)$, respectively. $(\alpha, \beta) = (5, 1)$.

Example 4.4. Suppose $\epsilon \sim U(-1, 1)$, $Z \sim Exp(0, 19)$, $Y \sim Exp(0, 25)$, $(\alpha, \beta) = (0, 2)$.

Example 4.5. Suppose $\epsilon \sim N(0, 1)$, $Y \sim N(0, 6)$, $Z \sim N(2, 1)$, $(\alpha, \beta) = (0, 1)$.

The results of the above examples are summarized in Table 1. One can see that the SMLE is better than the BJE under the exponential distribution, but vice versa under the normal distribution.

Property 2 is proved for the following cases (1) complete-data and (2) right-censored discrete data (Yu and Wong (2003a,b)). The proof for right-censored discontinuous data is under preparing. The following example illustrates a proof under a simple assumption.

Example 4.6. Consider the simple linear regression ($p = 1$), with complete data. Suppose that ϵ and Z have the binomial distribution, $Bin(1, 0.5)$ and $\beta = 1$. The possible values of the observation (Z, X) are $(0, 0)$, $(0, 1)$, $(1, 1)$ and $(1, 2)$, denoted by (Z_i, X_i) , $i = 1, 2, 3, 4$. Thus there are 4 possible values of $T(b)$, say, $T_i(b) = X_i - bZ_i$, $i = 1, \dots, 4$. Suppose a random sample of size n contains N_1, N_2, N_3 and N_4 of them. The empirical df

$$\hat{f}_b(T_i(b)) = \begin{cases} \frac{N_i}{n} & \text{if } T_i(b) \neq T_j(b) \forall j \neq i, \\ \frac{N_i + N_j}{n} & \text{if } T_i(b) = T_j(b) \text{ for only one } j \neq i, \end{cases} \quad (4.2)$$

where $i, j = 1, 2, 3, 4$. The possible finite solutions b to the equations $T_i(b) = T_j(b)$ are 0, 1 and 2.

For n large enough, $N_i \approx n/4$, and likelihood function (1.2) or (4.1)

$$\mathcal{L} = \prod_{i=1}^n \hat{f}_b(T_i(b)) =$$

$$\begin{cases} \left(\frac{N_1 + N_3}{n} \right)^{N_1 + N_3} \left(\frac{N_2 + N_4}{n} \right)^{N_2 + N_4} \approx (0.5)^n & \text{if } b = 1, \\ \left(\frac{N_1}{n} \right)^{N_1} \left(\frac{N_2 + N_3}{n} \right)^{N_2 + N_3} \left(\frac{N_4}{n} \right)^{N_4} \approx (0.5)^{n/2} (0.25)^{n/2} & \text{if } b = 0, \\ \left(\frac{N_1 + N_4}{n} \right)^{N_1 + N_4} \left(\frac{N_2}{n} \right)^{N_2} \left(\frac{N_3}{n} \right)^{N_3} \approx (0.5)^{n/2} (0.25)^{n/2} & \text{if } b = 2, \\ \left(\frac{N_1}{n} \right)^{N_1} \left(\frac{N_2}{n} \right)^{N_2} \left(\frac{N_3}{n} \right)^{N_3} \left(\frac{N_4}{n} \right)^{N_4} \approx (0.25)^n & \text{otherwise,} \end{cases} \quad (4.3)$$

is maximized by $b = 1$. Thus the SMLE of β is $\hat{\beta} = 1 = \beta$ for all large n (i.e., (1.3) holds). The LSE $\tilde{\beta}_{LSE} = \beta + \frac{N_1 N_4 - N_2 N_3}{(N_1 + N_4)(N_2 + N_3)}$. Thus,

$$P(\tilde{\beta}_{LSE} \neq \beta \text{ i.o.}) = 1.$$

The LSE $\tilde{\beta}_{LSE}$ satisfies

$$\sqrt{n}(\tilde{\beta}_{LSE} - \beta) \xrightarrow{D} N(0, 1), \text{ and } n\sigma_{\tilde{\beta}_{LSE}}^2 \rightarrow 1 \text{ as } n \rightarrow \infty,$$

$$\text{as the asymptotic variance } \sigma_{\tilde{\beta}_{LSE}}^2 = \frac{1}{n}.$$

It can be shown that $\sqrt{n}(\tilde{\beta} - \beta) \xrightarrow{D} N(0, 0)$ and $n\sigma_{\tilde{\beta}}^2 \rightarrow 0$. \square

Example 4.7. (Magazine advertising (Chatterjee and Price ((22), p. 257)). In a study of revenue from advertising, data were collected for 41 magazines in 1986. There was no censoring. Let Z denote the number of pages of advertising and X the advertising revenue.

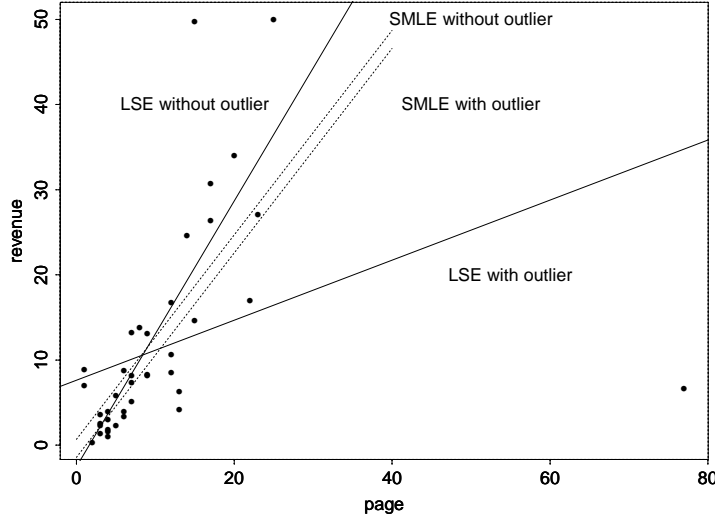
The 41 data are plotted in Figure 1. Roughly speaking, there are three outliers in the data set. They are $(25, 50)$, $(15, 49.7)$, $(77, 6.6)$. The SMLE of β is unique for this data set. The SMLE and the LSE are significantly different (see the first block of Table 2). The entries in the second block of Table 2 are results after deleting the three outliers. From Table 2, it is seen that the SMLE of β does not change after deleting outliers, though the estimate of α changes.

In Figure 1, we also plot the fitted straight lines with and without deleting those three outliers. We further plot the fitted line by the SMLE method without deleting the outliers. From Figure 1, it is seen that the fitted line by the SMLE approach using the original data is very close to the least squares fitted line after deleting outliers. This suggests that the

SMLE is robust while the LSE is not.

Table 2. Results on estimating (α, β)			
		SMLE (SE)	LSE (SE)
with outliers	β	1.200 (0.196)	0.353 (0.1449)
	α	-1.427 (3.178)	7.604 (2.3466)
without outliers	β	1.200 (0.1379)	1.238 (0.138)
	α	-0.642 (1.410)	-0.962 (1.409)

Fig. 1. SMLE v.s. LSE (Journal Data)



Consistency is proved for the mixed case IC model (Yu and Wong (2006)) and is under preparation for the RC model. The latter case is quite complicated. Let $\hat{f}_{\mathbf{b}}(t) = \hat{S}_{\mathbf{b}}(t-) - \hat{S}_{\mathbf{b}}(t)$. If ϵ is continuous and there is no censoring, then

$$\hat{f}_{\mathbf{b}}(T_i(\mathbf{b})) = 1/n \forall i, \text{ except perhaps for two, at which } \hat{f}_{\mathbf{b}}(T_i(\mathbf{b})) = 2/n. \quad (2.2)$$

