SURVIVAL ANALYSIS (Lecture Notes) by Qiqing Yu Version 7/3/2024

This course will cover parametric, non-parametric and semi-parametric maximum likelihood estimation under 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 and the mixture censorship model will be used to formulate the censored data. Multivariate censorship models and the two-sample problem may also be introduced.

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.

The lecture notes would be the textbook, posted in my website.

Math 557 MWF 8-9:30 am

Classroom : WH 329

Office: WH 132

Office hours: M 4:05-5:05pm in my office;

Tu 7:00pm-8:00pm through zoom

https://binghamton.zoom.us/j/8265526594?pwd=d3l6OGx1cmZ4M3cxZEJwVGd1RGcrUT09 Meeting ID: 826 552 6594

Passcode: 031320

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

B = 75

Homework due Wednesday, **email to** qiyu@binghamton.edu

Homework assigned during a week is due next Wednesday in class.

use pdf format if it is pictures, do not zip

Remind me if you do not see it by Saturday morning !

There is a homework due this Friday.

Sample of the Latex homework is in **report and hw-solution** items in my website under Data Analysis item.

Quiz: Once a week on Friday,

Quiz formulas on Math 447 and 448 (at my personal website).

 $https://www2.math.binghamton.edu/p/people/qyu/qyu_personal$

Midterm Exam: Oct. 21(M)

Final: Dec. 12 Th, 8-10am

You can bring one page with R commands and formulas in exams, bring a simple calculator.

The lecture note is also on my website "Course materials for Survival Analysis"

First item is the recent lecture note, and item 2 is the cumulative one.

Chapter 1. Introduction.

$\S1$. Two main characters of survival analysis.

Suppose X is a random variable,

with the cumulative distribution function (cdf),

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

Assume $X_1, ..., X_n$ are independently and identically distributed (i.i.d.) from F. Statistical inferences: $\begin{cases} estimation \\ testing hypotheses, \end{cases}$ on S = 1 - F (called the survival function) 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, *e.g.*, U(a, b); nonparametric, *i.e.*, F is only known to be a c.d.f. semi-parametric, *i.e.*, in between the above two *e.g.* $Y = \alpha + \beta Z + \epsilon$, F_{ϵ} unknown. 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.

Two main character of survival analysis:

(1) $X \ge 0$,

(2) incomplete data.

(1) $X \ge 0$, referred as survival time or failure time.

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

S(2 years) - - - - - the chance of surviving more than 2 years.

F(2 years) - - - - - the chance of dying within 2 years.

the larger the survival probability S, the better.

(2) Incomplete data.

Definition: An observation on X_i is called

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

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

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) .

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

$\S1.2.$ Right censoring.

§1.2.1. Representations of an RC observation:

$$(L_i, R_i)$$
 — a vector,
 $I_i = (L_i, \infty)$ — an interval,
 $(L_i, 0)$ — a vector,
 L_i+ ,
of an exact observation:

Representations of an exact observation:

$$(L_i, R_i) = (X_i, X_i)$$

$$I_i = [X_i, X_i] -- \text{ an interval,}$$

$$(L_i, 1) -- \text{ a vector,}$$

$$L_i.$$

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

 $(L_i, R_i), i = 1, ..., n, --$ random vectors,

where $(L_i, R_i) = \begin{cases} (X_i, X_i) & \text{if the observation is exact,} \\ (L_i, \infty) & \text{if the observation is RC,} \end{cases}$ $I_i, i = 1, ..., n, ---$ random intervals, where $I_i = \begin{cases} [X_i, X_i] & \text{if the observation is exact,} \\ (L_i, \infty) & \text{if the observation is RC.} \end{cases}$ $(L_i, \delta_i), i = 1, ..., n, ---$ random vectors, where $\delta_i = \mathbf{1}_{\text{(the i-th observation is exact)}}$, 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. $L_i + \text{ or } L_i, i = 1, ..., n$.

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 congregative either linear X_i if the 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. (A testing on the lifetimes of n tires in a lab). Each individual was followed by a fixed time c. Each $X_{(i)}$ was recorded unless $X_{(i)} > c$. We observe

$$X_{(1)}, ..., X_{(i)}, c+, ..., c+,$$
 or $c+, ..., c+,$

where $X_{(1)} \leq \cdots \leq X_{(i)}$ are order statistics of observed exact values of $X_1, ..., X_n$. 3. Type II censoring.

Observation ceases after d failures.

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

where d is a predetermined number and $X_d = c$.

4. Progressive Type II censoring.

Select *n* samples and determine *d* and $r_1, ..., r_d$, where $\sum_{i=1}^d r_i + d = n$. Observation ceases after *d* failures, at the i-th failure, withdraw r_i experiments randomly, $1 \le i \le d$. Bandom censoring

5. 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



Leukaemia data

Gehan, 1965 recorded times of remission (not worsen) 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+, 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) Assume:

X – survival time, Y — censoring variable.

 $\begin{array}{l} X \text{ and } Y \text{ are independent } (X \perp Y) \text{ i.e., } P(X \leq x, Y \leq y) = P(X \leq x)P(Y \leq y). \\ \text{Observable random vector } (L, R) = \begin{cases} (X, X) & \text{if } X \leq Y \\ (Y, \infty) & \text{if } X > Y. \end{cases} \\ \text{Equivalent observations: } (Z, \delta) \ (= (\min(X, Y), \mathbf{1}_{(X \leq Y)})). \\ \text{Graphical illustration for } Y \text{ and } X \end{array}$



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

(2) This model is not applicable to type II censoring, as well as progressive type II censoring.

§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?

Intuition:

Treating living people as dead ones, does it shorten or extend life expectancy ? Only keeping death data shortens the true life expectancy.

Rigorous Reasoning:

For complete data, $X_1, ..., X_n$ (i.i.d. from X), if $\mu = E(X)$ exists, an estimator of μ is \overline{X} $(=\sum_{i=1}^n X_i/n)$. The properties of \overline{X} : $E(\overline{X}) = \mu$ (unbiased);

 $\overline{X} \to \mu$ with probability one (w.p.1) (strongly consistent); Nonparametric estimators of F and S are

$$\hat{F}(x) = \sum_{i=1}^{n} \mathbf{1}_{(X_i \le x)} / n \ (= \overline{\mathbf{1}_{(X \le x)}}) \quad and \quad \hat{S}(x) = \sum_{i=1}^{n} \mathbf{1}_{(X_i > x)} / n \ (= \overline{\mathbf{1}_{(X > x)}})$$

--Empirical distribution function (edf) and Empirical survival function, respectively. Their properties:

 $\hat{f}(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}_{(X_i=x)} = \overline{\mathbf{1}_{(X=x)}}.$ $\hat{\mu} = \sum_{x} x \hat{f}(x) = \overline{X}.$ $E(\hat{F}(x)) = F(x)$ (unbiased) Why ? (see properties of \overline{X} or $\overline{\mathbf{1}_{(X < x)}}$) $\hat{F}(x) \to F(x)$ w.p.1. (strongly consistent) Why ? Using method 1, we have $m = \sum_{i=1}^{n} \mathbf{1}_{(X_i \leq Y_i)}$ exact observations. So the "sample mean" is

$$\tilde{\mu} = \frac{\sum_{i=1}^{n} X_i \mathbf{1}_{(X_i \le Y_i)}}{\sum_{i=1}^{n} \mathbf{1}_{(X_i \le Y_i)}}? \qquad or = \frac{\sum_{i=1}^{m} X_i^*}{m}?$$

Question: (1) $E(\tilde{\mu}) = \mu$? (2) $\tilde{\mu} \to \mu$ a.s. ?

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

Or

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

Counter-example: Suppose $n = 2, X_i, Y_i$ are i.i.d. from bin(1,1/2)+1. $\mu = E(X) = ?$ $\vdash: E(\tilde{\mu}) \neq \mu.$ $E(\tilde{\mu}) = \begin{cases} \int t f_{\tilde{\mu}}(t) dt & ?\\ \sum_{t} t f_{\tilde{\mu}}(t) & ? \end{cases}$ $\tilde{\mu} = g(X_{1}, X_{2}, Y_{1}, Y_{2}) = \frac{X_{1} \mathbf{1}_{(X_{1} \le Y_{1})} + X_{2} \mathbf{1}_{(X_{2} \le Y_{2})}}{\mathbf{1}_{(X_{1} \le Y_{1})} + \mathbf{1}_{(X_{2} \le Y_{2})}}$ $E(g(\mathbf{X})) = \sum_{\mathbf{X}} g(\mathbf{x}) f_{\mathbf{X}}(\mathbf{x})$

 $f_{X_1,X_2,Y_1,Y_2}(x,y,u,v) = f_{X_1}(x)f_{X_2}(y)f_{Y_1}(u)f_{Y_2}(v)$ by the RC model.

case	X_1	Y_1	X_2	Y_2	$\tilde{\mu} = t = \frac{X_1 1_{(X_1 \le Y_1)} + X_2 1_{(X_2 \le Y_2)}}{1_{(X_1 \le Y_1)} + 1_{(X_2 \le Y_2)}}$	$f_{\tilde{\mu}}(t)$	$t \times f_{\tilde{\mu}}(t)$	
1	1	1	1	1	$\frac{1+1}{1+1} = 1$	1/16	1/16	
2	2	1	2	1	$\frac{0+0}{0+0} = \infty$	1/16	∞	
3	1	2	1	1	0 0	1/16		
4	2	2	2	1		1/16		
•	•	•	•	•		•	•	
16	•	•	•	•		•	•	
sum							∞	$=\sum_{t} t f_{\tilde{\mu}}(t)$

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

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

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

(2) The "edf" is

(1) $\mu = 7/6$

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

W.p.1. its limit is

$$\frac{E(\mathbf{1}_{(X \le Y, X \le x)})}{E(\mathbf{1}_{(X \le Y)})} \ (= \frac{P(X \le 1, X \le x)}{P(X \le 1)}) \tag{Why ?}$$

When x = 1, $\lim_{n \to \infty} \tilde{F}(1) = 1 \neq 1/3 = F(1)$. Thus both estimators are inconsistent.

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. Verify (1) and (2) above.

2. Are the modified edf and sample mean consistent if using the second method ? Justify your statement.

3. Using simulation to check whether the modified edf and sample mean consistent if using the first and second methods.

$\S1.3.$ 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).

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 $\delta = \mathbf{1}_{(X \le Y)}$

Given a sample from the C1 model, is it possible all are right censored ? For example, for a sample of size 2, say (Y_1, δ_1) , (Y_2, δ_2) , can they all be 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;

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

$\S1.5.$ 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 $Y_1 < \cdots < Y_{\mathcal{N}}$. 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

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

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

 $y_1 < y_2 < \cdots < y_k,$

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

This results in inspection times $Y_1 < \cdots < Y_N$, together with their distribution, where N is again a random integer $(N \ge k \text{ or } N \le k ?)$ Define $Y_0 = -\infty$ and $Y_{N+1} = \infty$.

Assume that X and $(Y_1, Y_2, ...)$ 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 for C2 data (mixed case IC model, Schick and Yu (2000)): Assume:

N is a random positive integer;

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

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

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

The observable random vector is

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

Remark .

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

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

3. When Y_i 's are discrete and $N \le k$, then the 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 $Y_1, ..., Y_N \sim f_{\mathbf{Y}}$: find j such that $Y_{j-1} < X \leq Y_j$ to obtain (L, R). repeat 100 times.

Two examples of generating Y_i 's:

(1) Generate $N Z_i \sim exp(\lambda), Y_1 = Z_1, Y_2 = Y_1 + Z_2, ..., Y_N = Y_{N-1} + Z_N.$ (2) generate $Y_j \sim U(2j, 2j+2), j = 1, ..., N$. A simulation example using (2). Assume $N \sim \text{Poisson}(5)+1$, $X \sim exp(3)$ (E(X) = 3). 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) + 1X = rexp(1,1/3) # generate 1 random variable from exp(3) 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 $(X \le Y[1])$ { L[i] = 0R[i] = Y[1]} else { if (X > Y[N]) { L[i] = Y[N]R[i] = 1000} else { j = length(Y[Y < X]) + 1L[i] = Y[j-1]R[i] = Y[j]U=c(L,R) $\dim(U) = c(100,2) \#$ matrix of dimension 100×2 U # print the matrix

Remark. An incorrect approach for dealing with C2 data is to treat the midpoints of IC observations as exact observations.

Question: What is wrong with the following way of generating C2 data?

 $\begin{array}{l} X = rexp(100) \\ L = X-1 \end{array}$

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:

$$\begin{split} N & \text{ is a random integer;} \\ T, Y_1 < Y_2 < \cdots < Y_k < \cdots \text{ are inspection times, } Y_0 = -\infty; \\ X & \text{ and } (N, T, Y_1, \dots, Y_k, \dots) \text{ are independent;} \\ & (\text{ or conditional on } N, X \text{ and } (T, Y_1, \dots, Y_k, \dots) \text{ are independent;} \\ & \mathbf{P}(N=0) > 0 \text{ and } \mathbf{P}(N>1) > 0; \\ & \text{ The observable random vector is} \end{split}$$

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

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> :	LC	SIC	RC	exact
$RC \ data$			+	+
$LC \ data$	+			+
$C1 \ data$	+		+	
$DC \ data$	+		+	+
$C2 \ data$	+	+	+	
\land MIC data	+	+	+	+ /

§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, ..., 8 and $P(N \le 8) = 1$. What will you do if you want to estimate the F(x)? (For example, one may consider the following $\int \frac{L_i + R_i}{2}$ if SIC

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

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

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 singlecensored 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 intervalcensored 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 intervalcensored data. *Scand. J. Statist.*, 27, 45-55.
- * Sun, Jianguo (2006). The statistical analysis of interval-censored failure time data *Spring*.
- * 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

In survival analysis, X is the age at death or failure time.

§2.1. Hazard. Suppose that X is a r.v. with cdf F and density function f.

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

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

where $S(t-) = \lim_{u \uparrow t} S(u)$. S(x) = 1 - F(x).

The d.f. of X,
$$f(t) = \begin{cases} S(t-) - S(t) & \text{if } X \text{ is discrete (why ?)} \\ -S'(t) & \text{if } X \text{ is continuous (why ?)} \end{cases} \begin{cases} F(t) - F(t-) & \dots \\ F'(t) & \dots \end{cases}$$

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

Definition. h(t) = f(t)/S(t-) is called the *hazard function* of X. (wei2xian3 han2su4). $H(t) = -\log S(t)$ is called the *integrated (or cumulative) hazard* of X.

Interpretation of *h* and *H*: $h(x) = \begin{cases} P(X = x | X \ge x) & \text{if } X \text{ is discrete (why ?)} \\ P(X \in [x, x + A) | X \ge x) & \text{if } X \text{ is discrete (why ?)} \end{cases}$

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

If X is continuous, $H'(t) = (-\log S(t))'$ = f(t)/S(t) =h(t) if f(t) exists. That is, $H(t) = \int_{-\infty}^{t} h(x) dx$ integrated hazard.

If X is discrete,

$$S(t) = \prod_{x_i \le t, x_i \in D_f} \frac{S(x_i)}{S(x_i)} \quad (= \frac{S(x_1)}{S(x_1)} \frac{S(x_2)}{S(x_2)} \frac{S(x_3)}{S(x_3)} \cdots \frac{S(t)}{S(t)}) \quad x_1 < x_2 < \cdots$$
$$= \prod_{x_i \le t, x_i \in D_f} (1 - h(x_i))$$

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

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

The distribution of X refers to one of f, h, F, S and H, as they are equivalent if f exists. Recall $F(t) = P(X \le t)$,

$$S(t) = P(X > t),$$

$$f(t) = \begin{cases} -S'(t) & \text{cts} \\ S(t-) - S(t) & \text{dis} \end{cases}$$

$$h(t) = f(t)/S(t-),$$

$$H(t) = -\ln S(t).$$

$\S2.2.$ 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) $(F(x), h(x), \text{ or } H(x)), x \ge 0$ implies that f(x) = 0 etc. for x < 0;

 $S(x), x \ge 0$ implies that S(x) = ? for x < 0.