Consequently,

$$\begin{aligned} L(\hat{S}_{\hat{\beta}}, \hat{\beta}) &= 2^2/n^n \text{ and } L(\hat{S}_{\hat{\beta}}, \mathbf{b}) = 1/n^n \text{ if } \mathbf{b} \text{ is not an SMLE} \\ \left| \frac{1}{n} \ln L(\hat{S}_{\hat{\beta}}, \hat{\beta}) - \frac{1}{n} \ln L(\hat{S}_{\hat{\beta}}, \mathbf{b}) \right| &= \begin{cases} \frac{1}{n} \ln 4 & \text{if } \mathbf{b} \text{ is not an SMLE} \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (2.3)$$

In other words, if ϵ is continuous, the inconsistent SMLE $\hat{\beta}$ also satisfies that

$$\left| \frac{1}{n} \ln L(\hat{S}_{\hat{\beta}}, \hat{\beta}) - \frac{1}{n} \ln L(\hat{S}_{\hat{\beta}}, \beta) \right| \rightarrow 0$$

In a consistency proof, one may want to establish

$$\lim_{n \rightarrow \infty} \left[\frac{1}{n} \ln L(\hat{S}_{\hat{\beta}}, \hat{\beta}) - E \left\{ \frac{1}{n} \ln L(S_o, \beta) \right\} \right] = 0 \text{ a.s..}$$

However, it does not work here, because $E \left\{ \frac{1}{n} \ln L(S_o, \beta) \right\} = -\infty$ due to

$$L(S_o) = (S_o(t-) - S_o(t))^\delta (S_o(t))^{1-\delta} = 0.$$

We consider a modification of the above equality in the proof.

Property 5 is still very difficult to prove. We find that a variant of the SMLE with the RC data has similar properties with the SMLE. The estimator is a value of \mathbf{b} that maximizes

$\prod_{i=1}^n (\hat{S}_{\mathbf{b}}(T_i(\mathbf{b})))^{1-\delta_i}$
instead of

$\prod_{i=1}^n (\hat{f}_{\mathbf{b}}(T_i(\mathbf{b}))^{\delta_i} (\hat{S}_{\mathbf{b}}(T_i(\mathbf{b})))^{1-\delta_i}$,
and is called the partial likelihood SMLE (PSMLE). Table 3 presents simulation studies when ϵ , Y and $Z (= \mathbf{Z})$ have distributions $Exp(3, 1)$, $Exp(1, 1)$ and $Exp(0, 1)$, respectively. $(\alpha, \beta) = (3, 1)$.

Table 3. Simulation Results on estimating (α, β)

n		β ($\sigma_{\hat{\beta}}$)	SMLE(SE)	PSMLE(SE)	BJE(SE)
200	sample mean	1	1.06	1.18	0.58
	sample SE	(0.33)	(0.47)	(1.47)	(1.23)
	est. of SE		0.48	0.51	
800	sample mean	1	1.00	1.01	0.86
	sample SE	(0.17)	(0.20)	(0.22)	(0.64)
	est. of SE		0.20	0.20	
1000	sample mean	1	1.02	1.03	0.91
	sample SE	(0.15)	(0.16)	(0.17)	(0.39)
	est. of SE		0.17	0.18	

Remark. Under the semiparametric model, $X = \beta' \mathbf{Z} + \epsilon$, the location parameter $\alpha = E(\epsilon)$ is not identifiable under censoring. This is the major reason why people do not consider the model $X = \alpha + \beta' \mathbf{Z} + e$, where $E(e) = 0$. In fact, if b is fixed, the likelihood function $\prod_{i=1}^n \mu_{\hat{F}_{a,b}}(I_i - a - b' \mathbf{z}_i)$ is constant in $a \in R$, where $\hat{F}_{a,b}$ is the GMLE of F_o based on $I_i - a - b' \mathbf{z}_i$'s

§5.4.4.2. Homework.

1. There are 4 observations $(M_i, \delta_i, \mathbf{z}_i)$'s: $(3, 1, 2)$, $(4, 0, 1)$, $(1, 1, 1)$, $(7, 1, 2)$. Find the SMLE and PSMLE of β under the linear regression model.

§5.4.5. A modified SMLE approach.

It is a little bit disappointed that the SMLE approach is not efficient when F_o is continuous, though it is super efficient if F_o is discontinuous.

It may due to the reason that $f(t) = F(t) - F(t-)$ in the nonparametric likelihood.

A modification is to replace f in \mathbf{L} by a smooth version, a kernel estimator of f_o , i.e.,

$$f(t) = \frac{1}{h} \int K((x-t)/h) dF(x),$$

where $K(\cdot)$ is a kernel. However, for most versions, it is difficult to find directly the maximum point of \mathbf{L} , one can only find a critical point of \mathbf{L} , or a root to $\frac{\partial \mathbf{L}}{\partial b}$. Such approach is called M-estimation approach, which requires that $\{x : K'(x) \neq 0\}$ is not a null set.

Yu and Wong (2005) propose a different modification for RC data. Let

$$f_F(x) = \int_{-\infty}^{\infty} \frac{1}{h_n} K\left(\frac{x-t}{h_n}\right) dF(t), \quad K(x) = \frac{1}{2} \mathbf{1}_{(-1 < x \leq 1)}, \quad h_n > 0, \quad \lim_{n \rightarrow \infty} h_n = 0 \quad (5.1)$$

(e.g., $h_n = O(n^{-1/5})$, as suggested in Härdle (1990, p.59 or p.91)), and $\hat{F}_{\mathbf{b}} = 1 - \hat{S}_{\mathbf{b}}$ be the PLE based on $T_i(\mathbf{b}) = M_i - \mathbf{b} \mathbf{z}_i$, $i = 1, \dots, n$. Then

$$\mathbf{L}(\hat{S}_{\mathbf{b}}, \mathbf{b}, f_{\hat{F}_{\mathbf{b}}}) = \prod_{i=1}^n \left[\frac{[\hat{S}_{\mathbf{b}}(T_i(\mathbf{b}) - h_n) - \hat{S}_{\mathbf{b}}(T_i(\mathbf{b}) + h_n)]}{2h_n} \right]^{\delta_i} (\hat{S}_{\mathbf{b}}(T_i(\mathbf{b})))^{1-\delta_i}. \quad (5.2)$$

We propose to estimate β by $\hat{\beta}$ which maximizes $\mathbf{L}(\hat{S}_{\mathbf{b}}, \mathbf{b}, f_{\hat{F}_{\mathbf{b}}})$ over $\mathbf{b} \in \mathcal{R}^p$.

We call $\hat{\beta}$ the MSMLE of β . Then $\hat{S}_{\hat{\beta}}(t)$ is an MSMLE of $S_o(t)$, where $\epsilon \sim F_o$.

Let $\alpha = E(\epsilon)$ if $E(\epsilon)$ exists.

Even if it does, it is well known (see Buckley and James (1979)) that there is no consistent estimator of α under right censoring, unless some further assumptions are made.

Nevertheless, a natural estimator of α is

$$\hat{\alpha} = a(\hat{\beta}), \text{ where } a(\hat{\beta}) = \frac{\int_{x \leq T_{(n)}(\hat{\beta})} x d\hat{F}_{\hat{\beta}}(x)}{\int_{x \leq T_{(n)}(\hat{\beta})} 1 d\hat{F}_{\hat{\beta}}(x)}, \quad (5.3)$$

where $T_{(1)} \leq \dots \leq T_{(n)}$ are order statistics of T_i 's.

Though the MSMLE is motivated for continuous F_o ,

it has some nice properties for arbitrary F_o (continuous or discontinuous).

Since $\frac{1}{2h_n}$ in (5.2) does not depend on \mathbf{b} , it suffices to maximize for $\mathbf{b} \in \mathcal{R}^p$,

$$l(\mathbf{b}) = \prod_{i=1}^n [[\hat{S}_{\mathbf{b}}(T_i(\mathbf{b}) - h_n) - \hat{S}_{\mathbf{b}}(T_i(\mathbf{b}) + h_n)]^{\delta_i} (\hat{S}_{\mathbf{b}}(T_i(\mathbf{b})))^{1-\delta_i}]. \quad (5.4)$$

By a similar argument as for the SMLE, the new likelihood function takes on finitely many values and is constant on intervals of the form (a_i, a_{i+1}) , where

$$a_1 < \dots < a_m \text{ are all the distinct values of } b = \frac{M_i - M_j + kh_n}{z_i - z_j}, z_i \neq z_j, k = 0, \pm 1, \pm 2. \quad (5.5)$$

The argument is similar to that for the SMLE.

Thus it can be obtained by a non-iterative algorithm.

We need the following notation.

Let \mathcal{A} be the set of points satisfies (5.5).

Let \mathcal{A}_1 be the set of ordered distinct elements of \mathcal{A} ;

Let $a_0 = -\infty$ and $a_{m+1} = \infty$;

Let \mathcal{A}_2 be the set consisting of $a_1 - 1$, $a_m + 1$ and points $\frac{a_{i-1} + a_i}{2}$, $i = 2, \dots, m$.

Non-iterative algorithm:

1. Obtain \mathcal{A} , \mathcal{A}_1 and \mathcal{A}_2 ;
2. Compute $l(b)$ (see (5.4)) for each $b \in \mathcal{A}_1 \cup \mathcal{A}_2$.
3. Each b that maximizes $l(b)$ over $b \in \mathcal{A}_1 \cup \mathcal{A}_2$ is an MSMLE of β . Moreover, if b is an MSMLE and $b \in (a_i, a_{i+1})$, then each point in (a_i, a_{i+1}) is also an MSMLE.

Since $\mathcal{A}_1 \cup \mathcal{A}_2$ is finite and $l(\mathbf{b})$ has a closed-form expression, the algorithm is non-iterative.

The MSMLE $\hat{\beta}$ is consistent under certain regulation conditions. Simulation suggests that it is also efficient.

In particular, it suggests that

if F_o is continuous, it attains the efficient lower bound of the variance, namely, $-\left(\frac{\partial \log \mathbf{L}}{\partial \beta \partial \beta'}\right)^{-1}$.

If F_o is discontinuous, then $P\{\hat{\beta} \neq \beta \text{ for large enough } n\} = 0$.

The following are simulation results supporting the above conjectures. In our simulation, we assume that ϵ , \mathbf{Z} and Y are independent.

We compare $\hat{\beta}$ to the BJE in several cases as follows.

- (A) F_o is continuous (Examples 5.1, 5.3, 5.5 - 5.7), or is neither discrete nor continuous (Example 5.2).
- (B) F_o is continuous but the regularity conditions in the Cramer-Rao theorem do not hold (Examples 5.1, 5.5 and 5.6), or all underlying distributions are exponential distributions so that they allow exchange of differentiation and integration (Example 5.7), or F_o is a normal distribution function (Examples 5.3 and 5.5).
- (C) There is no censoring (Examples 5.1- 5.3), or there is censoring (Examples 5.4-5.7).

In our simulation, for each case, we repeated 1000 times and computed the sample mean and sample standard error (SE) of the 1000 estimates.

Example 5.1. Suppose $\epsilon \sim U(-1, 1)$ (the uniform distribution), $Z \sim U(0, 9)$ and $(\alpha, \beta) = (0, 2)$.

Example 5.2. Suppose ϵ is a mixture of $U(0, 0.9)$ and a constant 0.45, with probabilities (w.p.) 0.9 and 0.1, respectively, $Z \sim U(1, 2)$ and $(\alpha, \beta) = (0.45, 1)$.

Example 5.3. Suppose $\epsilon \sim N(0, 0.09)$, $Z \sim U(0, 9)$ and $(\alpha, \beta) = (0, 1)$.

The results of the above examples are summarized in Table 1. In the two rows of each block of Table 1, we present the sample averages and sample standard errors (SE) of the MSMLE and the BJE. It is seen that $\hat{\beta}$ dominates the LSE in the sense that $SE_{\hat{\beta}} \leq SE_{LSE}$ in general, and $SE_{\hat{\beta}} < SE_{LSE}$ unless $\epsilon \sim N(\mu, \sigma^2)$, provided $n \geq 200$. In the next 4 cases, there are right-censored data. Define $Y^c = Y - \beta'Z$ and $\tau = \sup\{t : P(Y^c < t) < 1\}$.

Table 1. Simulation Results on estimating β without censoring.				
	SMLE	parameter β	LSE	MSMLE $\hat{\beta}$
Example 5.1. (continuous F_o)				
n=32	Sample mean	2	1.996	1.993
	SE		0.040	0.045
n=200	Sample mean	2	1.999	1.998
	SE		0.016	0.014
Example 5.2. (discontinuous but not discrete F_o)				
n=32	1.001	1	1.000	1.003
	0.099		0.156	0.121
n=200	1.000	1	0.997	1.000
	0.000		0.060	0.000
Example 5.3. ($N(\mu, \sigma^2)$)				
n=32	Sample mean	1	1.000	0.998
	SE		0.022	0.025
n=200	Sample mean	1	1.0000	1.0000
	SE		0.0083	0.0083

Table 2. Simulation Results on estimating β with censoring.				
	β	SMLE (SE)	BJE (SE)	MSMLE (SE) $\hat{\beta}$
Example 5.4 (discontinuous). $F_o(\tau) < 1$.				
n=32	1		0.290	0.981
			(0.700)	(0.181)
n=200	1		0.750	1.002
			(0.226)	(0.061)
Example 5.5 ($N(\mu, \sigma^2)$). $F_o(\tau-) = 1$.				
n=32	1		1.000	0.995
			(0.030)	(0.042)
n=200	1		1.000	0.994
			(0.011)	(0.013)
Example 5.6 (continuous). $F_o(\tau-) = 1$.				
n=32	2		2.001	1.999
			(0.022)	(0.029)
n=200	2	2.077	2.000	2.000
		(0.080)	(0.007)	(0.006)
Example 5.7 ($\text{Exp}(\mu, \sigma)$). $F_o(\tau-) = 1$.				
n=32	1		0.930	1.298
			(1.565)	(1.832)
n=200	1	1.012	0.957	1.004
		(0.693)	(0.751)	(0.274)

Example 5.4. Suppose ϵ is a mixture of $U(0, 0.5)$ and 51, w.p. 0.5 and 0.5, respectively, $Z \sim U(1, 2)$, $Y \sim U(4, 4.1)$, and $(\alpha, \beta) = (25.625, 1)$.

Example 5.5. Suppose $\epsilon \sim N(0, 0.09)$, $Z \sim U(0, 9)$ and Y equals 0.5 and 39 w.p. 0.5 and 0.5, respectively. $(\alpha, \beta) = (0, 1)$.

Example 5.6. Suppose $\epsilon \sim U(-1, 1)$, $Z \sim \text{Exp}(0, 19)$, $Y \sim \text{Exp}(0, 25)$, $(\alpha, \beta) = (0, 2)$.

Example 5.7. Suppose ϵ , Y and \mathbf{Z} have distributions $\text{Exp}(5, 2)$, $\text{Exp}(3, 4)$ and $\text{Exp}(2, 2)$, respectively. $(\alpha, \beta) = (5, 1)$.

The simulation results of Examples 5.4 - 5.7 are summarized in Table 2. In Table 3, we compare the sample variance of the MSMLE to the ELB under the exponential distribution (Example 5.7). We do not compare $\hat{\beta}$ to the ELB in other examples, as the ELB is not valid in Examples 5.1, 5.2, 5.4 and 5.6, and as $\hat{\beta}_{BJE}$ is efficient in Examples 5.3 and 5.5. In Table 4, we give the empirical relative efficiency of $\hat{\beta}$ to the BJE, based on results not necessarily in Tables 3 and 4. The following are main observations from our simulation.

(1) All the 7 examples suggest that the MSMLE $\hat{\beta}$ is consistent, as the values of β are all within 2 SE's from the sample means and the SE's are decreasing in n .

(2) The results suggest that unless F_o is a normal distribution, in general the MSMLE $\hat{\beta}$ is asymptotically more efficient than the BJE as the SE's of $\hat{\beta}$ are uniformly smaller than those of the BJE when sample sizes are large in all but Examples 5.3 and 5.5, and the 7 examples include different types of distributions specified in cases (A), (B) and (C).