Note that $E(X) = \int_0^\infty S(t)dt$, if X is a survival time. Reason: If X is a survival time, then

 $E(X) = E(\int_0^X 1dt) = E(\int_0^\infty \mathbf{1}_{(t \le X)} dt)$ $= \int_0^\infty E(\mathbf{1}_{(t \le X)}) dt$ $= \int_0^\infty S(t) dt$ $= \int_0^\infty S(t) dt$ as the set *D* of discrete points of *S*(*t*) is countable and $\int_D dt = 0$. If *X* is not a survival time, do we have $E(X) = \int_0^\infty S(t-) dt$? Try X = -1, $S_X(t) = \mathbf{1}(t < -1)$, $\int_0^\infty S_X(t-) dt = ?$ E(X) = ?

Ans. ?

$\S2.2.1.$ Exponential distribution

$$\begin{split} f(x) &= \rho e^{-\rho x}, \, x > 0, \, \rho > 0; \qquad E(X) = 1/\rho. \\ S(x) &= e^{-\rho x}, \, x > 0; \, (\text{easy way to remember}) \end{split}$$

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

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

$$F(x) = ?$$

§2.2.2. Weibull distribution

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

$$f(x) = \kappa \rho(\rho x)^{\kappa-1} e^{-(\rho x)^{\kappa}}, x > 0, \rho, \kappa > 0;$$

$$h(x) = \kappa \rho(\rho x)^{\kappa-1}, x > 0;$$

$$H(x) = (\rho x)^{\kappa}, x > 0;$$

$$F(x) = ?$$

§2.2.3 Log normal distribution

X has a log normal distribution if $X = e^{Z}$, where $Z \sim N(\nu, \tau^2)$ (= $N(\log \rho^{-1}, \tau^2)$), $\rho, \tau > 0$,

$$f(x) = f_Z(\log x) \left| \frac{dz}{dx} \right| \qquad z = \ln x$$

= $\frac{1}{\sqrt{2\pi\tau}} \exp(\frac{-(\log x - \nu)^2}{2\tau^2}) \frac{1}{x} \qquad \nu = -\log \rho$
= $\frac{1}{\sqrt{2\pi\tau}x} \exp(\frac{-(\log x - (-\log \rho))^2}{2\tau^2})$
= $\frac{1}{\sqrt{2\pi\tau}x} \exp(\frac{-(\log(x\rho))^2}{2\tau^2}), \quad x > 0,$

 $S(t) = \int_t^{\infty} f(x) dx, \ h = f/S, \ H = -\ln S \text{ and } F(x) = ?$ They cannot be simplified. Just remember $\ln X \sim N(\nu, \tau^2)$. §2.2.4 Log-logistic distribution

X has the log-logistic distribution, if $X = e^T$, where $T = \ln X$ has a logistic distribution. The standard logistic survival function is $S_o(x) = \frac{1}{1 + \exp(x)}$, x > 0? It induces the location-scale parameter family

$$\begin{split} S_T(x) &= \frac{1}{1 + \exp(\frac{x-\nu}{\tau})}.\\ S(x) &= S_X(x)\\ &= P\{X > x\}\\ &= P\{T > \ln x\}\\ &= \frac{1}{1 + \exp(\frac{\ln x - \nu}{\tau})} \qquad \text{letting } \nu = -\log \rho \text{ and } \kappa = 1/\tau \text{ (reparameterization)},\\ &= \frac{1}{1 + \exp((\ln x + \ln \rho)\kappa)}\\ &= \frac{1}{1 + \exp((\ln (x\rho))\kappa)}\\ &= \frac{1}{1 + \exp((\ln (x\rho)))^{\kappa}}\\ \text{then } S(x) &= \frac{1}{1 + (\rho x)^{\kappa}}, \ x \ge 0 \text{ (easy way to remember)}.\\ H(x) &= \log(1 + (\rho x)^{\kappa}), \ x \ge 0. \qquad = -\ln S(x) \end{split}$$

$$\begin{split} h(x) &= \frac{\kappa \rho^{\kappa} x^{\kappa - 1}}{1 + (\rho x)^{\kappa}}, \ x \ge 0. \\ f(x) &= \frac{\kappa \rho^{\kappa} x^{\kappa - 1}}{[1 + (\rho x)^{\kappa}]^2}, \ x \ge 0. \\ F(x) &= ? \end{split} = H'(x) \end{split}$$

\S **2.2.5.** Gompertz-Makeham distribution.

$$h(t) = \rho_0 + \rho_1 e^{\rho_2 t}, \ t > 0, \ \rho_0, \rho_1 > 0 \ (\text{easy way to remember}),$$

$$H(t) = \int_0^t h(x) dx = [\rho_0 t + \frac{\rho_1}{\rho_2} (e^{\rho_2 t} - 1)], \ t > 0;$$

$$S(t) = e^{-H(t)} = e^{-[\rho_0 t + \frac{\rho_1}{\rho_2} (e^{\rho_2 t} - 1)]}, \ t > 0;$$

$$f(t) = -S'(t) = h(t)S(t) = \cdots, \ t > 0;$$

$$F = ?$$

Remark:

1. It is an exponential distribution if $\rho_1 = 0$.

2. It is called Gompertz distribution if $\rho_0 = 0$.

\S **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

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

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

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

then

$$\begin{split} f(x) &= \int_0^\infty p e^{-px} \frac{p^{\alpha-1} e^{-p/\beta}}{\Gamma(\alpha)(\beta)^{\alpha}} dp = \int_0^\infty \frac{p^{\alpha} e^{-p(x+1/\beta)}}{\Gamma(\alpha)(\beta)^{\alpha}} dp \\ &= \frac{\Gamma(\alpha+1)(\frac{1}{x+1/\beta})^{\alpha+1}}{\Gamma(\alpha)(\beta)^{\alpha}} \int_0^\infty \frac{p^{(\alpha+1)-1} e^{-p/\frac{1}{(x+1/\beta)}}}{\Gamma(\alpha+1)(\frac{1}{x+1/\beta})^{\alpha+1}} dp \\ &= \frac{\Gamma(\alpha+1)}{(x+1/\beta)^{\alpha+1}\Gamma(\alpha)(\beta)^{\alpha}} \\ &= \alpha(1/\beta)^{\alpha} (x+1/\beta)^{-\alpha-1} \dots? \end{split}$$

The latter df is called the **Pareto distribution**, with $\rho = \alpha\beta$ (the mean),

$$S(x) = (\alpha/\rho)^{\alpha} (x + \alpha/\rho)^{-\alpha}, \ x > 0.$$

$$H(x) = \alpha(-\ln(\alpha/\rho) + \ln(x + \alpha/\rho)), \ x > 0.$$

$$h(x) = \alpha/(x + \alpha/\rho), \ x > 0.$$
 (simplest).

\S **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. $X \sim bin(2, p)$. $f(x) = P(X - x) - {2 \choose 2} n^x a^{2-x} \quad x \in \{0, 1, 2\}, a = 1, \dots, n$

$$\begin{aligned} f(x) &= P(X = x) = \binom{x}{x} p^{x} q^{2-x}, & x \in \{0, 1, 2\}, & q = 1 - p. \\ S(x) &= P(X > x) = \begin{cases} 1 & \text{if } x < 0 \\ 1 - q^{2} & \text{if } x \in [0, 1), \text{ why } ? \\ p^{2} & \text{if } x \in [1, 2), \\ 0 & \text{if } x \ge 2. \end{cases} \end{aligned}$$

$$\text{For } h(x) = f(x)/S(x-),$$

$$\text{need } S(0-) = 1, & S(1-) = 1 - q^{2}, & S(2-) = p^{2}, \end{cases}$$

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

\S 2.2.8. Proportional hazards (PH) model or Cox's 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 $\mathbf{X} \in \mathcal{R}^p$. We say (\mathbf{X}, Y) follows a proportional hazards (PH) model or Cox's regression model (Cox, 1972), if given $\mathbf{X} = \mathbf{x}$, the hazard of $Y|\mathbf{x}$ is

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

where $c(\mathbf{x}) \ge 0$, $c(\cdot)$ takes at least two distinct values, and h_o is a hazard function.

Remark 1. It is common to set $c(\mathbf{x}) = exp(\beta \mathbf{x})$ so that $c(\mathbf{x}) \ge 0$, 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. *e.g.*, if $h_o(y) = 0.5$, $c(\mathbf{x}) = 3$, $P(Y = y|Y \ge y, \mathbf{X} = \mathbf{x}) = h(y|\mathbf{x}) = h_o(y)c(\mathbf{x}) > 1$. Two alternatives: either

(1) choose $c(\mathbf{x}) = exp(-e^{\beta \mathbf{X}})$ to ensure that $c(\mathbf{x})$ is between 0 and 1 ($e^0 = 1, e^{-\infty} = 0$),

(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 no restriction $y < \tau$. We shall show in Example 1 that if h_o corresponds to a discrete random variable, Eq. (1) without the restriction does not define a hazard function.

Hereafter, let 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

$$h_T(\tau) = f_T(\tau)/S_T(\tau) = f_T(\tau)/P(T \ge \tau) = f_T(\tau)/f_T(\tau) = 1$$

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

$$h(\tau|x) = 1 \neq 1 \times c(x) = h_o(\tau)c(x)$$
 as $c(x)$ takes at least two values.

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

$$S(t|x) = (S_o(t))^{e^{\beta x}} \text{ satisfies the PH model } h(t|x) = e^{\beta x} h_o(t), \ t < \tau.$$

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 of Eq. (2):

- 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.

$$h(y|\mathbf{x}) = h_o(y)c(\mathbf{x})$$
 for $y < \tau$, satisfies the PH model. (1)

$$S(t|x) = (S_o(t))^{e^{\beta x}}$$
(2)

satisfies the PH model if X is continuous.

Definition. The family of the survival functions $S(t|\mathbf{x})$ satisfies Eq. (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

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

The PH model and the PIH model are the same for cts random variables, but are different for discrete random variables. 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 the textook by Sun (2006) p.18). Example 3 shows that they are different.

Example 3. Suppose $Y_0 \sim bin(2, p)$. 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

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

$$f(2|x) = 1 - f(0|x) - f(1|x).$$

Verify that $f(\cdot|\mathbf{x})$ defines a discrete density function (for $c(x) \leq 1$. Why add it ??) It follows $h(y|x) = h_o(y)c(x), y \in \{0, 1\}$, that

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

However, if p = 0.2 and c(x) = 0.3, PIH model:

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

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).

\S **2.2.9.** Accelerated lifetime (AL) model.

Definition. If S_o is a survival function and given a p dimensional vector $\mathbf{V} = \mathbf{v}$, $X | \mathbf{v}$ has a survival function $S(y|\mathbf{v}) = S_o(y/\exp(\beta \mathbf{v}))$, $\beta \in \mathcal{R}^p$, then we say $X | \mathbf{v}$ is from an accelerated lifetime model (or the log-linear regression model). That is,

$$X/e^{\beta \mathbf{V}} = W$$
 and $\ln X = \beta \mathbf{v} + \ln W$, where $S_W = S_o$ (and $\alpha = E \ln W$).

Examples.

Weibull: $S(y|\mathbf{x}) = \exp(-y^{\kappa}e^{\alpha\mathbf{X}}), y > 0, \alpha = -\beta\kappa. \ (S_o(y) = exp(-(\rho y)^{\gamma})).$ (hw). Log-logistic: $S(y|\mathbf{x}) = \frac{1}{1+(y\exp(\alpha\mathbf{X}))^{\kappa}}. \ \alpha = -\beta. \ (S_o(y) = \frac{1}{1+(\rho y)^{\kappa}}).$

$\S2.3.$ Homework:

- A.1. Let $X_1, ..., X_n$ be independent continuous nonnegative random variables with hazard functions $h_1(\cdot), ..., h_n(\cdot)$. Prove that $X = \min\{X_1, ..., X_n\}$ has hazard function $\sum_{j=1} h_j(t)$.
- 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 $Var(X) = 2E(1/P^2) [E(1/P)]^2$.
- 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,

 $(L_i, R_i, X_i), i = 1, ..., n$, are i.i.d. copies of (L, R, X).

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

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

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

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

Note: θ (*e.g.* $\theta = 1$) is a parameter and ϕ is a variable in Θ , as we do not know θ . In particular,

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

as $P\{X \in [x, x + \Delta)\} \approx f(x)\Delta$. The MLE of θ maximizes $l(\phi), \phi \in \Theta$. Abusing notation, we may write $l(\theta)$, instead of $l(\phi)$. 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 $l(\phi)$ the likelihood function of the IC data $(L_i, R_i) = (l_i, r_i), i = 1, ..., n$, where

$$l(\theta) = \prod_{i=1}^{n} P((L_i, R_i) = (l_i, r_i)), \quad \theta \text{ is the true value.}$$

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

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

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 true value θ in the above expression for simplicity,

as $f = f_o(x; \theta)$, $S = S_o(x; \theta)$ and $F = F_o(x; \theta)$.

$$\begin{split} l(\phi) &= \prod_{i=1}^{n} \left\{ [f(l_i)S_Y(l_i-)]^{\mathbf{1}_{e,i}} [S(l_i)f_Y(l_i)]^{\mathbf{1}_{r,i}} \right\} & \text{ignoring } \phi \\ &= \left\{ \prod_{i=1}^{n} \{ [f_o(l_i;\phi)]^{\mathbf{1}_{e,i}} [S_o(l_i;\phi)]^{\mathbf{1}_{r,i}} \} \right\} \left\{ \prod_{i=1}^{n} \underbrace{\{ [(S_Y(l_i-)]^{\mathbf{1}_{e,i}} [f_Y(l_i)]^{\mathbf{1}_{r,i}} \}}_{why \text{ no } \phi?} \right\}. \end{split}$$

Example 2 (DC model).

$$\begin{split} P((L,R) = (l,r)) &= \begin{cases} P(X = l, Z < l \le Y) & \text{if exact} \\ P(X > l = Y) & \text{if } RC \\ P(X < r = Z) & \text{if } LC \end{cases} = \begin{cases} f(l)P(Z < l \le Y) & \text{if exact} \\ S(l)f_Y(l) & \text{if } RC \\ F(r)f_Z(r) & \text{if } LC. \end{cases} \\ \text{Note } \{Z \ge l\} \subset \{Y \ge l\} \text{ as } Z \le Y \text{ w.p.1.} \\ \text{Thus } P(Z < l \le Y) = P(Y \ge l) - P(Z \ge l) = S_Y(l-) - S_Z(l-), \\ axix & y & \uparrow & z < l < y & | & z > l & \ddots \\ & y > l & \ddots & \\ & 0 & | & y > l & \ddots \\ & 0 & | & 1 & 0 \end{cases} \text{ not part of sample space} \\ l & ---- & i & --- & --- & --- & ---- & z \text{ axis} \\ P((L,R) = (l,r)) = [f(l)(S_Y(l-) - S_Z(l-))]^{\mathbf{1}_e}[S(l)f_Y(l)]^{\mathbf{1}_r}[F(r)f_Z(r)]^{\mathbf{1}_l}. \\ \text{Thus } l(\phi) \\ = \prod_{i=1}^n \{[f(l_i)(S_Y(l_i-) - S_Z(l_i-))]^{\mathbf{1}_{e,i}}[S(l_i)f_Y(l_i)]^{\mathbf{1}_{r,i}}[F(r_i)f_Z(r_i)]^{\mathbf{1}_{r,i}} \} \\ = \{\prod_{i=1}^n ([f_o(l_i;\phi)]^{\mathbf{1}_{e,i}}[S_o(l_i;\phi)]^{\mathbf{1}_{r,i}}[F_o(r_i;\phi)]^{\mathbf{1}_{i,i}})\} \prod_{i=1}^n ([S_Y(l_i-) - S_Z(l_i-)]^{\mathbf{1}_{e,i}}[f_Y(l_i)]^{\mathbf{1}_{r,i}}[f_Z(r_i)]^{\mathbf{1}_{r,i}}). \\ \text{Example 3 (Case k model, $k \ge 2$). \\ P((L,R) = (l,r)) = \begin{cases} P\{l < X \le r, (Y_{j-1}, Y_j) = (l,r), j \in \{2, ..., k\}\} & \text{if SIC} \\ S(l)f_{Y_1}(l) & \text{if IC} \\ F(r)f_{Y_1}(r) & \text{if it is SIC} \\ S(l)f_{Y_1}(l) & \text{if it is right censored} \\ F(r)f_{Y_1}(r) & \text{if it is left censored} \\ F(r)f_{Y_1}(r) & \text{if it is left censored} \end{cases}$$

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

Thus $l(\phi)$

=

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

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.