(3) If F_o is discontinuous, then $SE_{\hat{\beta}} = 0$ for large sample sizes, while the SE of the BJE never equals 0. It suggests that (5.1) holds when F_o is neither discrete nor continuous, rather than only when F_o is discrete. Here we shall give a heuristic explanation as follows. For simplicity, let $h_n = 0$ in (2.4), and consider the case of complete data in Example 5.2. Then (2.4) becomes $l(b) = \prod_{i=1}^n \hat{f}_b(T_i(b))$, where \hat{f}_b is the empirical density based on $T_i(b)$'s. Note $T_i(b) = M_i - bZ_i$ and $\epsilon_i = T_i(\beta)$. Let $n_1 = \sum_{i=1}^n \mathbf{1}_{(\epsilon_i=0.45)}$. If n is large enough, one expects that $n_1 \geq 10$. If $b = \beta$, then there are n_1 $T_i(b)$'s that equal 0.45, thus $l(\beta) = \left(\frac{n_1}{n}\right)^{n_1} \left(\frac{1}{n}\right)^{n-n_1}$. If $b \neq \beta$, then one expects that $T_i(b)$'s are all distinct and thus $l(b) = \left(\frac{1}{n}\right)^n$. Now it is easy to see that $b = \beta$ is the MSMLE of β when n is large.

(4) Simulation results in Examples 5.3 and 5.7 further suggest that the MSMLE is efficient. In Example 5.3, data are not censored and $\epsilon \sim N(\mu, \sigma^2)$. Thus the LSE is efficient. Since $SE_{\hat{\beta}}^2 / SE_{\hat{\beta}_{LSE}}^2 = 1$ (see Table 1), the MSMLE should also be efficient. In example 5.7, the efficient lower bound (ELB) of $\hat{\beta}$ is $Var(\epsilon)/(n \cdot Var(Z)) = 2.5^2/n$ (see (3.1)). It is seen from Table 3 that when $n = 800$, the MSMLE practically attains the ELB. Note that, in Table 3, $\hat{\sigma}_{\hat{\beta}}^2$ stands for the sample variance in the simulation.

Table 3. Comparison between the SE of the MSMLE and the ELB

$n :$	32	100	200	400	800	$\sqrt{n \cdot ELB}$
$\sqrt{n} \hat{\sigma}_{\hat{\beta}} :$	(10.363	5.707	3.875	3.132	2.503)	$\left(\begin{array}{c} 2.5 \\ \end{array} \right)$

(5) It is well known that if $\epsilon \sim N(\mu, \sigma^2)$ such as in Examples 5.3 and 5.5, the BJE is efficient. From Table 2, we note that the BJE is still better than $\hat{\beta}$ when $n = 200$ in Example 5.5, while $SE_{\hat{\beta}} = SE_{\hat{\beta}_{BJE}}$ in Example 5.3. The results have two opposite interpretations:

(5.a) $\hat{\beta}$ with right-censored data is not efficient, (5.b) $\hat{\beta}$ may be efficient but the sample size is not large enough. In fact, from our simulation results, $\frac{SE_{\hat{\beta}}^2}{SE_{\hat{\beta}_{BJE}}^2} = 1.95, 1.40, 1.23, 1.16$ for $n = 32, 200, 300, 400$, respectively, in Example 5.5 and $\frac{SE_{\hat{\beta}_{BJE}}^2}{SE_{\hat{\beta}}^2} = 1$ in Example 5.3. Thus (5.b) is a more logical explanation. If so, it also suggests that $\hat{\beta}$ is efficient in general.

Table 4. Estimates of the relative efficiency of $\hat{\beta}$ to the BJE.

<i>Example :</i>	5.1	5.2	5.3	5.4	5.5	5.6	5.7
$\frac{\hat{\sigma}_{\hat{\beta}_{BJE}}^2}{\hat{\sigma}_{\hat{\beta}}^2} :$	$\left(\begin{array}{cccccc} 1.3 & \infty & 1.0 & 3.9 & ? & 1.8 & 3.7 \end{array} \right)$						
$F_o :$	$\left(\begin{array}{cccccc} \text{unif.} & \text{mixture} & \text{normal} & \text{mixture} & \text{normal} & \text{unif.} & \text{expon.} \end{array} \right)$						

See Observation (5) above for “?” in Table 4.

Example 5.8. (The Stanford heart transplant data). The data and detailed description can be found in Miller (1981, p.156). In this data, right-censored survival time, indicator of death, and five covariates including age of the recipient at the time of transplant were recorded. $n = 69$. For illustrated purpose, several methods are compared using the logarithm of time until death against age. The a priori guess H_1 under the AL model would be that younger patients fare better, that is $H_0 : \beta < 0$. In Table 5, we compare the Miller estimator, the BJE and the Cox procedure to the MSMLE. Note that the Cox model is $P(X > t|Z = z) = (S(t))e^{bz}$, where S is a baseline survival function. Thus we expect $H_0 : b > 0$ rather than $\beta < 0$ as in the simple linear regression model.

The entries in Table 5 related to the Miller estimator and the Cox procedure as well as their SE's are taken from Miller (1981, p.156). As commented by Miller (1981, p162), “The Cox method indicates there is a highly significant age effect. The Miller method says there is no effect due to age.” The three BJE's in Table 2 basically suggest that there is no effect due to age. On the other hand, there is a unique MSMLE and is significantly negative and confirms with both the a priori guess and the Cox procedure. For this data set, taking $h_n = n^{-1/5}$ yields $\hat{\beta} = 0$, which does not lead to a satisfactory estimate. We choose $h_n = 3n^{-1/5}$.

Table 5. Regression analysis on the Stanford heart transplant data				
age at transplant v.s. all death				
	Miller (SE)	BJE (SE)	MSMLE (SE)	Cox (SE)
$H_0 :$	$\beta < 0$	$\beta < 0$	$\beta < 0$	$b > 0$
	-0.006 0.004 (0.017) 0.002 (0.016)	-0.028 (0.015)	-0.036 (0.017)	0.058 (0.023)

Example 2.1 of the section on BJE (Insulation data (Nelson 1973)). In this data there is unique MSMLE, which is -0.018 , very close to the non-root zero-crossing point BJE solutions -0.02 . Thus it is a reasonable estimator. The SMLE is -0.0416 . It is also quite consistent with the trend.

§5.4.5.2. Homework.

1. There are 4 observations $(M_i, \delta_i, \mathbf{z}_i)$'s: $(3, 1, 2)$, $(4, 0, 1)$, $(1, 1, 1)$, $(7, 1, 2)$. Find the MSMLE of β under the linear regression model (with $h_n = n^{-1/5}$).
2. Prove that the likelihood in (5.2) is a constant on the interval (a_i, a_{i+1}) as defined above based on the data in problem 1, thus there are only finitely many values.

References

- * Finkelstein, D.M. (1986). A proportional hazards model for interval-censored failure time data. *Biometrics* 42 845-854.
- * Härdle W. (1990). Smoothing techniques, with implementation in S. *Springer-Verlag*, N.Y.
- * Huber, P.J. (1964). Robust estimation of a location parameter”. *Ann. Math. Statist.* 35, 73-101.
- * Lai, T.L. and Ying, Z.L. (1991). Large sample theory of a modified Buckley-James estimator for regression-analysis with censored data. *Ann. Stat.* 19 1370-1402.
- * Li, L.L. and Pu, Z.W. (1999). Regression models with arbitrarily interval-censored observations. *Comm. in Statist., Theory and Methods* 1547-1563.
- * Rabinowitz, D. Tsiatis, A. and Aragon, J. (1995). Regression with interval-censored data. *Biometrika*, 82, 501-513.
- * Yu, Q.Q. and Wong, G.Y.C. (2002a). How to find all Buckley-James estimates instead of just one ? *Journal of Statistical Computation and Simulation*, 72, 451-460.
- * Yu, Q.Q. and Wong, G.Y.C. (2002b). Asymptotic properties of a modified semi - parametric MLE in linear regression analysis with right-censored data. *Acta Mathematica Sinica*, 18 405-416.

- * Yu, Q.Q. and Wong, G.Y.C. (2003a). The semi-parametric MLE in linear regression with right-censored data. *Journal of Statistical Computation and Simulation*, 73 833-843.
- * Yu, Q.Q. and Wong, G.Y.C. (2003b). Asymptotic properties of the generalized semi-parametric MLE in linear regression. *Statistica Sinica*, 13 311-325.
- * Zhang, C.H. and Li, X. (1996). Linear regression with doubly censored data. *Ann. Statist.*, 24, 2720-2743.
- * Yu, Q.Q. and Wong, G.Y.C. (2003c). Semi-parametric MLE in simple linear regression with interval-censored data. *Communications in Statistics-Simulation and Computation*, 32 147-164.
- * Yu, Q.Q., Wong, G.Y.C. and Kong, F.H.(2006). Consistency of the semi-parametric MLE in linear regression with interval-censored data. *Scan. J. Statist.* 33 367-378.
- * Yu, Q.Q. and Wong, G.Y.C. (2005). A modified semi-parametric MLE in linear regression analysis with complete data or right-censored data. *Technometrics*, 47 34-42.
- * Kong, F.H. and Yu, Q.Q. (2006). Asymptotic distribution of the Buckley-James estimator under non-standard conditions *Statistica Sinica* (in press).
- * Yu, Q.Q and Wong, G.Y.C. (2006) A data set that there are two BJE's.
http://www.math.binghamton.edu/qyu.

Chapter 6. Testing (large sample tests)

Let X be the survival time and \mathbf{Z} the covariate. There are several objectives in hypothesis testing:

1. In order to apply certain regression model, check whether $F_{X|\mathbf{Z}}$ belongs to some regression model, e.g. a PH model, a Lehman model, or an accelerated lifetime model.
 $H_o: h_{X|Z}(x|z) = h_o(x)e^{\beta z}$, or $H_o: \log X = \beta Z + W$, where $F_{X|Z}(\cdot|0)$ is unknown.
2. In order to apply a certain parametric model, say $X \sim F(\cdot; \theta)$, test
 $H_o: F = F(\cdot; \theta)$, where the form $F(\cdot)$ is given.
3. $H_o: F = F_o$, where F_o is a given cdf.
4. $H_o: \theta = \theta_o$, where θ is a parameter of F_X , such as the mean and the variance, or the parameter in a certain parametric distribution.

In 502, we mainly deal with type 4 testing problems, where θ is the parameter in a certain parametric distribution.

In data analysis, we deal with all the 4 types of testing problems.

There are two common approaches in constructing a test of a parameter:

- (1) MLE approach: If $\hat{\beta}$ is the MLE, then under certain assumption, a test for
 $H_o: \beta = \beta_o$ vs. $H_1: \beta \neq \beta_o$
is $\phi = \mathbf{1}((\hat{\beta} - \beta_o)' J(\hat{\beta} - \beta_o)|_{\beta=\beta_o} > \chi_{\alpha,p}^2)$, where $J = -\frac{\partial^2 \ln \mathbf{L}}{\partial \beta \partial \beta'}$, $\beta \in \mathcal{R}^p$ and $\chi_{\alpha,p}^2$ is the $(1 - \alpha)100$ -th percentile of the χ^2 distribution with degree freedom p .
(Or $\phi_1 = \mathbf{1}(|\hat{\beta} - \beta_o|/\hat{\sigma}_{\hat{\beta}} > z_{\alpha/2})$).
- (2) Score test approach: Under certain assumptions,
a test for $H_o: \beta = \beta_o$ is $\phi = \mathbf{1}(U(\beta)' J^{-1} U(\beta)|_{\beta=\beta_o} > \chi_{\alpha,p}^2)$, where $U(\beta) = \frac{\partial \ln \mathbf{L}}{\partial \beta}$.

There are some common approaches in the first 3 types testing problems.

- (1) Kolmogorov test and Smirnov test.
- (2) Convert to type (4). For example, for testing $H_o: X = \beta Z + W$, convert it to $H_o^*: \theta = 0$, assuming $X = \beta Z + \theta Z^2 + W$.

§6.1. One sample nonparametric test

Hereafter denote \hat{F} (or \hat{S}) the GMLE of F_o (or S_o).

- (1) Two-sided level- α test for $H_o: S_o(t_0) = p_0$, where p_0 is known:
For RC data, under the RC model, $\phi = \mathbf{1}\left(\frac{|S_{pl}(t_0) - p_0|}{\hat{\sigma}_{\hat{S}_{pl}(t_0)}} \geq z_{\alpha/2}\right)$,

where $\hat{\sigma}_{\hat{S}_{pl}(t_0)}^2$ is the estimate of

$$\sigma_{\hat{S}_{pl}(t)}^2 \approx S_o(t)^2 \int_0^t \frac{1}{S_o(x-)S_Y(x-)S_o(x)} dF_o(x)/n.$$

For IC data, the situation varies. If there are exact observations, in general, a test is

$$\phi = \mathbf{1}_{\left(\frac{|\hat{S}(t_0) - p_0|}{\hat{\sigma}_{\hat{S}(t_0)}} \geq z_{\alpha/2}\right)},$$

where $\Phi(z_\alpha) = 1 - \alpha$, Φ is the cdf of $N(0, 1)$, and $\hat{\sigma}_{\hat{S}(t_0)}^2$ is the estimate of $\sigma_{\hat{S}(t)}^2$ given in §4.

If there is no exact observation, then there are three possible cases corresponding to Theorems 2 and 3 and the conjecture in §4.7.

The convergence rates are $n^{-1/2}$, $n^{-1/3}$ and $(n \ln n)^{-1/3}$, respectively.

For each of the three cases, a test can be constructed, which is introduced in §4.7.

- (2) Two-sided nonparametric level- α test for $H_0: F_X = F_o$, where F_o is a known cdf:

In complete data case, we often convert it to

$H_o: \mu = \mu_o$, where $\mu = E(X)$ and $\mu_o = \int t dF_o(t)$. Then the test is

$$\phi = \mathbf{1}_{(|T| > z_{\alpha/2})}, \text{ where } T = \frac{\bar{X} - \mu}{s/\sqrt{n}}, \text{ and } s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$$

Note $E(X) = \int S_o(x) dx$ for nonnegative X . Thus the test statistic

$$T = \frac{U_o}{\hat{\sigma}_{U_o}}, \text{ where } U_o = \int_0^\infty (\hat{S}_{pl}(x) - S_o(x)) dx.$$

This motivates, for RC data, the weighted Kaplan-Meier (WKM) statistic by Pepe and Fleming (1989, Biometrics and 1991, JRSS. A.).

WKM test

$$T = \frac{U}{\hat{\sigma}_U}, \text{ where } U = \int_0^\infty W(x)(\hat{S}_{pl}(x) - S_o(x)) dx, \quad (1.1)$$

and W is a weight function, which may depend on F_0 and S_Y . For example, let $W(t)$ be an estimate of $P(X \wedge Y \geq t)$ or $1 - S_Y(t)$ etc.. In such a case,

$$\sigma_U^2 \approx \frac{1}{n} \int_0^\tau \frac{[\int_t^\tau W(u) S_o(u) du]^2}{S_o(t-) S_o(t) S_Y(t-)} dF_o(t),$$

where $\tau = \sup\{t : F_o(t) < 1, S_Y(t) < 1\}$. σ_U^2 can easily be estimated. The test ϕ with T given by (1.1) is a location test for H_0 above.