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

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

$$\mathcal{L}(\phi) = \Big\{ \prod_{i=1}^{n} \big\{ [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}} \Big\},\$$

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

In an obvious way, we write

$$\begin{aligned} \mathbf{L}(\phi) &= \prod_{i:\ ex} f(l_i) \prod_{i:\ rc} S(l_i) \prod_{i:\ lc} F(r_i) \prod_{i:\ ic} (S(l_i) - S(r_i)) \\ &= \prod_{i=1}^{n} [(f_o(l_i;\phi))^{\delta_i} (S_o(l_i;\phi) - S_o(r_i;\phi))^{1-\delta_i}] \\ &(= \prod_{i=1}^{n} [(P(X = l_i))^{\delta_i} (P(l_i < X \le r_i))^{1-\delta_i}] \end{aligned}$$
if X is dicrete),

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

Remark. Hereafter ϕ can be replaced by θ , ρ etc..

§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 $L = \prod_{i=1}^n (S(x_i - 1) - S(x_i + 1))$? (prove or disprove it).

§3.3. MLE.

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

Example 1. Suppose that X has an exponential distribution,

- 1. f(t) = ? (1) e^{-t} Y,N,DNK; (2) $e^{-t}\mathbf{1}(t \ge 0)$ Y,N,DNK.
- 2. S(t) = ? (1) e^{-t} Y,N,DNK; (2) $e^{-t}\mathbf{1}(t \ge 0)$ Y,N,DNK; (3) $e^{-t\mathbf{1}(t\ge 0)}$ Y,N,DNK.
- 3. Find the MLE of ρ under the RC model.

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

$$\begin{split} \log \mathcal{L}(\rho) &= \sum_{i=1}^{n} \log\{[\mathbf{1}_{(Z_{i}>0)}\rho e^{-\rho Z_{i}}]^{\delta_{i}}[e^{-\rho Z_{i}}\mathbf{1}_{(Z_{i}>0)}]^{1-\delta_{i}}\} \quad (=\log(\prod_{i:\ ex}f(l_{i})\prod_{i:\ rc}S(l_{i})))\\ &= \sum_{i=1}^{n} \log\{[\rho e^{-\rho Z_{i}}]^{\delta_{i}}[e^{-\rho Z_{i}}]^{1-\delta_{i}}\} \qquad \qquad Z_{(1)}>0 \end{split}$$

$$=\sum_{i=1}^{n} \log\{\rho^{\delta_i} e^{-\rho Z_i}\} \qquad \qquad Z_{(1)} > 0$$

$$=\sum_{i=1}^{n} \delta_i \log \rho - \rho \sum_{i=1}^{n} Z_i.$$
 $Z_{(1)} > 0$

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 L, but it is not the MLE, as $0 \notin (0, \infty) = \Theta$. We define $\hat{\rho} = 0$ so that $\hat{\rho}$ is well defined.

Remark. Formally, we shall write L(p) instead of $L(\rho)$ why ? 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 $\frac{n}{\sum_{i=1}^n X_i} = h(\overline{X})$ 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,

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

2. By the central limit theorem and Corollary of the Slutsky's Theorem $(X_n + Y_n \xrightarrow{\mathcal{D}} X + c)$ and $X_n Y_n \xrightarrow{\mathcal{D}} X * c$ if $X_n \xrightarrow{\mathcal{D}} X$ and $Y_n \xrightarrow{\mathcal{D}} c$ (a constant)), if σ_X is finite, h' is continuous at μ_X and $h'(\mu_X) \neq 0$,

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

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

$$\frac{\sqrt{n}(h(\overline{X}) - h(\mu_X))}{\sqrt{(\frac{\partial h(\mu_X)}{\partial \mu})^t \Sigma \frac{\partial h(\mu_X)}{\partial \mu}}}$$
 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 common regression model s:

- (1) the Cox regression model $S(t|\mathbf{v}) = S_{X|\mathbf{v}}(t|\mathbf{v}) = (S_o(t))^{\exp(\mathbf{v}\beta)}$;
- (2) the accelerated lifetime model $\ln X = \beta \mathbf{V} + \ln W$ or $W = X/e^{\beta \mathbf{V}}$;
- (3) the linear regression model $X = \beta \mathbf{V} + W$ or $W = X \mathbf{V}\beta$,

where X is the survival time, **V** is a $p \times 1$ covariate vector, $\beta \mathbf{v} \stackrel{def}{=} \beta' \mathbf{v}$, $\mathbf{v} = (v_1, ..., v_p)'$. and $S_o(t) = S(t|0) = S_W$ is the baseline survival function.

Definition. The likelihood function of the IC regression data (l_i, r_i, \mathbf{v}_i) 's is defined to be

$$\mathbf{L} = \begin{cases} \prod_{i=1}^{n} (f_W(l_i - \beta \mathbf{v}_i))^{\delta_i} (S_W(l_i - \mathbf{v}_i\beta) - S_W(r_i - \mathbf{v}_i\beta))^{1 - \delta_i} & \text{LR model} \\ \prod_{i=1}^{n} (f_W(l_i/e^{\mathbf{v}_i\beta}))^{\delta_i} (S_W(l_i/e^{\mathbf{v}_i\beta}) - S_W(r_i/e^{\mathbf{v}_i\beta}))^{1 - \delta_i} & \text{AL model} \\ \prod_{i=1}^{n} (-d(S_W(l_i))^{\exp(\beta \mathbf{v}_i)})^{\delta_i} (S_W(l_i - q_i))^{\exp(\beta \mathbf{v}_i)} - S_W(r_i/e^{\mathbf{v}_i\beta}))^{1 - \delta_i} & \text{AL model} \end{cases}$$

$$\left(\prod_{i=1}^{n}\left(-\frac{d(S_W(l_i))^{\exp(\beta\mathbf{v}_i)}}{dl_i}\right)^{\delta_i}\left((S_W(l_i))^{\exp(\beta\mathbf{v}_i)}-(S_W(r_i))^{\exp(\beta\mathbf{v}_i)}\right)^{1-\delta_i}\quad\text{Cox model}$$

Example 2. Suppose that (Z_i, δ_i, V_i) , i = 1, ..., n are i.i.d. RC regression data from (Z, δ, V) , where $Z = X \wedge Y$, $\delta = \mathbf{1}(X \leq Y)$, $X = \beta V + W$, W, Y, V are independent random variables, $W \sim Exp(1)$ and $V = \pm 1$. Derive the MLE of β if n = 3, $\delta_1 = 1$, $\delta_2 = 0 = \delta_3$, $Z_2 = Z_3$, $V_1 = 1$ and $V_2 = -1 = V_3$.

In the real case, Z_i 's are given. Here is for exercises. Sol. Observations: $Z_1 = X_1 = W_1 + \beta$, $Z_2 = Y_2 < X_2 = W_2 - \beta$ and $Z_3 = \cdots$, $\delta_1 = \cdots$

In real data, either $Z_1 < -Z_2$ or $Z_1 > Z_2$. Only one case. To find the MLE of β , (1) $\frac{d}{d\beta}L(\beta) = 0 \dots$?? (2) $\frac{d \log L}{d\beta} = 0 \dots$?? Or (3) otherwise. Which way ??

$$\begin{split} &\text{If } Z_1 < -Z_2, \, \text{then } L(\beta) = \begin{cases} e^{-Z_1 + \beta} & \text{if } \beta \leq Z_1 \\ 0 & \text{if } \beta > Z_1, \end{cases} \quad L'(\beta) = \begin{cases} e^{-Z_1 + \beta} > 0 & \text{if } \beta \leq Z_1 \\ 0 & \text{if } \beta > Z_1, \end{cases} \quad \text{Conclusion ?} \\ &L(\beta) = \begin{cases} e^{-Z_1 + \beta} (\uparrow in \beta) & \text{if } \beta \leq Z_1 \\ 0 & \text{if } \beta > Z_1, \end{cases} = > \text{MLE } \hat{\beta} = Z_1. \\ &\text{If } Z_1 \geq -Z_2, \, \text{then } L(\beta) = \begin{cases} e^{-Z_1 + \beta} (\uparrow in \beta) & \text{if } \beta \leq Z_1 \\ e^{-Z_1 - 2Z_2 - \beta} & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} & \text{if } \beta \in (-Z_2, Z_1] \\ 0 & \text{if } \beta > Z_1 \end{cases} \\ &L'(\beta) = \begin{cases} e^{-Z_1 + \beta} & \text{if } \beta < -Z_2 \\ -e^{-Z_1 - 2Z_2 - \beta} & \text{if } \beta \in (-Z_2, Z_1] \\ 0 & \text{if } \beta > Z_1 \end{cases} \\ &L(\beta) = \begin{cases} e^{-Z_1 + \beta} & \text{if } \beta \leq -Z_2 \\ -e^{-Z_1 - 2Z_2 - \beta} & \text{if } \beta \in (-Z_2, Z_1] \\ 0 & \text{if } \beta > Z_1 \end{cases} \\ &L(\beta) = \begin{cases} e^{-Z_1 + \beta} (\uparrow) & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} & \text{if } \beta \in (-Z_2, Z_1] \\ 0 & \text{if } \beta > Z_1 \end{cases} \\ &L(\beta) = \begin{cases} e^{-Z_1 + \beta} (\uparrow) & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in (-Z_2, Z_1] \\ 0 & \text{if } \beta > Z_1 \end{cases} \\ &L(\beta) = \begin{cases} e^{-Z_1 + \beta} (\uparrow) & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in (-Z_2, Z_1] \\ 0 & \text{if } \beta > Z_1 \end{cases} \\ &L(\beta) = \begin{cases} e^{-Z_1 + \beta} (\uparrow) & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in (-Z_2, Z_1] \\ 0 & \text{if } \beta > Z_1 \end{cases} \\ &L(\beta) = \begin{cases} e^{-Z_1 + \beta} (\uparrow) & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in (-Z_2, Z_1] \\ 0 & \text{if } \beta > Z_1 \end{cases} \\ &L(\beta) = \begin{cases} e^{-Z_1 + \beta} (\uparrow) & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in (-Z_2, Z_1] \end{cases} \\ &= \begin{cases} e^{-Z_1 + \beta} (\uparrow) & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in (-Z_2, Z_1] \end{cases} \\ &= \begin{cases} e^{-Z_1 + \beta} (\uparrow) & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in (-Z_2, Z_1] \end{cases} \\ &= \begin{cases} e^{-Z_1 + \beta} (\uparrow) & \text{if } \beta \leq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \geq -Z_2 \end{cases} \end{cases} \\ &= \begin{cases} e^{-Z_1 + \beta} (\uparrow) & e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \geq -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in -Z_2 \\ e^{-Z_1 - 2Z_2 - \beta} (\downarrow) & \text{if } \beta \in -Z_2 \end{cases} \end{cases}$$

Remark. Both answers are correct. In the survival analysis context, $Z_i \ge 0$. Thus $Z_1 \ge -Z_2$ and the MLE is $\hat{\beta} = -Z_2$. §3.3.2. Homework:

1. Show that $\hat{\rho} = \overline{\delta}/\overline{Z}$ in Example 1 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

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

- 3. Give a 99% approximate confidence interval for ρ when n = 100. Try to give a 99% confidence interval for ρ when n = 1 with $\rho \ge 1/\ln 50$ and $Y \equiv c$.
- 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 ?
- 7. Suppose that (Z_i, δ_i, V_i) , i = 1, ..., n are i.i.d. RC regression data from (Z, δ, V) , where $Z = X \wedge Y$, $\delta = \mathbf{1}(X \leq Y)$, $X = \beta V + W$, W, Y, V are independent random variables, $W \sim Exp(1)$ and $V = \pm 1$. Derive the MLE of β if n = 2, $\delta_1 = 1$, $\delta_2 = 0$, $V_1 = 1$ and $V_2 = -1$.

\S **3.4.** Numerical methods for MLE.

3.4.1. Newton-Raphson method.

If $L(\phi)$ 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

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

This is an iterative algorithm:

Step 1. Assign an initial value $\phi_{(1)}$ to $\phi \ (\in \mathbb{R}^m)$. Step k + 1, $k \ge 1$. Given $\phi_{(k)}$, up-date ϕ by

$$\phi_{(k+1)} = \phi_{(k)} - \left(\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t}\Big|_{\phi = \phi_{(k)}}\right)^{-1} \frac{\partial \mathcal{L}}{\partial \phi}\Big|_{\phi = \phi_{(k)}}.$$
(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 $\left(\frac{\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t}\Big|_{\phi=\phi_{(k)}}\right)^{-1}$. 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 $U_i = \ln X_i$ be a $N(\mathbf{z}_i^t \beta, \sigma^2)$ variate, i = 1, ..., n, where \mathbf{z}_i is a $p \times 1$ covariate vector (non-random vectors), 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, ..., 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 $U_1, ..., U_n$ are observed,

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

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

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

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

$$\mathcal{L}(\theta) = \ln \prod_{i=1}^{n} (S(l_i - \mathbf{z}_i \beta; \sigma) - S(r_i - \mathbf{z}_i \beta; \sigma)) \qquad (\epsilon_i = U_i - \mathbf{z}_i \beta)$$
$$= \ln \prod_{i=1}^{n} (F(r_i - \mathbf{z}_i \beta; \sigma) - F(l_i - \mathbf{z}_i \beta; \sigma))$$
$$= \sum_{i=1}^{n} \ln(\Phi(\frac{r_i - \beta' \mathbf{z}_i}{\sigma}) - \Phi(\frac{l_i - \beta' \mathbf{z}_i}{\sigma})) \qquad (\text{if treating as i.i.d.}),$$

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

$$\frac{d}{dx}F(b(x)) = \left(\int_{-\infty}^{b(x)} f(t)dt\right)'_{x} = f(b(x))b'(x),$$

$$\frac{\partial\mathcal{L}(\theta)}{\partial\sigma} = \sum_{i=1}^{n} \frac{\frac{\partial}{\partial\sigma}\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \frac{\partial}{\partial\sigma}\Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma})}{\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma})} \tag{lng(x)}'_{x} = \frac{g'(x)}{g(x)}$$

$$= -\frac{\sigma^{-2}}{\sqrt{2\pi}} \sum_{i=1}^{n} \frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} (r_i - \beta^t \mathbf{z}_i) - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} (l_i - \beta^t \mathbf{z}_i)}{\Phi(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}) - \Phi(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma})}$$
$$= -\frac{\sigma^{-1}}{\sqrt{2\pi}} \sum_{i=1}^{n} \frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} \frac{(r_i - \beta^t \mathbf{Z}_i)}{\sigma} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} \frac{(l_i - \beta^t \mathbf{Z}_i)}{\sigma}}{\sigma}}{\Phi(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}) - \Phi(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma})} \quad formally. \quad (3.2)$$

Recall that $(l,r) \in \{(-\infty,r), (l,r), (l,\infty)\}$. If $l_i = -\infty$, then $\Phi(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}) = 0$ and thus $\Phi'_{\sigma}(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}) = 0$. But we can nevertheless write $\Phi'_{\sigma}(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}) = -\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{\sigma^2}e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}$, 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 as in Eq. (3.2).