With IC data, under a discrete assumption or in the case that there exist exact observations, U in (1.1) has the following form: Suppose there are $m+1$ finite innermost intervals with right end-points $b_1 < \dots < b_{m+1}$. WLOG, we can assume that $F_o(b_m) < 1$. Suppose that $W(t) = \mathbf{1}_{(t < \tau)}$, where τ can be estimated by the largest finite observation, denoted by b_{m+1} . Then

$$U = \int_0^{\hat{\tau}} (\hat{S}(t) - S_o(t)) dt = \sum_{i=0}^m \hat{S}(b_i)(b_{i+1} - b_i) - \int_0^\tau S_o(t) dt.$$

where $b_0 = 0$. U can be rewritten as

$$U = b_0 + (b_2 - b_1, \dots, b_{m+1} - b_m)(\hat{S}(b_1), \dots, \hat{S}(b_m))^t + \int_0^\tau S_o(t) dt.$$

$$\sigma_U^2 = (b_2 - b_1, \dots, b_{m+1} - b_m) \Sigma (b_2 - b_1, \dots, b_{m+1} - b_m)^t$$

where Σ is the covariance matrix of $(\hat{S}(b_1), \dots, \hat{S}(b_m))$.

For continuous IC data, the test statistic is not simple and needs to be investigated, as it is the integration of a stochastic process and the convergence rate of the GMLE varies at least in three different cases.

Kolmogorov test: The test is $\phi = \mathbf{1}_{(U > u_\alpha)}$, where

$$U = \sup_t |\hat{F}(t) - F_o(t)|,$$

and the critical values u_α can be computed from tables for RC data (Hall and Wellner, (1980). *Biometrika*).

Smirnov test: The test is $\phi = \mathbf{1}_{(U > u_\alpha)}$, where

$$U = \int (\hat{F}(t) - F_o(t))^2 dF_o(t) \text{ or } U = \int_0^\infty (\hat{F}(t) - F_o(t))^2 dW(t),$$

where W is a measure, and the critical values u_α can be computed from tables for RC data (Korjot and Green, (1976). *Technometrics*).

For IC data, one may use bootstrap method to find u_α for both the Kolmogorov test and the Smirnov test. For instance,

Given I_1, \dots, I_n ,

(1) resample I_1^*, \dots, I_n^* , and get U_1 ;

(2) repeat step 1 N times (including the first time), and get U_1, \dots, U_N ;

(3) U_1, \dots, U_N leads to an estimate of F_U , and u_α .

Another approach is simulation. Since S_o is known, we only need to estimate the censoring distribution. We can make the assumption that the follow-up time takes values among finite L_i 's and R_i 's with equal probability. Then we can generate n observations 100 times, and thus compute U 100 times and the $100(1 - \alpha)$ sample percentile of the 100 U values is our estimate of u_α .

§6.2. Two-sample problem

Suppose that there are two independent random samples for sizes n_1 and n_2 , from survival function S_1 and S_2 , respectively. $H_0: S_1 = S_2$, v.s. $H_1: S_1(t) \geq S_2(t)$ for all t and $S_1 \neq S_2$.

(1) **WKM statistic:**

For RC data, a test is $\phi = \mathbf{1}_{(T > z_\alpha)}$, where

$$T = \frac{U}{\hat{\sigma}_U}, \text{ where } U = \int_0^\infty \hat{W}(t)[\hat{S}_1(t) - \hat{S}_2(t)]dt,$$

$$\hat{\sigma}_U^2 = \int_0^\infty \frac{[\int_t^\infty \hat{W}(u)\hat{S}_o(u)du]^2}{\hat{S}_o(t)\hat{S}_o(t-)[1 - \hat{G}_1(t-)]} d\hat{F}_o(t)/n_1 + \int_0^\infty \frac{[\int_t^\infty \hat{W}(u)\hat{S}_o(u)du]^2}{\hat{S}_o(t)\hat{S}_o(t-)[1 - \hat{G}_2(t-)]} d\hat{F}_o(t)/n_2,$$

\hat{S}_o is PLE based on pooled-sample, G_1 and G_2 are censoring cdf of samples 1 and 2, respectively, and \hat{G}_1 and \hat{G}_2 are their PLEs based on samples 1 and 2, respectively.

For discrete IC data, we only consider

$$T = \frac{U}{\hat{\sigma}_U}, \text{ where } U = \int_0^\tau [\hat{S}_1(t) - \hat{S}_2(t)]dt,$$

where τ is a fixed constant. Since $\sigma_U^2 = \text{Var}(\int_0^\tau \hat{S}_1(t)dt) + \text{Var}(\int_0^\tau \hat{S}_2(t)dt)$, it can be estimated using an approach similar to the one discussed in §6.1.

(2) **Kolmogorov test:** The test is $\phi = \mathbf{1}_{(U > u_\alpha)}$, where

$$U = \sup_t |\hat{F}_1(t) - \hat{F}_2(t)|.$$

We can use simulation to estimate u_α . For instance, consider re-sample the pooled-sample (L_i, R_i) , $i = 1, \dots, n_1 + n_2$, to generate two independent samples of sizes n_1 and n_2 . Compute the value of U with these two samples. Repeat it 100 times, use the upper 100(1 - α) percentile to be an estimate of u_α .

(3) **Smirnov test** Geskus and Groeneboom's asymptotic results (1999) on smooth functionals with IC data can be applied to the two-sample version of the test.

(4) **Gehan's generalized Wilcoxon test.**

For RC data :

Notations: (Use M or $M+$ instead of $(x \wedge y, \delta)$ representation).

n_1 observations in the first sample: a_i or a_i+ 's;

n_2 observations in the second sample: b_i or b_i+ 's.

$$U_{ij} = \begin{cases} 1 & \text{if } a_i > b_j \text{ or } a_i+ \geq b_j \\ & \text{(we know for sure that obs-}i \text{ in sample 1} > \text{obs-}j \text{ in sample 2),} \\ -1 & \text{if } a_i < b_j \text{ or } a_i \leq b_j+ \\ & \text{(we know for sure that obs-}i \text{ in sample 1} < \text{obs-}j \text{ in sample 2),} \\ 0 & \text{if not sure.} \end{cases}$$

$$U = \sum_{j=1}^{n_2} \sum_{i=1}^{n_1} U_{ij}. \text{ A test is } \phi = \mathbf{1}_{(\frac{U}{\sigma_U} > z_\alpha)}.$$

Question: $\sigma_U^2 = ?$

σ_U^2 is derived (Gehan, 1965, Biometrika), but is very complicated. Thus it is not presented here.

Mantel (1967, Biometrics) considered a different sample space and derived a different but simpler variance:

1. Pool two samples together, and denoted by c_i or c_i+ , $i = 1, \dots, n_1 + n_2$.

The first n_1 c_i 's are the first sample, follows by the second sample. That is,

$c_1 = a_1$ (or $c_1+ = a_1+$), ..., $c_{n_1} = a_{n_1}$ (or $c_{n_1}+ = a_{n_1}+$), ...,

$c_{n_1+1} = b_1$ (or $c_{n_1+1}+ = b_1+$), ..., $c_{n_1+n_2} = b_{n_2}$ (or $c_{n_1+n_2}+ = b_{n_2}+$).

2. Denote

$$V_{kh} = \begin{cases} 1 & \text{if we know for sure obs-}k > \text{obs-}h, \\ -1 & \text{if we know for sure obs-}k < \text{obs-}h, \\ 0 & \text{if not sure,} \end{cases} \quad (2.1)$$

(for RC data, we know for sure obs- $k >$ obs- h iff $c_k > c_h$ or $c_k+ \geq c_h$.)

$$V_k = \sum_{h=1}^{n_1+n_2} V_{kh}, \quad (2.2)$$

Let W be a random variable taking values $\sum_{i=1}^{n_1} V_{k_i}$, where $\{k_1, \dots, k_{n_1}\}$ is a selection of n_1 distinct integers from $\{1, \dots, n_1 + n_2\}$.

$$V = \sum_{k=1}^{n_1} V_k \quad (2.3)$$

is a value of W . Under this sample space (permutation sample space, each permutation has equal probability),

$$E(W) = 0 \text{ and } \sigma_W^2 = n_1 n_2 \sum_{k=1}^{n_1+n_2} \frac{V_k^2}{(n_1 + n_2)(n_1 + n_2 - 1)}. \quad (2.4)$$

$$\left(E(W^2) = \sum_{k_1, \dots, k_{n_1}} \frac{1}{\binom{n_1+n_2}{n_i}} \left(\sum_{i=1}^{n_1} V_{k_i} \right)^2 \right).$$

Mantel suggests that

$$\frac{W}{\sigma_W} \xrightarrow{D} N(0, 1) \text{ as } n \rightarrow \infty.$$

We observe $W = V$.