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

Here

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

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 $\begin{pmatrix} \beta \\ \sigma \end{pmatrix}$:

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

One can pretend β and \mathbf{z}_i as $\in \mathcal{R}$ first, and adjust later on. $\beta z_i \rightarrow \beta^t \mathbf{z}_i, z_i^2 \rightarrow \mathbf{z}_i \mathbf{z}_i^t, etc.$ Their dimensions ?

$$\begin{split} &\frac{\partial^{2}\mathcal{L}(\theta)}{\partial\sigma\partial\beta} = \frac{\partial}{\partial\sigma} \Big[-\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(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma} - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}} \mathbf{z}_{i} \right) \Big] \\ &= -\left[\frac{\partial}{\partial\sigma} \frac{\sigma^{-1}}{\sqrt{2\pi}} \right] \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(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i})} - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma}} \mathbf{z}_{i} \right) - \frac{\sigma^{-1}}{\sqrt{2\pi}} \frac{\partial}{\partial\sigma} \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(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}} - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{2\sigma^{2}}} \mathbf{z}_{i} \right) \\ &= -\left[\frac{\partial}{\partial\sigma} \frac{\sigma^{-1}}{\sqrt{2\pi}} \right] \sum_{i=1}^{n} \left(\frac{e^{-\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma}} - e^{-\frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{2\sigma^{2}}}}{\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i})} - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma})} \mathbf{z}_{i} \right) - \frac{\sigma^{-1}}{\sqrt{2\pi}} \sum_{i=1}^{n} \frac{\mathbf{z}_{i}}{\left[\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) \right]^{2}} \left\{ \left(\frac{\partial}{\partial\sigma} \left[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}}} \right] - \frac{\partial}{\partial\sigma} \left[\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) \right] \right) \right\} \end{split}$$

$$= \frac{\sigma^{-2}}{\sqrt{2\pi}} \sum_{i=1}^{n} \mathbf{z}_{i} \Big\{ \frac{e^{-\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{2\sigma^{2}}} - e^{-\frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{2\sigma^{2}}}}{\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma})} - \frac{\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma^{2}} e^{-\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{2\sigma^{2}}}}{\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma})} - \frac{\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma^{2}} e^{-\frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma^{2}}}}{\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma})} - \frac{\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma^{2}}}{(\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}))^{2}} \cdot \frac{e^{-\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{2\sigma^{2}}} - e^{-\frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{2\sigma^{2}}}}{\sqrt{2\pi}} \Big\}.$$

$$\frac{\partial^{2} \mathcal{L}(\theta)}{\partial \beta \partial \beta^{t}} = \frac{\partial}{\partial \beta} \left[-\frac{\sigma^{-1}}{\sqrt{2\pi}} \sum_{i=1}^{n} \left(\frac{e^{-\frac{(r_{i} - \beta^{t} \mathbf{Z}_{i})^{2}}{2\sigma^{2}}} - e^{-\frac{(l_{i} - \beta^{t} \mathbf{Z}_{i})^{2}}{2\sigma^{2}}}{\Phi(\frac{r_{i} - \beta^{t} \mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i} - \beta^{t} \mathbf{Z}_{i}}{\sigma})} \mathbf{z}_{i}^{t} \right) \right] \qquad \qquad = \frac{\partial^{2} \mathcal{L}(\theta)}{\partial \beta^{t} \partial \beta} ?$$

$$= -\sum_{i=1}^{n} \left[\frac{\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma^2} e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - \frac{l_i - \beta^t \mathbf{Z}_i}{\sigma^2} e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{(\Phi(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}) - \Phi(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}))\sqrt{2\pi\sigma^2}} + \left(\frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}) - \Phi(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma})} \right)^2 \frac{1}{2\pi\sigma^2} \right] \mathbf{z}_i \mathbf{z}_i^t$$
$$= -\sum_{i=1}^{n} \left[\frac{\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma} e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - \frac{l_i - \beta^t \mathbf{Z}_i}{\sigma} e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{(\Phi(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}) - \Phi(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma}))\sqrt{2\pi}} + \left(\frac{e^{-\frac{(r_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}} - e^{-\frac{(l_i - \beta^t \mathbf{Z}_i)^2}{2\sigma^2}}}{\Phi(\frac{r_i - \beta^t \mathbf{Z}_i}{\sigma}) - \Phi(\frac{l_i - \beta^t \mathbf{Z}_i}{\sigma})} \right)^2 \frac{1}{2\pi} \right] \frac{\mathbf{z}_i \mathbf{z}_i^t}{\sigma^2}.$$

$$\begin{split} \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}}}{\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma})} \\ &-\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}}} \\ &-(\frac{\sigma^{-2}}{\sqrt{2\pi}})^{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}}}{\Phi(\frac{r_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma}) - \Phi(\frac{l_{i}-\beta^{t}\mathbf{Z}_{i}}{\sigma})} \\ &-(\frac{\sigma^{-2}}{\sqrt{2\pi}})^{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}}} \\ &=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}}} \\ &-\frac{\sigma^{-2}}{\sqrt{2\pi}}\sum_{i=1}^{n}\frac{\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{3}}{\sigma}e^{-\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{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}}} \\ &-\frac{\sigma^{-2}}{\sqrt{2\pi}}\sum_{i=1}^{n}\frac{\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{3}}{\sigma}e^{-\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{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}}} \\ &-\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}}{\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}}} \\ &-\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}}{\sigma^{2}}} - \frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})}{\sigma}e^{-\frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma^{2}}}} \\ &-\frac{\sigma^{-2}}{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}}{\sigma^{2}}} - \frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})}{\sigma}e^{-\frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma^{2}}}} \\ &-\frac{\sigma^{-2}}{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}}{\sigma^{2}}} - \frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})}{\sigma}e^{-\frac{(l_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}{\sigma^{2}}} \\ &-\frac{\sigma^{-2}}{2\pi}\sum_{i=1}^{n}\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})}{\sigma^{2}}e^{-\frac{(r_{i}-\beta^{t}\mathbf{Z}_{i})^{2}}} - \frac{(l_{i}-$$

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 $\theta \in \Theta \subset \mathbb{R}^m$.

fidence 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}_{(-2\ln \frac{\mathbf{L}(\tilde{\theta}_0)}{\mathbf{L}(\hat{\theta})} > \chi^2_{\alpha,d})} where \hat{\theta} and \tilde{\theta}_0 are the MLEs of \theta in the
      space \Theta and \Theta_0, respectively, d = ||\Theta|| - ||\Theta_0|| and ||\Theta_0|| is the dimension of \Theta_0;
      The score test is based on \frac{\partial \mathcal{L}}{\partial \phi}\Big|_{\phi=\tilde{\theta}_0}. In particular, if \Theta_o = \{\theta_o\}, then the test is \mathbf{1}_{(T>\chi^2_{\alpha,d})}, where T = \left(\frac{\partial \mathcal{L}}{\partial \phi}^t (\frac{-\partial^2 \mathcal{L}}{\partial \phi \partial \phi^t})^{-1} \frac{\partial \mathcal{L}}{\partial \phi}\right)\Big|_{\phi=\theta_o}.
3.4.2. Monte-Carlo method.
       Generate a large number N of potential values of \theta, say t_1, ..., t_N.
Let \hat{\theta} = \arg \max_{t_i} \{ \mathcal{L}(t_i) \}.
Remark. The drawback:
      1. Time-consuming,
      2. Do not guarantee to approximate the true MLE.
Example 2 A simulation study.
Generate RC data from S_Y(y) = \exp(-(y/\tau)^{\gamma}), y > 0, \tau = e^6, \gamma = 1/3. \mu = \tau \Gamma(1 + 1/\gamma).
Then pretend we do not know (\tau, \gamma). Find its MLE.
       #in R
      library(survival)
      #RC data
      set.seed(1)
      b=exp(6) \# \tau = e^{6}
      g=1/3 \# \gamma = 1/3
      y=rweibull(100,g,b)
                                                            \# mean(y) [1] 4145.353 \# median(y) [1] 287.3097
      c=runif(100,0,780) \# censoring variable
      d=as.numeric(y < =c)
      m = y^*d + c^*(1-d)
       #Monte Carlo
      set.seed(1)
\# Either use loop
      N = 200
      L=0
      M = c(0,0)
for(i in 1:N)
                              ł
      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))) \# (m,d) given
      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# Or use function $> b = \exp(6)$ > g = 1/3> bb=rnorm(N,b,100)> gg=rnorm(N,g,1/9) > a = 1:N> a = sapply(1:N, function(i))prod(d*dweibull(m,gg[i],bb[i])+(1-d)*(1-pweibull(m,gg[i],bb[i]))))> bb[a==max(a)]> gg[a = max(a)]**3.4.3.** R commands. R package provides some program to derive the MLE. Probability Distribution Functions in R. Let X be a random variable (rv). Its cdf $F(t) = P\{X \le t\},\$ density function (df) $f(t) = \begin{cases} F'(t) & \text{if } X \text{ is continuous,} \\ F(t) - F(t-) & \text{if } X \text{ is discrete.} \end{cases}$ quantile $Q(u) = F^{-1}(u) = \min\{t : F(t) \ge u\}.$ **Example 2.** $X \sim$ Weibull distribution with cdf $F(x|\gamma,\tau) = 1 - \exp(-(x/\tau)^{\gamma}), \ x > 0$ 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 the next table.

Dist	$S \ name$	parameters
beta	beta	shape1, shape2
binomial	binom	size, prob
Cauchy	cauchy	location, scale
chi-square	chisq	$d\!f$
exponential	exp	rate
F	f	$d\!f1, d\!f2$
gamma	gamma	shape, rate
geometric	geom	prob
hypergeometric	hyper	m,n,k
log-normal	lnorm	meanlog, sdlog
logistic	log is	location, scale
negative binomial	nbinom	size, prob
normal	norm	mean, sd
Poisson	pois	lambda
T	t	$d\!f$
uniform	unif	min, max
Weibull	weibull	shape, scale
Wilcox	wilcox	m,n

?dwilcox

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

- 1. weibull distribution $S(t) = \exp((-(\alpha t)^{\gamma}) = exp(-(\frac{t}{\beta})^{\gamma}) = exp(-(\frac{t}{\beta})^{1/\sigma}).$ Standard form $S(t) = exp(-t^{\gamma}/\theta), t > 0. \ \mu = \theta^{1/\gamma}\Gamma(1+1/\gamma). \ E(X^2) = \theta^{2/\gamma}\Gamma(1+2/\gamma).$ With covariate in R, reparametrization: $S_Y(y|\mathbf{x}, \beta, \sigma) = exp(-(y/e^{\beta'\mathbf{x}})^{1/\sigma}), \ y > 0, \ \mathbf{x}' = (1, x_1, ..., x_p)$ $\ln Y = \beta'\mathbf{x} + \sigma \ln T, \ T \sim Exp(1).$ $T = (Y/e^{\beta'\mathbf{x}})^{1/\sigma}. = Y^{1/\sigma}/e^{\beta'\mathbf{x}/\sigma} = Y^{\gamma}/\theta = (Y/\tau)^{\gamma}$ 2. exponential distribution Standard form $S(t) = exp(-t/\theta), \ t > 0$ With provide the provided of the set of t
 - With covariate in R, reparametrization: $S_Y(y|\mathbf{x},\beta) = exp(-y/e^{\beta'\mathbf{X}}), y > 0.$
 - $\ln Y = \beta' \mathbf{x} + \ln T, \ T \sim Exp(1). \qquad \qquad T = Y/e^{\beta' \mathbf{x}} = Y/\theta$
- 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: $S_Y(y|\mathbf{x}, \beta, \tau) = \frac{1}{1+exp(\frac{y-\beta'\mathbf{X}}{\tau})},$

$$Y = \beta' \mathbf{x} + \tau T, \ T \sim logistic(0, 1), \text{ with } \sigma_T = \pi/\sqrt{3}.$$

5. lognormal distribution

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

6. loglogistic distribution

$$\ln Y = \beta' \mathbf{x} + \tau T, \ T \sim logistic(0, 1).$$

Remark. With complete data,

Lognormal distribution: $E(\ln Y | \mathbf{x}) = \beta' \mathbf{x}$, and MLE = LSE.

Weibull dist.: $E(\ln Y | \mathbf{x}) = \beta' \mathbf{x} + \sigma \underbrace{E(\ln T)}_{\approx -0.59 \neq 0}$, $T \sim Exp(1)$ and $\ln(E(T)) = 0$. MLE \neq LSE.

Logistic distribution: $E(Y|\mathbf{x}) = \beta'\mathbf{x}$ and MLE \neq LSE. Log-logistic distribution: $E(\ln Y|\mathbf{x}) = \beta'\mathbf{x}$ and MLE \neq LSE. **B** command:

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 codes as follows.

library(survival)

 $zz=survreg(Surv(m,d)\sim x, dist="exponential")$ #survreg in R =survReg in Splus dist: (default: weibull), gaussian, logistic, lognormal and loglogistic predict(zz,data.frame(x=130),se=T) summary(zz)

Example 2. Notice that the scale in rweibull() is different from the scale in survreg(). 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($\cdot \sim 1$): $\ln Y = \alpha + \sigma \ln T$, where $T \sim Exp(1)$, $\alpha =$ intercept and $\sigma =$ scale (in survreg($\cdot \sim x$): $\ln Y = \alpha + \beta x + \sigma \ln T$ for regression). E(T) = 1, $E(\ln T) = ?$ Q: Relation between (γ, τ) and (σ, α) ? $Y = \exp(\alpha)T^{\sigma}$ $S_Y(t) = P(e^{\alpha}T^{\sigma} > t) = \exp(-e^{-\alpha/\sigma}t^{1/\sigma}) = \exp(-(t/e^{\alpha})^{1/\sigma}) = \exp(-(t/\tau)^{\gamma})$ $= > \gamma = 1/\sigma$ and $\tau = e^{\alpha}$.

Example 3. Simulation under Exponential distribution. Understanding the parameter and output.

> n = 10mean(Y)=0.5> y = rexp(n,2)> c = runif(n, 0, 2)> d=as.numeric(y < =c) $> m = y^* d + c^* (1-d)$ > (zz=survreg(Surv(m,d)~1,dist="exponential")) (Intercept) -0.5628376 $\# (= \hat{\alpha}, \alpha = \ln E(Y) = \ln \frac{1}{2} \approx -0.693)$ > predict(zz,data.frame(x=0),se=T) \$fit 0.5695905 v.s. mean(Y)=0.5\$se.fit 0.215285 $\# \ln Y = \beta' \mathbf{x} + \ln T \ (T \sim Exp(1) \text{ or } Y = e^{\beta' \mathbf{X}} T$ $> \exp(zz \text{scoef})$ $\ln Y = \alpha \times 1 + 0 \times x + \ln T = \beta' \mathbf{x} + \sigma \ln T$ (Intercept) 0.5695905 > predict(zz,data.frame(x=2),se=T) =?? **Remark.** What does predict() estimate ? predict(zz,data.frame(x=3)) = estimate of $\begin{cases} E(Y|X=3) & \text{if logis or guassian or Exp} \\ e^{E(\log Y|X=3)} & \text{if log-normal or loglogis} \\ e^{\alpha+3\beta} & \text{if Exp or weibul!} \end{cases}$ $zz=survreg(Surv(y,d)\sim x, dist="")$ \neq estimate of E(Y|X=3) in general. Thus for log-normal or loglogis, predict(zz,data.frame(x=3)) = exp($\hat{E}(\log Y | X = 3)) \not\approx E(Y | X = 3)$ > summary(zz) Value Std.Error zp-1.49(Intercept) -0.5630.3780.14 for H_o : intercept $\alpha = 0$. valid ? isitQuestion: Does the p-value support H_o ? $\alpha = ?$ Why does the survreg suggest such a p-value ? **Remark.** 1. Std.Error $\neq \sigma_{\hat{\alpha}}$ 2. Std.Error $\approx \sigma_{\hat{\alpha}}$ if *n* is very large. > (u=mean(m)/mean(d))[1] 0.5695905 # MLE of E(Y) with RC data (= 1/2) > mean(y)

[1] 0.4213111 # MLE \overline{y} of E(Y) with complete data.

$$\begin{split} S(t) &= \exp(-2t) = \exp(-t/0.5) = \exp(-\frac{t}{e^{\alpha}}) = \exp(-t/e^{-\ln 2}) \\ &\ln Y = \alpha + \ln T, \ T \sim Exp(1), \ e^{\alpha} = 1/2, \ \alpha = -\ln 2 \\ &> \log(u) \\ & [1] - 0.5628376 \qquad \log(\hat{E}(Y)) \ (\text{see Summary}(zz)) \\ &> \log(1/2) \\ & [1] - 0.6931472 \qquad \log(E(Y)) \\ &\textbf{Example 2 in §3.4.2.} \ (\text{continued}). \\ & (\text{Generate 100 RC data from Monte Carlo example with } \tau = e^6 \ \text{and } \gamma = 1/3.) \\ & ((\text{m,d}) - 100 \text{ RC data from Monte Carlo example with } \tau = e^6 \ \text{and } \gamma = 1/3.) \\ & ((\text{m,d}) - 100 \text{ RC data from Monte Carlo example with } \tau = e^6 \ \text{and } \gamma = 1/3.) \\ & ((\text{m,d}) - 100 \text{ RC data from Monte Carlo example with } \tau = e^6 \ \text{and } \gamma = 1/3.) \\ & ((\text{m,d}) - 100 \text{ RC data from Monte Carlo example with } \tau = e^6 \ \text{and } \gamma = 1/3.) \\ & ((\text{m,d}) - 100 \text{ RC data from Monte Carlo example with } \tau = e^6 \ \text{and } \gamma = 1/3.) \\ & ((\text{m,d}) - 100 \text{ RC data from Monte Carlo example with } \tau = e^6 \ \text{and } \gamma = 1/3.) \\ & ((\text{m,d}) - 100 \text{ RC data from Monte Carlo example with } \tau = e^6 \ \text{and } \gamma = 1/3.) \\ & ((\text{m,del}) = -142.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.78 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.78 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.78 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.78 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.78 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.78 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.78 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.78 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.78 \text{ Loglik}(\text{model}) = -342.3 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.2079198 \text{ Loglik}(\text{intercept only}) = -342.3 \text{ Scale} = 2.20$$

Now fit the data to another model with unrelated covariate

$$\begin{split} &> {\rm x=rpois(100,1)} \ \# \ {\rm unrelated\ covariate} \\ &> ({\rm zz=survreg(Surv(m,d)}{\sim}{\rm x}{+}{\rm I(x^2)})) \ \# \ {\rm log\ } y = 6 + 0 * x + 0x^2 + 3 * lnT, \ T \sim {\rm exp(1)} \\ &\quad (Intercept) \qquad x \qquad I(x^2) \\ &\quad 5.9362703 \quad 1.4134716 \quad -0.6472259 \\ &\quad {\rm Scale=\ } 2.731045 \\ &\quad {\rm Loglik(model)=\ -340.1\ Loglik(intercept\ only)=\ -342.3} \\ &\quad {\rm Chisq=\ } 4.45 \ {\rm on\ } 2 \ {\rm degrees\ of\ freedom,\ p=\ } 0.11 \\ &> {\rm summary(zz)} \end{split}$$

Value Std.Error \boldsymbol{z} p(Intercept) 5.9360.53311.138.98e - 291.83e - 011.0621.33x1.413 $I(x^2)$ 7.19e - 02-0.6470.360-1.804.81e - 19Log(scale)1.0050.113 8.92 Scale = 2.73Loglik(model) = -340.1 Loglik(intercept only) = -342.3Chisq= 4.45 on 2 degrees of freedom, p = 0.11Number 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 \boldsymbol{z} p8.975185e - 29(Intercept) 5.9362703 0.5333632 11.129884 1.4134716 1.06158121.331478 1.830319e - 01x $I(x^2)$ -1.7994487.194787e - 02-0.64722590.3596803Log(scale)4.808864e - 191.0046844 0.11267598.916585 > z var # covariance matrix $I(x^2)$ Log(scale)(Intercept)x0.0950058340.016751330(Intercept) 0.32635298-0.38842957-0.388429571.08207114 -0.345080622-0.015565844x $I(x^2)$ -0.345080620.1258692930.09500583 0.004110281 Log(scale)0.01675133 -0.015565840.0041102810.012808764 **Remark.** In parametric analysis, we want to make statistical inferences: 1. $\theta = ?$ (estimate θ with SD) (often use the MLE) 2. $F(t; \theta) = ?$ (estimate $F(t; \theta)$ with SD) (often use the MLE. How ?) 3. Y|(X = x) = ? (predict Y with given X = x with SD) 4. Test $H_o: \theta = \theta_o$. MLE solutions under assumptions in Ex. 2. (Weibull) Value Std.Error zp(Intercept) 5.990.37116.141.33e - 58Log(scale)9.021.94e - 191.020.113Scale = 2.781. Formula: $T = Y^{\gamma}/\theta$, or $= Y^{1/\sigma}/e^{\beta' \mathbf{X}/\sigma}$ (with covariates), $\overset{h(\hat{\eta})-h(\eta)}{\longrightarrow} \overset{\mathcal{D}}{\longrightarrow} N(0,1) \ (=> \sigma_{h(\hat{\eta})} \approx |h'(\hat{\eta})|\sigma_{\hat{\eta}}).$ $|\dot{h'}(\eta)|\sigma_{\hat{\eta}}$ $h(\hat{\eta}) - h(\eta) \longrightarrow \mathcal{D} N(0,1)$ $\sqrt{\frac{\partial \overline{h(\eta)}}{\partial \eta}^t \Sigma_{\hat{\eta}} \frac{\partial h(\eta)}{\partial \eta}}$ Parameters are (β, σ) or (γ, θ) . surveg yields $(\hat{\beta}, \hat{\sigma}) = (5.99, 2.78)$ with $\mathbf{x} = 1$, $(\beta, \sigma) = (6, 3).$ The SE of β is 0.37 and the SE of $\sigma = ?$ $\sigma = scale = \exp(log(scale)) = h(\log(scale)).$ $h'(x) = (e^x)' = e^x = scale = 2.78,$

the SE of the estimate of $\log(\text{scale})$ is 0.113. the SE of is estimate of scale $0.113 \times 2.78 = 0.314$. 2. parameters in pweibull: $(\gamma, \tau) = (1/\sigma, e^{\beta' \mathbf{X}})$, estimated by $(1/2.78, \exp(5.99))$. pweibull(t, 1/2.78, exp(5.99)). $\hat{F}(t;\hat{\gamma},\hat{\tau})$ 3. $\hat{Y}|(X = x) = ?$ (predict Y with given X = x with SE) predict(zz, x=1, se=T)4. Test H_o : $\sigma = 3$. or $\beta = ?$ $\mathbf{1}_{(\frac{|1.02-\ln 3|}{0.113}>1.96)}$ > abs(1.02 - log(3))/0.113[1] 0.695684conclusion of the test ? **Example 4** (simulation study under weibull distribution). #complete data (uncensored) > n = 100 $\tau=e^6$ $> b = \exp(6)$ > g = 1/3> y=rweibull(n,g,b) # $S(y) = \exp(-(y/\tau)^{\gamma}) = \exp(-(y/\tau)^{\gamma}) = \exp(-(y/\tau)^{1/3}) = \exp(-(y/\tau)^{1/3}) = \exp(-(y/\tau)^{1/3}/e^2)$ > $\exp(2^*3)^{*3*2} = \mu = \theta^{1/\gamma} \Gamma(1+1/\gamma), \qquad \theta = e^2, \ \gamma = 1/3$ [1] 2420.573> mean(y)[1] 2224.555 > z = rexp(n) $> y = \exp(6 + 3 \log(z)) \# \log y = 6 + 3\ln(z)$, 2nd way to generate rweibull > mean(y)[1] 2117.473 > (zz=survreg(Surv(y)~1, dist="exponential")) # Is it a correct model ? (Intercept) 7.657979 Scale fixed at 1 > (zz=survreg(Surv(y)~1)) # Is it a correct model ? (Intercept) 6.259506 Scale = 2.570908> y = rweibull(100, 1/3, exp(6))> zz=survreg(Surv(y)~1) > summary(zz) Std.Error Value zp(Intercept) 5.7470.2810 $20.5 \quad 5.77e - 93$ $12.7 \quad 3.94e - 37$ Log(scale)0.9790.0769Scale = 2.66# compare to previous 6.259506 and 2.570908 > zz=survreg(Surv(y)~1, dist="lognormal") # Is it a correct model ? > summary(zz) Std.ErrorValue zp0.3014 16.1 < 2e - 16(Intercept) 4.838315.6 < 2e - 16Log(scale)0.0707 1.1034Scale = 3.01

compare to above outcomes and the outcome below. Conclusion ? $> \ln(\log(y) \sim 1)$ LSE (Intercept) 4.838does not fall in the CI based on the MLE of weibull. $> 5.75 - 2.66 \times 0.56$ $(= \hat{\alpha} + \hat{\sigma} * E(\ln(T))$ in weibull) [1] 4.263623 > y = rnorm(100, 8, 4)> (zz=survreg(Surv(y)~1,dist="gaussian")) (Intercept) 7.929682 Scale = 3.741224 $> u = \exp(y)$ $> (zz=survreg(Surv(u)\sim 1, dist="lognormal"))$ (Intercept) 7.929682 Scale = 3.741224Value Std.Error zp> summary(zz) (Intercept) 7.9297 0.3741 21.2< 2e - 16Log(scale)1.31940.0707 18.7< 2e - 16Does the CI include the true value ? > n = 400> z = rexp(n)> x = runif(n)> y=exp(6+2*x+3*log(z)) # lny = $\alpha + \beta_1 * x + 3 * \ln z$ > mean(y) [1] 4714.647 > (zz=survreg(Surv(y)~x, dist="exponential")) (Intercept) xIs the MLE consistent ? 0.8000083 8.0369979 Scale fixed at 1 > summary(zz) Std.Error Value zp(Intercept) 8.0370 0.093486.01 < 2e - 16x0.8000 0.1616 4.957.4e - 07Scale fixed at 1 Does the CI include the true value ? > (zz=survreg(Surv(y)~x)) (Intercept) x6.127987 1.750514Scale = 2.775253> summary(zz) Value Std.Error \boldsymbol{z} p(Intercept 6.12800.2702 22.7< 2e - 161.75050.4605 3.80.00014xLog(scale)1.02070.0392 26.0 < 2e - 16Scale = 2.78

Does the CI include the true value ?
Is the MLE consistent ? > (zz=survreg(Surv(y)~x, dist="weibull")) (Intercept) x6.127987 1.750514 Scale = 2.775253> (zz=survreg(Surv(y)~x, dist="lognormal")) > summary(zz) Value Std.Error \boldsymbol{z} < 2e - 16 is the MLE consistent ? 0.3481 11.87 (Intercept) 4.13323.3e - 052.52320.60814.15xLog(scale)1.28920.035436.46 < 2e - 16Scale= 3.63 Does the CI of α include the true value ? Is it expected ? $> 6.13 + 2.78^{*}(-0.56)$ $(= \hat{\alpha} + \hat{\sigma} * E(\ln(T))$ in weibull) [1] 4.57 \in CI of α based on LSE. Does the CI of β_1 include the true value ? Is it expected ? > y = rnorm(100, 8, 4) $> u = \log(y)$ > (zz=survreg(Surv(u)~x,dist="gaussian")) Is the MLE consistent ? (Intercept) x4.133211 2.523191 Scale = 3.62989> (zz=survreg(Surv(y)~x,dist="gaussian")) compare to the MLE in the above case. (Intercept) 2737.883 4047.005 Scale = 12263.48

§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. With the above data, estimate the MLE of the parameter using the Newton-Raphson method

(you could use the command in R:

y=rweibull(100,1/2,1) # check whether the parameters are right

c = runif(100, 0, 2)

 $d=as.numeric(y \le c)$

 $m = d^*y + (1-d)^*c$

 $yy=survreg(Surv(m,d) \sim 1, dist="weibull")$

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? Do a simulation to check your answer.
- 8. Compute $\frac{\partial^2 \mathcal{L}}{\partial \beta \partial \beta^t}$ for a weibull distribution $S(t|\beta, \sigma) = \exp(-(t/e^{\beta' \mathbf{X}})^{1/\sigma}), t > 0$, with $\mathbf{x} \in \mathcal{R}^6$ and n RC data.

\S **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 $\S3.3$. 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) $X \sim Exp(\rho_o)$,

- 2) (U, V) has a joint c.d.f. G(u, v), where $0 \le U < V$ w.p.1.,
- 3) X and (U, V) are independent; $\int (U, V)$ if $X \in (U, V]$,

4)
$$(L, R) = \begin{cases} (V, +\infty) & \text{if } X > V, \\ (-\infty, U) & \text{if } X \le U. \end{cases}$$

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

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

Since $S_o(t, \rho) = e^{-\rho t}$, Eq. (1.1) yields that

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

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

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

There is no closed form expression for the MLE of ρ_o . Thus we need to use the Newton-Raphson method to derive it.

$\S3.5.1$. 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 $\hat{\rho} \in (0, +\infty) \text{ unless } \sum_{i,rc} 1 = \sum_{i} \mathbf{1}_{(R_i = \infty)} = n \text{ or } \sum_{i,lc} 1 = \sum_{i} \mathbf{1}_{(L_i = -\infty)} = n.$ 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) $\lim_{\rho \to 0^+} \mathcal{L}(\rho) = -\infty;$

(ii) $\lim_{\rho \to +\infty} \mathcal{L}(\rho) = -\infty;$

(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

 $\mathcal{L}(\rho) < \mathcal{L}(1)$ 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) ?) \Box Does the MLE always exist ?

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 but we can 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 but we can define $\hat{\rho} = +\infty$. In this way, $\hat{\rho}$ is properly defined in the whole sample space and thus we can study its properties.

$\S3.5.2$. Consistency of the MLE.