If H_1 is true, V should be large. Thus a test is $\psi = \mathbf{1}_{(\frac{W}{\sigma_W} > z_\alpha)}$. That is $\psi = \mathbf{1}_{(\frac{V}{\sigma_W} > z_\alpha)}$. Under such a set-up, U is a value of W , in fact $U = V$.

$$\begin{aligned} V &= \sum_{k=1}^{n_1} \sum_{h=1}^{n_1+n_2} V_{kh} \\ &= \sum_{k=1}^{n_1} \sum_{h=1}^{n_1} V_{kh} + \sum_{k=1}^{n_1} \sum_{h>n_1}^{n_1+n_2} V_{kh} \\ &= \sum_{k=1}^{n_1} \sum_{h=1}^{n_1} V_{kh} + \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} U_{ij} \\ &= 0 + U \quad (\text{as } V_{kh} = -V_{hk}). \end{aligned}$$

Thus σ_W^2 can be viewed as a pseudo variance of U (**not really** a variance of U) and the Gehan's generalized Wilcoxon test is

$$\psi = \mathbf{1}_{(\frac{U}{\sigma_W} > z_\alpha)}. \quad (2.5)$$

For IC data we are sure that observation- $k >$ observation- h iff

$$\text{either } L_k > R_h \text{ or } R_k > L_k = R_h.$$

Using this interpretation in (2.1), (2.4) still holds. Thus the generalized Wilcoxon test can be extended to the IC data by using (2.5) directly.

Remark. For a random sample of RC data, say (M_i, δ_i) , $i = 1, \dots, n$.

$$\begin{aligned} \hat{S}_X(t) &= \prod_{i: M_{(i)} \leq t} \left(1 - \frac{\delta_{(i)}}{n-i+1} \right) \rightarrow S_X(t) \text{ a.s.}, \\ \hat{S}_M(t) &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}(M_i > t) \rightarrow S_M(t) \text{ a.s.} \end{aligned}$$

$$\hat{S}_M(t) = \hat{S}_X(t) \hat{S}_Y(t) \text{ and } P(M > t) = P(X > t)P(Y > t).$$

$$\hat{S}_M(t) / \hat{S}_X(t) \rightarrow S_Y(t) \text{ a.s.},$$

$$\hat{S}_M(t) / \hat{S}_X(t) = \hat{S}_Y(t) ?$$

$$\hat{S}_Y(t) = \prod_{i: M_{(i)} \leq t} \left(1 - \frac{1-\delta_{(i)}}{n-i+1} \right) \rightarrow S_Y(t) \text{ a.s.},$$

Thus $1 - \delta = \mathbf{1}(Y < X) \neq \mathbf{1}(Y \leq X)$ does not matter to the consistency of $\hat{S}_Y(t)$.

(5) **Logrank test** (Mantel (1966)) The test is a common test in medical research.

Notations:

c_1, \dots, c_{m+n} are pooled RC observations (sample 1, sample 2). Among the pooled-sample

$$\begin{cases} t_1 < \dots < t_k & \text{are distinct exact times,} \\ d_{\cdot,i} & \text{is the \# of deaths at } t_i, i = 1, \dots, k, \\ \mathcal{R}(t_i) & \text{is the set of individuals in risk at time } t_i. \end{cases}$$

$d_{j,i}$ = # of deaths from group j at time t_i .

$d_{j\cdot}$ = # of deaths in group j .

$r_{j,i}$ = # of elements of group j in $\mathcal{R}(t_i)$.

$r_{\cdot,i} = r_{1,i} + r_{2,i}$.

$U = \sum_{j=1}^k [d_{1,j} - d_{\cdot,j} \frac{r_{1,j}}{r_{\cdot,j}}]$.

$$\frac{U}{\hat{\sigma}_U} \xrightarrow{D} N(0,1), \text{ where } \hat{\sigma}_U^2 = \sum_{i=1}^k d_{\cdot,i} \frac{r_{2,i} r_{1,i} (r_{\cdot,i} - d_{\cdot,i})}{r_{\cdot,i}^2 (r_{\cdot,i} - 1)} \quad (2.6)$$

The level- α logrank test is $\phi = \mathbf{1}_{(\frac{U}{\hat{\sigma}_U} < -z_\alpha)}$.

Derivation of logrank test

Under continuous assumption and PH model,

$z_i = z_i(t) = w(t) \mathbf{1}_{(\text{observation } i \text{ is from group } 1)}$, where $w(t)$ is a weight function.

$\mathcal{D} = \{i : \text{observation } i \text{ in pooled-sample died}\}$.

$$S_1 = S(t|0) \text{ and } S_2 = S(t|1), \text{ where } S(t|z) = (S_o(t)) e^{\beta z_i(t)}, \quad (2.7)$$

$H_0: S_1 = S_2$ is equivalent to $\beta = 0$.

$$lik = \prod_{i \in \mathcal{D}} \frac{e^{\beta z_i(c_i)}}{\sum_{h \in \mathcal{R}(c_i)} e^{\beta z_h(c_i)}}$$

$$U(\beta) = \frac{d \ln lik}{d\beta} = \sum_{i \in \mathcal{D}} [z_i(c_i) - \frac{\sum_{j \in \mathcal{R}(c_i)} e^{\beta z_j(c_i)} z_j(c_i)}{\sum_{h \in \mathcal{R}(c_i)} e^{\beta z_h(c_i)}}], \text{ --- the score function.}$$

$$U(0) = \sum_{i \in \mathcal{D}} [z_i(c_i) - \frac{\sum_{j \in \mathcal{R}(c_i)} z_j(c_i)}{|\mathcal{R}(c_i)|}] = \sum_{i \in \mathcal{D}} [z_i(c_i) - \sum_{j \in \mathcal{R}(c_i)} \frac{z_j(c_i)}{|\mathcal{R}(c_i)|}]$$

$$U(0) = \sum_{j=1}^k w(t_j) [d_{1,j} - d_{\cdot,j} \frac{r_{1,j}}{r_{\cdot,j}}], \text{ as } z_i = 0 \text{ for } i > n_1.$$

If $w(t_i) = 1$, $U(0) = \sum_{i=1}^{n_1} [d_{1,j} - d_{\cdot,j} \frac{r_{1,j}}{r_{\cdot,j}}]$, corresponding to the logrank test.

The asymptotic variance of $U(0)$ can be obtained by

$$\begin{aligned} \sigma_U^2 &\approx - \frac{dU(\beta)}{d\beta} \Big|_{\beta=0} = \sum_{i \in \mathcal{D}} \left[\frac{\sum_{j \in \mathcal{R}(c_i)} e^{\beta z_j} z_j^2}{\sum_{h \in \mathcal{R}(c_i)} e^{\beta z_h}} - \frac{\sum_{j \in \mathcal{R}(c_i)} e^{\beta z_j} z_j \sum_{l \in \mathcal{R}(c_i)} e^{\beta z_l} z_l}{(\sum_{h \in \mathcal{R}(c_i)} e^{\beta z_h})^2} \right] \Big|_{\beta=0} \\ &= \sum_{i=1}^k (w(t_i))^2 d_{\cdot,i} \left[\frac{r_{1,i}}{r_{\cdot,i}} - \frac{r_{1,i}^2}{r_{\cdot,i}^2} \right] \\ &= \sum_{i=1}^k [(w(t_i))^2 d_{\cdot,i} \frac{r_{1,i} r_{2,i}}{r_{\cdot,i} r_{\cdot,i}}] \end{aligned} \quad (2.8)$$

Note that the logrank test is a score test, not a test based on MLE. Thus the Fisher information matrix is the variance of the score function.

$$\hat{\sigma}_U^2 = - \frac{dU(\beta)}{d\beta} \Big|_{\beta=0} = \sum_{i=1}^k [d_{\cdot,i} \frac{r_{1,i} r_{2,i}}{r_{\cdot,i} r_{\cdot,i}}] \text{ if } w(t_i) = 1,$$

which equals (2.6) as $d_{.,i} = 1$ by the continuity assumption.