We assume:

- 1) $X \sim Exp(\rho_o)$,
- 2) (U, V) has a joint c.d.f. G(u, v), where $0 \le U < V$ w.p.1.,
- 3) X and (U, V) are independent;

4)
$$(L, R) = \begin{cases} (U, V) & \text{if } X \in (U, V] \\ (V, +\infty) & \text{if } X > V, \\ (-\infty, U) & \text{if } X \le U, \end{cases}$$

Hereafter, abusing notations, we write $S = S(t) = e^{-\rho t}$ etc.. The normalized log likelihood function

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

Theorem 2. Suppose that $X \sim Exp(\rho_o)$. Under the C2 model with P(0 < U < V) = 1, $\hat{\rho}$ defined in the remark behind Theorem 1 is strongly consistent. **Proof.** Let $\mu(\rho) = E(l(\rho))$ and

 $\rho^*(\omega)$ be a limiting point of $\hat{\rho}(\omega)$, where $\omega \in \Omega$, the sample space. We shall show

$$\mu(\rho) \le \mu(\rho_o)$$
 with equality iff $\rho = \rho_o$. (2.2)

$$\rho^* = \rho^*(\omega) \in (0,\infty) \text{ and } \mu(\rho^*) \ge \mu(\rho_o) \text{ for all } \omega \in \Omega^* \text{ with } P(\Omega^*) = 1. \quad \Omega^* = \Omega?$$
(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 $\hat{\rho} \to \rho_o$ 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 0 or ∞ .

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

$$\begin{split} & \mu(\rho) \quad (= E(l(\mathbf{L}, \mathbf{R}, \rho))) \\ = E(E(l(\mathbf{L}, \mathbf{R}, \rho) | \mathbf{U}, \mathbf{V})) & (\mathbf{U}, \mathbf{V}) = (U, V)? \quad \text{see } (2.1)) \\ = E(E[\mathbf{1}_{(X \le U)} \log F(U; \rho) + \mathbf{1}_{(X > V)} \log S(V; \rho) + \mathbf{1}_{(U < X \le V)} \log(S(U; \rho) - S(V; \rho)) | U, V])? \\ = E[F(U; \rho_o) \log F(U; \rho) + S(V; \rho_o) \log S(V; \rho) + (S(U; \rho_o) - S(V; \rho_o)) \log(S(U; \rho) - S(V; \rho))]? \end{split}$$

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

(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 $S(v; \rho) = S(v; \rho_o)$ (see Remark below). (b) statement (a) implies $\rho = \rho_o$ as $S(x, \rho) = e^{-\rho x}, x > 0$. (c) $0 \ge \mu(\rho_o) \ge -3/e.$ $(\mu(\rho) = E(l(\mathbf{L}, \mathbf{R}, \rho)))$ before Eq.(2.2). To prove (c) note that $\mu(\rho_o) = g(p_1, p_2)$ (given above), each summand in $g(p_1, p_2)$ is of the form $x \ln x$, $0 \ge x \log x \ge -1/e$, for $x \in [0, 1]$, (check yourself why !) and thus $g(p_1, p_2) \geq -3/e$, Thus (c) follows. (a), (b) and (c) imply

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

 $\mu(\rho) = E(l(\rho)) = \sum_{i} f(i;\theta_o) \log f(i;\theta) \leq \sum_{i} f(i;\theta_o) \log f(i;\theta_o)$ with equality iff $\theta = \theta_o$ if $\mu(\rho_o) < \infty$, where f is the d.f. of trinomial distribution $Multi(1, p_1, p_2, p_3)$

$$\sum_{i} f(i;\theta_o) \log \frac{f(i;\theta_o)}{f(i;\theta)} = \int f(x;\theta_o) \log \frac{f(x;\theta_o)}{f(x;\theta)} d\nu(x) \ge 0, \text{ with equality iff } \theta = \theta_o$$

where ν is the counting measure.

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

 $g(q_1, q_2) = p_1 \log q_1 + p_2 \log q_2 + p_3 \log q_3$, (where $q_i, p_i \ge 0$, $\sum_i p_i = \sum_i q_i = 1$) is uniquely maximized by $q_1 = p_1$ and $q_2 = p_2$. Class exercise.

Prove (2.3). We shall now construct Ω^* . To emphasize that $\hat{\rho}$ is a function of n, we write $\hat{\rho} = \hat{\rho}_n$. By definition of the modified MLE $\hat{\rho}$ (why MMLE ?) $l(\hat{\rho}(\omega))(\omega) \ge l(\rho_o)(\omega) \forall \omega \in \Omega$. Thus

$$\liminf_{n \to \infty} l(\hat{\rho}) \ge \lim_{n \to \infty} l(\rho_o) = \mu(\rho_o) \quad \text{a.s. (by SLLN)}.$$
(2.4)

$$l(\rho) = \sum_{i=1}^{n} \frac{\log(1 - e^{-\rho U_i}) \mathbf{1}_{(X_i \le U_i)}}{n} + \sum_{i=1}^{n} \frac{\log(1 - e^{-\rho(V_i - U_i)}) \mathbf{1}_{(X_i \in (U_i, V_i])}}{n} \qquad by \ (2.1)$$

$$-\rho \sum_{i=1}^{n} \frac{U_i \mathbf{1}_{(X_i \in (U_i, V_i])}}{n} - \rho \sum_{i=1}^{n} \frac{V_i \mathbf{1}_{(X_i > V_i)}}{n} \quad (= \overline{Z}_1(\rho) + \overline{Z}_2(\rho) - \rho \overline{Z}_3 - \rho \overline{Z}_4).$$
(2.5)

For each ρ , the four summations in (2.5) all converges a.s. to their means, respectively.

$$P(l(\rho) \to \mu(\rho)) = 1$$
 for all $\rho > 0$.

However, it is not clear that $P(l(\rho) \to \mu(\rho) \text{ for all } \rho > 0) = 1$ and $P(l(\hat{\rho}) \to \mu(\rho^*)) = 1$ (ρ^* as in (2.3)). Thus, we let K be the set of all positive rational numbers and ρ_o , Ω_ρ be the event such that (2.4) holds for all $\rho \in K$ and the four summations in (2.5) converges a.s. to their means, respectively, and let $\Omega^* = \bigcap_{q \in K} \Omega_q$. Then $P(\Omega^*) \ (= P(\bigcap_{q \in K} \Omega_q)) = 1$ as K is countable.

For each ω in Ω^* , let ρ^* be a limiting point of $\hat{\rho} = \hat{\rho}_n(\omega)$ in the sense that

 $\hat{\rho}_{n_j}(\omega) \to \rho^*$ for a subsequence of $\{\hat{\rho}_n\}_{n\geq 1}$, where ρ^* may be $+\infty$ or 0. In order to prove inequality (2.3) (see (2.3) and (2.2)), it suffices to prove

$$\mu(\rho_o) \le \mu(\rho^*). \tag{2.6}$$

We shall show (2.6) hereafter. Notice that

$$l(\hat{\rho}(\omega)) = \sum_{i=1}^{n} \frac{\log(1 - e^{-\hat{\rho}(\omega)U_{i}})\mathbf{1}_{(X_{i} \leq U_{i})}}{n} + \sum_{i=1}^{n} \frac{\log(1 - e^{-\hat{\rho}(\omega)(V_{i} - U_{i})})\mathbf{1}_{(X_{i} \in (U_{i}, V_{i}])}}{n} - \hat{\rho}(\omega) \sum_{\substack{i=1\\ext{i} \\ ext{i} \\ ext{i}$$

We first show that it is impossible that $\rho^* = -\infty$ or 0. Assume $\omega \in \Omega^*$. If $\rho^* = +\infty$, then $\overline{Z}_1 < 0$, $\overline{Z}_2 < 0$, $\overline{Z}_3 > 0$ (see (2.5)) and Eq. (2.5^{*}) implies that

$$l(\hat{\rho}_{n_j}(\omega)) \le -\hat{\rho}_{n_j}\overline{Z}_4 = -\hat{\rho}_{n_j}(\omega)\left[\sum_{i=1}^n V_i \frac{\mathbf{1}_{(X_i > V_i)}}{n}\right](\omega) \to -\infty$$
(2.7)

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

$$-\infty = \lim_{n_j \to \infty} l(\hat{\rho}(\omega)) \ge \mu(\rho_o)$$
 (which is finite by (2.2))

It reaches a contradiction. Thus $\rho^* = +\infty$ is impossible.

If $\rho^* = 0$, last 2 terms in (2.5) equal 0 and the first two sums in Eq. (2.5^{*}) tend $\log(0+) = -\infty \ge \mu(\rho_o)$ as $n \to \infty$. It leads to a contradiction again. Thus $\rho^* = 0$ 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 $m < \hat{\rho}_{n_j} < M$ (why ?) Since $\log(1 - e^{-\rho x}) \uparrow$ in ρ, \overline{Z}_1 (in (2.5))

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

In an obvious way, rewrite the above inequalities as

$$\begin{aligned}
\Psi_{n_j}(m,\omega) &\leq \Psi_{n_j}(\hat{\rho}_{n_j}(\omega),\omega) \leq \Psi_{n_j}(M,\omega). \end{aligned}$$

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

Since $m, M \in K$,

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

 $\operatorname{limsup}_{i \to \infty} \Psi_{n_i}(M, \omega) = \operatorname{lim}_{n \to \infty} \Psi_n(M, \omega) = E(\log(1 - e^{-MU}) \mathbf{1}_{(X < U)}).$ (2.11)

Recall the monotone convergence theorem:

If f_n is a monotone convergent sequence, and $f_n \to f$ and f_n are all integrable, $\begin{array}{l} then \int f_n(x)d\mu(x) \to \int f(x)d\mu(x). \\ \text{Since } g(\rho,\omega) = \log(1-e^{-\rho U(\omega)})\mathbf{1}(X(\omega) \leq U(\omega)) \text{ is a monotone function of } \rho, \end{array}$

 $g(m, \cdot)$ (or $g(M, \cdot)$) is an increasing function in m (or M) as $m \uparrow \rho^*$ (or $M \downarrow \rho^*$), and $E(q(m, \cdot)) = \int q(m, \omega) dP(\omega),$

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

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

Then it follows from (2.8) through (2.12) that the first summand \overline{Z}_1 in (2.5^{*})

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

(Remark: In (2.13), we are dealing $\lim_{n\to\infty}\sum_{i=1}^n g_i(\hat{\rho}_n)/n$, does SLLN work ? Without the arguments from (2.8) through (2.12), we cannot conclude directly,

$$\lim_{n \to \infty} \sum_{i=1}^{n} g_i(\hat{\rho}_n)/n = \lim_{n \to \infty} \sum_{i=1}^{n} \lim_{n \to \infty} g_i(\hat{\rho}_n)/n,$$

even if all the limits make sense.)

Since the 2nd summand in (2.5^{*}) $\log(1 - e^{-\rho(V_i - U_i)})$ is also a monotone function of ρ , it can be shown that

$$\sum_{i=1}^{n_j} \log(1 - e^{-\hat{\rho}_{n_j}(\omega)(V_i - U_i)(\omega)}) \mathbf{1}_{(X_i(\omega) \in (U_i(\omega), V_i(\omega)])} / n_j$$

$$\rightarrow E(\log(1 - e^{-\rho^*(V - U)}) \mathbf{1}_{(X \in (U, V])}).$$
(2.14)

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

Thus, $\liminf_{n_j \to \infty} l_{n_j}(\hat{\rho}_{n_j}(\omega)) = \mu(\rho^*)$, 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

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

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

1. $P(U < V) = 1 \not\Rightarrow U \not\perp V$. e.g. If $U \sim U(0, 1)$, $V \sim U(1, 2)$, and $U \perp V$, then P(U < V) = 1. 2. P(U < V) = 1 and $\sup\{t: F_U(t) < 1\} > \inf\{t: F_V(t) > 0\} \Rightarrow U \not\perp V$. **3.5.2.1. Homework:** (1) Prove the two statements in the previous Remark.

(2) Prove (2.14) by mimicking the proof of (2.13).

In this section, let f_o and f be two densities w.r.t. a measure μ (which does not have to be continuous). Denote

 $I(f_o, f) = \int f_o(t) \ln(f_o/f)(t) d\mu(t).$

SK Inequality. If $\int f_o(t) \ln f_o(t) d\mu(t)$ is finite, then

 $\int f_o(t) \ln f_o(t) d\mu(t) \geq \int f_o(t) \ln f(t) d\mu(t); \text{ with equality iff } \int |f(t) - f_o(t)| d\mu(t) = 0.$ **KL inequality**. $I(f_o, f) \geq 0;$ with equality iff $\int |f(t) - f_o(t)| d\mu(t) = 0.$

It is worth mentioning that Kullback and Leibler (1951) proved that $I(f_o, f)$ exists, though it may be ∞ . Are these two inequalities equivalent ?

3.5.2.2. Homework: Given a counterexample to that $\int f_o(t) \ln f_o(t) d\mu(t)$ is finite.

In his classical textbook, Ferguson (1996, p.114) showed that the MLE of θ is consistent if the following conditions hold:

(A1) $X_1, ..., X_n$ are i.i.d. observations from $f(\cdot; \theta), \theta \in \Theta$ and θ_o is the true value of θ ; (A2) $\int |f(x; \theta) - f(x; \theta_o)| d\mu(x) = 0$ implies that $\theta = \theta_o$ (identifiability);

(A3) $\overline{\lim}_{\theta_n \to \theta} f(x; \theta_n) \leq f(x; \theta)$ or $\lim_{\delta \to 0^+} \sup_{|\theta - \theta'| < \delta} f(x; \theta') = f(x; \theta), \forall x;$

(A4) Θ is compact;

(A5) \exists a function K(x) such that $E_{\theta_o}(|K(X)|) < \infty$ and $\log \frac{f(x;\theta)}{f(x;\theta_o)} \leq K(x) \ \forall \ (x,\theta);$

(A6) for all $\theta \in \Theta$, and sufficiently small $\delta > 0$, $\sup_{|\theta - \theta_o| < \delta} f(x; \theta)$ is measurable in x. **Remark.** Suppose that $Y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$, $i, j \in \{1, 2\}$, and $\epsilon_{ij} \sim N(0, 1)$. Then (A2) fails and $(\mu + \alpha_1, \alpha_2, \beta_1, \beta_2)$ is not identifiable. Counterexample: if $\theta = (\mu, \alpha_1, \alpha_2, \beta_1, \beta_2) = (0, 1, 2, 1, 2)$ and $\theta^* = (\mu^*, \alpha_1^*, \alpha_2^*, \beta_1^*, \beta_2^*) = (3, 1, -1, 1, 2)$, then $\int |f(\mathbf{x}; \theta) - f(\mathbf{x}; \theta^*)| d\mathbf{x} = 0$. **Remark.** The standard theory for proving the consistency of the MLE does not work here, as it requires that Θ is compact, which is not the case here.

Casella and Berger (2001, p.516) presented a set of somewhat simpler but not stronger sufficient conditions for the consistency of the MLE of θ in their popular textbook as follows: In addition to the aforementioned (A1) and (A2), the following conditions also hold:

(A7) The densities $f(x;\theta)$ have common support, and $f(x;\theta)$ is differentiable in θ .

(A8) The parameter space Θ contains an open set A and the true parameter $\theta_o \in A$. **Remark.** (A8) fails for bin(n, p) with parameter $(n, p) \in \{1, 2, ...\} \times [0, 1]$.

(A7) fails for $U(0,\theta)$ with parameter $\theta > 0$.

But it works for the current case.

Theorem 2^{*}. Suppose that $X \sim \text{Exp}(\rho_o)$. Suppose that (L_1, R_1) , ..., (L_n, R_n) are i.i.d. from the mixed case IC model, $\hat{\rho}$ defined in the remark after Theorem 1 is consistent. **Proof.** It suffices to verify assumption (A1), (A2), (A7) and (A8).

The df is $f(l,r;\rho) \propto (S(l;\rho) - S(r;\rho)) = (\exp(-\rho l) - \exp(-\rho r))\mathbf{1}(0 \le l < r \le \infty).$ (A1) is obviously true.

(A2) holds: If $f(l,r;\rho) = f(l,r;\theta)$ for all (l,r), then it is true also for $r = \infty$. $f(l,\infty;\rho) = f(l,\infty;\theta)$ for all l yields $e^{-\rho l} = e^{-\theta l}$ $=> \rho l = \theta l$ for all l $=> \rho = \theta$.

(A7) holds: (1) The support is $\{(l,r): 0 \le l < r \le \infty\}$ and it does not depends on the parameter ρ . So the support is common. (2) $\frac{\partial}{\partial a} f(l,r;\rho) = (-l \exp(-\rho l) + r \exp(-\rho r)) \mathbf{1} (0 \le l < r \le \infty)$ is also differentiable.

(A8) holds: $\Theta = (0, \infty)$ and $\rho_o \in \Theta$. \Box

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

Proof. We shall show in 3 steps that for each t, $P\{\frac{\sqrt{n}(\hat{\rho}-\rho_o)}{\sigma} \leq t\} \rightarrow \Phi(t)$, as $n \rightarrow \infty$. **Step 1** (preliminary). Let Ω^o be the subset of the sample space Ω such that

$$\hat{\rho} \to \rho_o, \lim_{n \to \infty} \overline{\mathbf{1}(X < U)} < 1 \text{ and } \lim_{n \to \infty} \overline{\mathbf{1}(X \ge V)} < 1.$$
 (3.0)

Then $P(\Omega^o) = 1$, as P(0 < U < V) = 1 (is it possible that $U = \infty$ or V = 0?)

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

$$\frac{\partial l(\rho)}{\partial \rho} = \frac{1}{n} \sum_{i=1}^{n} \frac{U_i e^{-\rho U_i}}{1 - e^{-\rho U_i}} \mathbb{1}_{\{X_i \le U_i\}} + \frac{1}{n} \sum_{i=1}^{n} \frac{(V_i - U_i) e^{-\rho (V_i - U_i)}}{1 - e^{-\rho (V_i - U_i)}} \mathbb{1}_{\{U_i < X_i \le V_i\}}$$

$$- \frac{1}{n} \sum_{i=1}^{n} U_i \mathbb{1}_{\{U_i < X_i \le V_i\}} - \frac{1}{n} \sum_{i=1}^{n} V_i \mathbb{1}_{\{X_i > V_i\}}$$
(next expression is for $\frac{\partial^2 l(\rho)}{\partial \rho^2}$)
$$= \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n} (-U_i - U_i) \mathbb{1}_{\{U_i < X_i \le V_i\}} + \frac{1}{n} \sum_{i=1}^{n}$$

$$= \frac{1}{n} \sum_{i=1}^{n} (-U_i + \frac{U_i}{1 - e^{-\rho U_i}}) \mathbf{1}_{(X_i \le U_i)} + \frac{1}{n} \sum_{i=1}^{n} (-(V_i - U_i) + \frac{(V_i - U_i)}{1 - e^{-\rho (V_i - U_i)}}) \mathbf{1}_{(U_i < X_i \le V_i)}$$

$$-\frac{1}{n}\sum_{i=1}^{n}U_{i}1_{(U_{i}< X_{i}\leq V_{i})} - \frac{1}{n}\sum_{i=1}^{n}V_{i}1_{(X_{i}>V_{i})}, \qquad how ?$$

$$\frac{\partial^2 l(\rho)}{\partial \rho^2} = -\frac{1}{n} \sum_{i=1}^n \frac{U_i^2 e^{-\rho U_i}}{(1 - e^{-\rho U_i})^2} \mathbb{1}_{\{X_i \le U_i\}} - \frac{1}{n} \sum_{i=1}^n \frac{(V_i - U_i)^2 e^{-\rho (V_i - U_i)}}{(1 - e^{-\rho (V_i - U_i)})^2}) \mathbb{1}_{\{U_i < X_i \le V_i\}}.$$
(3.1)

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

For each $\omega \in \Omega^{o}$, for *n* large enough, $\sum_{lc} 1/n < 1$ and $\sum_{rc} 1/n < 1$ (due to (3.?)), thus the MLE $\hat{\rho}(\omega) \in (0, +\infty)$ by Theorem 1. Since $\frac{\partial l(\rho)}{\partial \rho}(\omega)$ exists,

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

Since $\frac{\partial^2 l(\rho)}{\partial \rho^2}(\omega)$ is continuous in ρ , by the first order Taylor expansion,

$$\frac{\partial l(\rho)}{\partial \rho}(\omega)\Big|_{\rho=\rho_o} - \frac{\partial l(\rho)}{\partial \rho}(\omega)\Big|_{\rho=\hat{\rho}(\omega)} = \frac{\partial^2 l(\rho)}{\partial \rho^2(\omega)}\Big|_{\rho=\rho^*}(\rho_o - \hat{\rho}(\omega)) \text{ if } \omega \in \Omega^o, \tag{3.3}$$

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

$$\rho^* \to \rho_o, \ \forall \ \omega \in \Omega^o \ (\text{due to } \hat{\rho} \to \rho_o).$$

Write
$$l(\rho) = \frac{1}{n} \sum_{i=1}^{n} \log p(X_i, U_i, V_i, \rho)$$
 (see (2.1)), where
 $p(X, U, V, \rho) = F(U; \rho)^{\mathbf{1}_{(X \le U)}} \times (S(U; \rho) - S(V; \rho))^{\mathbf{1}_{(U < X \le V)}} \times S(V; \rho)^{\mathbf{1}_{(V < X)}};$
define $Z = \frac{\partial}{\partial \rho} \log p(X, U, V, \rho) \Big|_{\rho = \rho_o},$
then $\frac{\partial l}{\partial \rho} \Big|_{\rho = \rho_o} = \overline{Z}.$ (3.4)

By (3.2), (3.3) and (3.4),

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

but not just if $\omega \in \Omega^o$!) By the CLT,

$$\sqrt{n}(\overline{Z} - E(Z))/\sigma_Z \xrightarrow{\mathcal{D}} N(0, 1).$$
 (3.6)

Step 2. $\vdash: \sqrt{n}(\hat{\rho} - \rho_o) \xrightarrow{\mathcal{D}} N(0, \tau^2), \text{ for } \tau^2 = \sigma_Z^2 / (E(\frac{\partial^2 l(\rho_o)}{\partial \rho^2}))^2, \text{ if }$

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

and
$$\frac{\partial^2 l(\rho)}{\partial \rho^2}\Big|_{\rho^*} \to E(\frac{\partial^2 l(\rho_o)}{\partial \rho^2})$$
 a.s.. (3.8)

Notice that (3.5), (3.7) and (3.8) yield

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

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

$$thus \{-\sqrt{n}(\hat{\rho}-\rho_{o})\frac{\partial^{2}l(\rho^{*})}{\partial\rho^{2}} \neq \sqrt{nZ}\} \subset \{\hat{\rho} \notin (0,\infty)\},$$

$$(3.9)$$

$$|P\{\sqrt{n}(\hat{\rho}-\rho_{o})/\tau_{n} \leq t\} - P\{\sqrt{n}\frac{\overline{Z}}{\sigma_{Z}} \leq t\}|$$

$$(\tau_{n} = -\sigma_{Z}/\frac{\partial^{2}l(\rho^{*})}{\partial\rho^{2}})$$

$$= |P\{\sqrt{n}(\hat{\rho}-\rho_{o})/\tau_{n} \leq t\} - P\{\sqrt{n}\frac{\overline{Z}}{\sigma_{Z}} \leq t\}|$$

$$(\tau_{n} = -\sigma_{Z}/\frac{\partial^{2}l(\rho^{*})}{\partial\rho^{2}})$$

$$= |P\{\sqrt{n}(\hat{\rho}-\rho_{o})/\tau_{n} \leq t\} - P\{\sqrt{n}(\hat{\rho}-\rho_{o})/\tau_{n} \leq t, \hat{\rho} \in (0,\infty)\} + P\{\sqrt{n}\frac{\overline{Z}}{\sigma_{Z}} \leq t, \hat{\rho} \in (0,\infty)\}$$

$$(= 0 \text{ by (3.5)})$$

$$- P\{\sqrt{n}\frac{\overline{Z}}{\sigma_{Z}} \leq t\}|$$

$$\leq |P\{\sqrt{n}(\hat{\rho}-\rho_{o})/\tau_{n} \leq t\} - P\{\sqrt{n}(\hat{\rho}-\rho_{o})/\tau_{n} \leq t, \hat{\rho} \in (0,\infty)\}|$$

$$+ |P\{\sqrt{n}\frac{\overline{Z}}{\sigma_{Z}} \leq t, \hat{\rho} \in (0,\infty)\} - P\{\sqrt{n}\frac{\overline{Z}}{\sigma_{Z}} \leq t\}|$$

$$\leq P\{\hat{\rho} \notin (0,\infty)\} + P\{\hat{\rho} \notin (0,\infty)\} \quad (|P(A) - P(A \cap B)| \leq P(B^{c}))$$

$$(by (3.9))$$

$$i.e., |P\{\sqrt{n}(\hat{\rho}-\rho_{o})/\tau_{n} \leq t\} - P\{\sqrt{n}\frac{\overline{Z}}{\sigma_{Z}} \leq t\}| \rightarrow 0, \text{ as } n \rightarrow \infty \text{ due to } \hat{\rho} \rightarrow \rho_{o} \text{ a.s.}$$

$$(3.10)$$

$$P\{\sqrt{n}\frac{\overline{Z}}{\sigma_{Z}} \leq t\} = P\{\sqrt{n}\frac{\overline{Z} - E(Z)}{\sigma_{Z}} \leq t\} \rightarrow \Phi(t) \text{ as } n \rightarrow \infty,$$

where Φ is the cdf of N(0, 1) (due to CLT),

thus

$$P\{\sqrt{n}(\hat{\rho} - \rho_o)/\tau_n \leq t\}$$

$$= \underbrace{P\{\sqrt{n}(\hat{\rho} - \rho_o)/\tau_n \leq t\} - P\{\sqrt{n}\frac{\overline{Z}}{\sigma_Z} \leq t\}}_{see \ (3.10)} + P\{\sqrt{n}\frac{\overline{Z}}{\sigma_Z} \leq t\}$$

 $\rightarrow \Phi(t)$

Moreover, by Slustky's theorem

 $W_n \xrightarrow{\mathcal{D}} W$ and $T_n \xrightarrow{\mathcal{D}} b$ imply that $W_n T_n \xrightarrow{\mathcal{D}} Wb$, letting $W_n = \sqrt{n}(\hat{\rho} - \rho_o)/\tau_n$ 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)).

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

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

$$-\frac{\partial^2 l(\rho)}{\partial \rho^2}\Big|_{\rho^*} = \frac{1}{n} \sum_{i=1}^n \frac{-U_i^2}{1 - e^{-\rho U_i}} \mathbf{1}_{\{X_i \le U_i\}}\Big|_{\rho = \rho^*} + \frac{1}{n} \sum_{i=1}^n \frac{U_i^2}{(1 - e^{-\rho U_i})^2} \mathbf{1}_{\{X_i \le U_i\}}\Big|_{\rho = \rho^*}$$
(3.12)
$$+ \frac{1}{n} \sum_{i=1}^n \frac{-(V_i - U_i)^2}{1 - e^{-\rho (V_i - U_i)}} \mathbf{1}_{\{U_i < X_i \le V_i\}}\Big|_{\rho = \rho^*} + \frac{1}{n} \sum_{i=1}^n \frac{(V_i - U_i)^2}{(1 - e^{-\rho (V_i - U_i)})^2} \mathbf{1}_{\{U_i < X_i \le V_i\}}\Big|_{\rho = \rho^*}$$

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

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

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

$$\begin{split} &\frac{\partial^2 l}{\partial \rho^2} = \int \int \int \frac{\partial^2}{\partial \rho^2} \mathrm{ln} p(x, u, v, \rho) dF_n(x, u, v) \text{ and} \\ &E(\frac{\partial^2 l}{\partial \rho^2}) = \int \int \int \frac{\partial^2}{\partial \rho^2} \mathrm{ln} p(x, u, v, \rho) dF_o(x, u, v), \end{split}$$

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

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

Finally, verify that the Fisher information number (homework)

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

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

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 $|\frac{d^2l(\rho;x,u,v)}{d\rho^2}|$ 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 $X \sim F_X$ (cdf),

 $(X_i, L_i, R_i), i = 1, ..., n$ are i.i.d. from an extended random vector (X, L, R). **Question**: Observed (L_i, R_i) s,

$F_X = ?$ without further restriction !

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 set the parameter space as a compact space

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

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 $(L_i, R_i), i = 1, ..., n$, is

$$\mathcal{L}(F) = \prod_{i=1}^{n} \mu_F(I_i), \ F \in \Theta.$$
(1.1)

Eq. (1.1) 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_X maybe continuous. **Definition**. The generalized (or nonparametric) maximum likelihood estimator (GMLE) of F_X is an $F = \hat{F} \in \Theta$ such that

$$\hat{F}$$
 maximizes $L(F)$ over Θ .

Remark.

1. The GMLE is also called the nonparametric MLE (NPMLE).

2. If F is discrete, L(F) is the same as the definition of parametric likelihood $L(\phi)$ 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_X is given by

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

 \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 $\hat{F}(t) = \sum_{x \le t} \hat{f}(x)$.
4. $\overline{X} = \sum_t t \hat{f}(t)$.

> n = 100

> x = rnorm(100)

- > x = sort(x),
- > plot(x, pnorm(x), type='l')
- > lines(x,ecdf(x)(x),type='s'),

> lines(x,p+1.96*sqrt(p*(1-p)/n),type='s', lty=3) p=?

> lines(x,p-1.96*sqrt(p*(1-p)/n),type='s', lty=3)





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

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

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

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

where $Z_{(1)} \leq \cdots \leq Z_{(n)}$ 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)).