A two-sided test $\psi = \mathbf{1}_{(|\frac{U(0)}{\sigma_U}| > z_{\alpha/2})}$ is called a linear rank test, where

$w(t_i) = r_{.,i}$ — generalized Wilcoxon test, Gehan (1965);

$w(t_i) = 1$ — logrank test Mantel (1965);

$w(t_i) = n\hat{S}_{pl}(t_i-)$ — Prentice (1978);

$w(t_i) = n(\hat{S}_{pl}(t_i-))^k$ — Harrington and Fleming (1982).

Example in R.

`x=coxph(Surv(time)~ ag+log(wbc),data=leuk)`

`summary(x)`

	<i>coef</i>	<i>exp(coef)</i>	<i>se(coef)</i>	<i>z</i>	<i>Pr(> z)</i>	
<i>agpresent</i>	-1.0691	0.3433	0.4293	-2.490	0.01276	*
<i>log(wbc)</i>	0.3677	1.4444	0.1360	2.703	0.00687	**

	<i>exp(coef)</i>	<i>exp(-coef)</i>	<i>lower.95</i>	<i>upper.95</i>
<i>agpresent</i>	0.3433	2.9126	0.148	0.7964
<i>log(wbc)</i>	1.4444	0.6923	1.106	1.8857

Concordance= 0.726 (se = 0.065)

Rsquare= 0.377 (max possible= 0.994)

Likelihood ratio test= 15.64 on 2 df, p=4e-04

Wald test = 15.06 on 2 df, p=5e-04

Score (logrank) test = 16.49 on 2 df, p=3e-04

Remark. The test statistics of the logrank test as all existing tests for the PH model are based on the assumption that the data are from a model larger than the model in H_0 . In particular, it assumes (2.7), that is, $S_1 = S(t|0)$ and $S_2 = S(t|1)$, where $S(t|z) = (S_o(t))^{e^{\beta z(t)}}$ and $z(t)$ is given. Then it tests $H_o^*: \beta = 0$ v.s. $H_1^*: \beta \neq 0$, based on Eq. (2.7).

The size of the test is true as long as n is large, that is, when H_o is true.

If (2.7) does not hold, then it is not true that $\frac{U}{\sigma_U} \approx N(0, 1)$ and the test is not valid.

When does the assumption fail ?

1. $z(t)$ is mis-specified;
2. the data is from another regression model, *e.g.* a log linear regression model;
3. the data is not from any common regression model.

If the assumption fails, then the logrank test becomes a random guessing. One can check by simulation study that it is possible that the logrank test rejects H_0 with a probability 0.5, a large probability, or with a small probability.

(6) The marginal distribution test (in two-sample problem) (Dong and Yu (2018)).

Let $(M_1, Z_1, \delta_1), \dots, (M_n, Z_n, \delta_n)$ be a random sample, where $Z_i \in \{1, 2\}$.

H_o : $S_1 = S_2$ v.s. $S_1 \neq S_2$.

A test statistics is

$$T = \int |\hat{S}_X(t) - \hat{S}_{X^*}(t)| d\hat{S}_X(t),$$

where \hat{S}_X is the PLE based on (M_i, δ_i) 's,

$\hat{S}_{X^*}(t) = \frac{1}{n} \sum_{i=1}^n (\hat{S}_1(t))^{\beta Z_i}$ and

$\hat{S}_1(t)$ is the PLE based on the first sample. In particular, if the first sample is complete, $\hat{S}_1(t) = \frac{1}{\sum_{i=1}^n \mathbf{1}_{(Z_i=1)}} \sum_{j=1}^n \mathbf{1}_{(M_j > t, Z_j = 1)}$.

Notice that $\hat{F}_1 = 1 - \hat{S}_1$ (the edf based on the first sample).

The critical value can be obtained by a modified bootstrap method.

Justification. $F_{X,Z}$ is a joint cdf, which does not need to be from any PH model.

F_X is the marginal distribution.

If $F_{X|Z}$ is from a PH model, then

$$F_X = E(F_{X|Z}(\cdot|Z)) = E((S_1(t))^{\beta Z}),$$

which can be estimated by $\hat{S}_{X^*}(t) = \frac{1}{n} \sum_{i=1}^n (\hat{S}_1(t))^{\beta Z_i}$.
 Otherwise, $F_{X^*} = E((S_1(t))^{\beta Z})$ satisfies the PH model.

Homework Solution.

§5.4.3.2. Homework.

Derive the expressions of Φ when f is the density of $U(0, \theta)$ and when $f(x) = e^{-x}$, $x > 0$.

Solution. There are two expressions for Φ :

$$\Phi(x) = \sum_{i=1}^n (z_i - \bar{z}) \left(\delta_i \left(\frac{f'}{f} \right) (M_i - \beta' z_i) - (1 - \delta_i) \left(\frac{f}{S} \right) (M_i - \beta' z_i) \right). \quad (3.1)$$

$$\Phi(x) = \sum_{i=1}^n (z_i - \bar{z}) \left(\delta_i \left(\frac{f'}{f} \right) (M_i - \beta' z_i) - (1 - \delta_i) \frac{\int_{t > T_i(\beta)} \left(\frac{f'}{f} \right) (x) dS(x)}{S(T_i(\beta))} \right) \quad (2)$$

under certain regularity conditions (which is not applicable to $(U(0, \theta)$, as f is discontinuous
Why we say so ?

Now if $f(x) = e^{-x} = S(x)$, $x > 0$.

$f'/f = -1$. $f/S = 1$, thus

$$\Phi(x) = \sum_{i=1}^n (z_i - \bar{z}) (-\delta_i - (1 - \delta_i)) = 0.$$

Does it make sense (as $\Phi = (-\log \mathbf{L})'$) ?

$$\begin{aligned} L &= \prod_{i=1}^n (f(T_i(\beta)))^{\delta_i} (S(T_i(\beta)))^{1-\delta_i} / \\ f(x) &= S(x) ?? \\ f(x) &= e^{-x} \mathbf{1}(x > 0), \\ S(x) &= \mathbf{1}(x < 0) + \mathbf{1}(x \geq 0) e^{-x} = e^{-x} \mathbf{1}(x \geq 0) \end{aligned}$$

$$\begin{aligned} \mathbf{L} &= \prod_{i=1}^n [\mathbf{1}(M_i - \beta' z_i > 0) \exp(-(M_i - \beta' z_i))]^{\delta_i} \\ &\quad \times \prod_{i=1}^n \exp(-(M_i - \beta' z_i)) \mathbf{1}(M_i - \beta' z_i > 0) (1 - \delta_i) \end{aligned}$$

$$\begin{aligned} f'/f &= \begin{cases} -1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \\ dne & \text{if } x = 0 \end{cases} \stackrel{D}{=} -\text{sign}(x). \\ f/S &= \mathbf{1}(x \geq 0). \end{aligned}$$

$$\begin{aligned} \Phi(x) &= \mathbf{1}(\min_i (X_i - \beta' z_i) \leq 0) \left[\sum_{i=1}^n (-\delta_i \text{sign}(M_i - \beta' z_i) + (1 - \delta_i) \text{sign}(M_i - \beta' z_i)) \right] \text{ (by Eq. (2))} \\ &= \mathbf{1}(\min_i (X_i - \beta' z_i) \leq 0) \left[\sum_{i=1}^n (1 - 2\delta_i) \text{sign}(M_i - \beta' z_i) \right]. \end{aligned}$$