Example.
$$n = 6$$
, $\begin{pmatrix} order: 1 & 2 & 3 & 4 & 5 & 6 \\ data: 3+, 2, 3, 5 & 3 & 0.5+ \\ Z_{(i)}: & Z_{(5)} & Z_{(2)} & Z_{(3)} & Z_{(6)} & Z_{(4)} & Z_{(1)} \\ \delta_{(i)}: & 0 & 1 & 1 & 1 & 0 \end{pmatrix}$

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

$$d_k = \sum_{i=1}^n \mathbf{1}_{\{L_i = R_i = a_k\}}, \text{ $\#$ of deaths at a_k;} (in \text{ Ex.1}, \begin{pmatrix} a_i : 2 & 3 & 5 \\ d_i : 1 & ? & 1 \\ r_i : 5 & ? & 1 \end{pmatrix})$$

- a. The PLE $\hat{S}_{pl}(t)$ is constant at $[a_{k-1}, a_k), k = 1, ..., m$, 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 ,

$$\Rightarrow \hat{S}_{pl} = 1 - \frac{d_1}{r_1} \ (= \hat{P}(X > a_1)). \qquad (\hat{S}_{pl}(2) = 1 - 1/5 = 4/5).$$

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 $\frac{\hat{S}_{pl}(a_{k-1})}{r_k}$

In Ex1. if
$$Z_{(6)} = 5+$$
 then
$$\begin{pmatrix}
a_i : & 2 & 3 & (3,\infty) \\
d_i : & 1 & 2 \\
r_i : & 5 & 4 \\
S(2-) = & 1 \\
f(2) = & 1/5 \\
S(2) = & 1 - \frac{1}{5} \\
S(3-) = & 4/5 \\
f(3) = & \frac{452}{54} \\
S(3) = & \frac{45(1-24)}{5} \\
S(1) = & \frac{25}{5}
\end{pmatrix}$$

2. The above algorithm results in the expression

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

Note that

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

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

$$\lim_{t \to \infty} \hat{F}_{pl}(t) = \hat{F}_{pl}(Z_{(n)}) < 1 \text{ if } \delta_{(n)} = 0.$$
(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. There are several conventions:

$$\hat{S}_{pl}(t) = 0 \begin{cases} \text{if } t > Z_{(n)} \text{ (convention 1); (is it a survival function ?)} \\ \text{if } t \ge Z_{(n)} \text{ (convention 2); (is it a survival function ?)} \\ \text{if } t \ge Z_{(n)} + c \text{ where } c > 0 \text{ (convention 3). (is it a survival function ?)} \end{cases}$$

However, it can be shown that only the last definition results in an GMLE and has optimal asymptotic properties.

$$\hat{S}_{pl}(t) = \prod_{t \ge Z_{(i)}} (1 - \frac{1}{n - i + 1})^{\delta_{(i)}} = \prod_{k: t \ge a_k} (1 - \frac{d_k}{r_k}).$$

The following table calculates the PLE using the Leukaemia data Group 0 (6-MP): 6+, 6,

	Tal	ole 1. Calculation	ı of PLE	
Remission	Reverse	$\left(1-\frac{1}{n-i+1}\right)^{\delta_{(i)}}$	$\left(1-\frac{d_k}{r_k}\right)$	$\hat{S}_{pl}(a_k)$
Time	Order (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		$\frac{18}{21} \cdot \frac{16}{17}$
10	15	14/15	14/15	$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15}$
10+	14	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15}$
11+	13	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15}$
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	$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{10}{12}$
17+	10	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{10}{12}$
19+	9	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{10}{12}$
20+	8	1		$\frac{18}{21} \cdot \frac{16}{17} \cdot \frac{14}{15} \cdot \frac{10}{12}$
22	7	6/7	6/7	$\frac{6.10.14.16.18}{7.12.15.17.21}$
23	6	5/6	5/6	$\frac{5 \cdot \overline{10} \cdot \overline{14} \cdot \overline{16} \cdot \overline{18}}{7 \cdot 12 \cdot 15 \cdot 17 \cdot 21}$
25+	5	1		$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{7 \cdot 12 \cdot 15 \cdot 17 \cdot 21}$
32+	4	1		
32+	3	1		$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{7 \cdot 12 \cdot 15 \cdot 17 \cdot 21}$
34+	2	1		$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{7 \cdot 12 \cdot 15 \cdot 17 \cdot 21}$
35+	1	1		$\frac{5 \cdot 10 \cdot 14 \cdot 16 \cdot 18}{7 \cdot 12 \cdot 15 \cdot 17 \cdot 21}$

Which of these two is correct ?

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

§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 and use **two** ways!)
- 2. Suppose that $X \sim Bin(3, 1/3)$, $Y \sim Bin(1, 0.4)$. There are two observations (Z_1, δ_1) and (Z_2, δ_2) 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 $S_X(t)$ for t > 1. However, we can estimate S_X for $t \le 1$.

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

Theorem 1. (Johanson (1978)). The PLE \hat{S}_{pl} is a GMLE of S_X (= 1 - F_X), 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 $c_k = \sum_{i=1}^n \mathbf{1}_{\{X_i \neq Z_i = a_k\}} \ge 1$, # 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 $r_k = \sum_{i=1}^n \mathbf{1}_{\{Z_i \ge a_k\}}$. Then $[-\infty, a_1), [a_1, a_1], (a_1, a_2), [a_2, a_2], (a_2, a_3), [a_3, a_3] \text{ and } (a_3, \infty] \text{ is a partition of } [-\infty, +\infty],$ 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{split} \mathbf{L}(F) &= \prod_{i=1}^{n} (f(Z_{i}))^{\delta_{i}} (S(Z_{i}))^{1-\delta_{i}} \\ &= \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}^{a_{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}^{a_{k}})(0 + s_{4} + 0 + s_{6} + s_{7})^{c_{1}}(0 + s_{6} + s_{7})^{c_{2}}(s_{7})^{c_{3}} \text{ (if } p_{1} + p_{2} = s_{2}, \\ &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) \\ &= (\prod_{k=1}^{3} s_{2k}^{a_{k}})(s_{4} + s_{6} + s_{7})^{c_{1}}(s_{6} + s_{7})^{c_{2}}s_{7}^{c_{3}} \\ &= (note \ s_{2} + s_{4} + s_{6} + s_{7} = 1) \\ \mathbf{L}(F) \leq [s_{2}^{d_{1}}(s_{4} + s_{6} + s_{7})^{c_{1}}(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}}(\frac{s_{4}}{s_{4} + s_{6} + s_{7}})^{d_{2}}(\frac{s_{6} + s_{7}}{s_{4} + s_{6} + s_{7}})^{c_{2}}(\frac{s_{6}}{s_{6} + s_{7}})^{d_{3}}(\frac{s_{7}}{s_{6} + s_{7}})^{c_{3}} \\ &= s_{2}^{d_{1}}(1 - s_{2})^{c_{1}}(\frac{s_{4}}{s_{4} + s_{6} + s_{7}})^{d_{2}}(\frac{s_{6} + s_{7}}{s_{4} + s_{6} + s_{7}})^{c_{2}}(\frac{s_{6}}{s_{6} + s_{7}})^{d_{3}}(\frac{s_{7}}{s_{6} + s_{7}})^{c_{3}} \\ &= s_{2}^{d_{1}}(1 - s_{2})^{c_{1}}(\frac{s_{4}}{s_{4} + s_{6} + s_{7}})^{d_{2}}(\frac{s_{6} + s_{7}}{s_{4} + s_{6} + s_{7}})^{c_{2}}(\frac{s_{6}}{s_{6} + s_{7}})^{d_{3}}(\frac{s_{7}}{s_{6} + s_{7}})^{c_{3}} \\ &= s_{2}^{d_{1}}(1 - s_{2})^{c_{1}}(\frac{s_{4}}{s_{4} + s_{6} + s_{7}})^{d_{2}}(\frac{s_{6} + s_{7}}{s_{4} + s_{6} + s_{7}})^{c_{2}}(\frac{s_{6}}{s_{6} + s_{7}})^{d_{3}}(\frac{s_{7}}{s_{6} + s_{7}})^{c_{3}} \\ &\times (s_{4} + s_{6} + s_{7})^{d_{2}+c_{2}+d_{3}+c_{3}}(\frac{s_{6} + s_{7}}{s_{4} + s_{6} + s_{7}})^{d_{3}+c_{3}}} \end{bmatrix}$$

$$=s_{2}^{d_{1}}(1-s_{2})^{c_{1}+d_{2}+c_{2}+d_{3}+c_{3}} \text{ why}? \left(\frac{s_{4}}{1-s_{2}}\right)^{d_{2}} \left(1-\frac{s_{4}}{1-s_{2}}\right)^{c_{2}+d_{3}+c_{3}}$$

$$\times \left(\frac{s_{6}}{1-s_{2}-s_{4}}\right)^{d_{3}} \left(1-\frac{s_{6}}{1-s_{2}-s_{4}}\right)^{r_{3}-d_{3}}$$

$$=s_{2}^{d_{1}}(1-s_{2})^{r_{1}-d_{1}} \left(\frac{s_{4}}{1-s_{2}}\right)^{d_{2}} \left(1-\frac{s_{4}}{1-s_{2}}\right)^{r_{2}-d_{2}}$$

$$\times \left(\frac{s_{6}}{1-s_{2}-s_{4}}\right)^{d_{3}} \left(1-\frac{s_{6}}{1-s_{2}-s_{4}}\right)^{r_{3}-d_{3}} \quad (\text{as } r_{1}=n, \ r_{2}=d_{2}+c_{2}+d_{3}+c_{3}$$

$$\text{and } r_{3}=d_{3}+c_{3}, \text{ where } s_{2}+s_{4}+s_{6}+s_{7}=1 \text{ and } s_{i} \ge 0)$$

$$=(s_{2}^{*})^{d_{1}} \left(1-s_{2}^{*}\right)^{r_{1}-d_{1}} (s_{4}^{*})^{d_{2}} \left(1-s_{4}^{*}\right)^{r_{2}-d_{2}} (s_{6}^{*})^{d_{3}} \left(1-s_{6}^{*}\right)^{r_{3}-d_{3}}; \qquad (2.3)$$

where s_i^* 's are defined in an obvious way, satisfying $s_i^* \in [0, 1]$.

Since the transformation from (s_2, s_4, s_6) to (s_2^*, s_4^*, s_6^*) is one-to-one and onto. Thus, $\mathcal{L}(F)$ is maximized by setting $s_{2i}^* = d_i/r_i$ for i = 1, 2, 3. $step \quad (-\infty, a_1) \quad \{a_1\} \quad (a_1, a_2) \quad \{a_2\} \quad (a_2, a_3) \quad \{a_3\} \quad (a_3, \infty)$

bvcp	(∞, a_1)	lαŢ	(a_1, a_2)	[^a 2]	(a_2, a_3)	լացյ	(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^*	$s_{2k}^*, \ s_7^* \in (0,1)$

Relation between s_i and s_i^* :

$$s_2 = s_2^*,$$

$$\frac{s_4}{1-s_2} = s_4^*,$$

$$\frac{s_6}{1-s_2-s_4} = s_6^*.$$

It follows from the relation between s_i and s_i^* that L(F) is maximized by

$$s_2 = d_1/r_1 = \frac{d_1}{n},$$

$$s_4 = (1 - s_2)s_4^* = \frac{r_1 - d_1}{r_1} \frac{d_2}{r_2},$$

$$s_6 = (1 - s_2 - s_4)s_6^* = \cdots$$

Consequently, $\hat{S}_{pl}(a_1) = \frac{r_1 - d_1}{r_1}$, $\hat{S}_{pl}(a_2) = \frac{r_1 - d_1}{r_1} \frac{r_2 - d_2}{r_2}$, $\hat{S}_{pl}(a_3) = \frac{r_1 - d_1}{r_1} \frac{r_2 - d_2}{r_2} \frac{r_3 - d_3}{r_3}$. $\hat{S}_{pl}(t) = \prod_{t \ge a_k} (1 - \frac{d_k}{r_k}) = \hat{S}_{pl}(a_i)$ if $t \in [a_i, a_{i+1})$, i = 0, ..., 3. $a_0 = ? a_4 = ?$ §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 $d_k + c_k \ge 1$, but $c_k = 0$ is possible now.

§4.2.2. Properties of the PLE \tilde{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 $\tau = \sup\{t: F_Z(t) < 1\}, \text{ where } Z = X \land Y, \text{ and } D_Z^* = \begin{cases} \{t: t \le \tau\} & \text{if } P(X = \tau \le Y) > 0\\ \{t: t < \tau\} & \text{otherwise.} \end{cases}$

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_X}). (2) $G_1(\cdot)$ does not depend on $F_X(\cdot)$. Then $\sup_{t \in D_X^*} |\hat{F}(t) - F_X(t)| \xrightarrow{a.s.} 0$ and $\sup_t |\hat{F}(t) - F_*(t)| \xrightarrow{a.s.} 0$ where

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

Note. $F_X(t), t > \tau$, is not estimable as there is no observation beyond τ , unless $F_X(\tau) = 1$. Examples that G depends on F_X :

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

	t value :	2	3		t value :	2	3	4	
case (1)	$f_{Y X}(1 t) f_{Y X}(2 t)$	$\begin{pmatrix} 1/3\\ 1/3 \end{pmatrix}$	$\frac{1/2}{1/6}$	and case (2)	$f_{Y X}(1 t) f_{Y X}(2 t)$	$\begin{pmatrix} 1/5\\ 1/5 \end{pmatrix}$	$\frac{1}{5}{3}{5}$	$\frac{1/5}{3/5}$	
	$f_{Y X}(3 t)$	$\sqrt{1/3}$	1/3		$f_{Y X}(3 t)$	$\sqrt{3/5}$	1/5	1/5	

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_X is discrete, or S_X 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

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

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,

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

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) $S_X(\tau -) > 0$ or (3) $S_X(\tau -) = 0$, $\tau < \infty$ and $\sigma_\tau = 0$.

$$\sigma_t^2 = \int_0^t \frac{(S_*(t))^2}{S_Z(x-)S_*(x)} dF_*(x) = \begin{cases} (S_*(t))^2 \int_0^t \frac{F'_*(x)}{S_Z(x-)S_*(x)} dx & \text{if } cts \\ (S_*(t))^2 \sum_{x \le t} \frac{S_*(x-)-S_*(x)}{S_Z(x-)S_*(x)} & \text{if } discrete & \text{if } X \perp Y. \\ \dots & \dots & \dots \end{cases}$$

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

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

Note

$$\hat{S}_{pl}(t) = \prod_{i: \ Z_{(i)} \le t} (1 - \frac{1}{n - i + 1})^{\delta_{(i)}} = \prod_{i: \ Z_{(i)} \le t} (1 - \frac{\delta_{(i)}}{n - i + 1}).$$
$$\hat{S}_{Z}(t) = \hat{S}_{Y}(t)\hat{S}_{pl}(t) = \prod_{i: \ Z_{(i)} \le t} (1 - \frac{1}{n - i + 1}) = \sum_{i=1}^{n} \frac{\mathbf{1}_{(Z_{i} > t)}}{n}.$$

The GMLE of S_Y is also a PLE of S_Y .

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

Explain using RTR method.

$$\begin{aligned} \mathbf{Example.} \ n &= 6, \left(\, data: \ 0.5+, \ 2, \ 3, \ 3, \ 3+ \ 5 \, \right) \\ \hat{S}_X(t) &= \prod_{i:Z_{(i)} \leq t} (1 - \frac{\delta_{(i)}}{n - i + 1}) = \begin{cases} 1 & \text{if } t < 2 \\ (1 - \frac{0}{6})(1 - \frac{1}{5}) & \text{if } t \in [2, 3) \\ (1 - \frac{0}{6})(1 - \frac{1}{5})(1 - \frac{1}{4})(1 - \frac{1}{3})(1 - \frac{0}{2}) & \text{if } t \in [3, 5) \\ (1 - \frac{0}{6})(1 - \frac{1}{5})(1 - \frac{1}{4})(1 - \frac{1}{3})(1 - \frac{0}{2})(1 - \frac{1}{1}) & \text{if } t \geq 5 \end{cases} \end{aligned}$$

$$= \begin{cases} 1 & \text{if } t < 2 \\ \frac{4}{5} & \text{if } t \in [2,3) \\ \frac{4}{5}\frac{3}{4}\frac{2}{3} & \text{if } t \in [3,5), \\ 0 & \text{if } t \ge 5 \end{cases}$$
$$\hat{S}_Y(t) = \prod_{i:Z_{(i)} \le t} (1 - \frac{1 - \delta_{(i)}}{n - i + 1})$$

$$= \begin{cases} 1 & \text{if } t < 0.5 \\ (1 - \frac{1}{6})(1 - \frac{0}{5}) & \text{if } t \in [0.5, 2) \\ (1 - \frac{1}{6})(1 - \frac{0}{5})(1 - \frac{0}{4})(1 - \frac{0}{3})(1 - \frac{1}{2}) & \text{if } t \in [3, 5) \\ (1 - \frac{1}{6})(1 - \frac{0}{5})(1 - \frac{0}{4})(1 - \frac{0}{3})(1 - \frac{1}{2})(1 - \frac{0}{1}) & \text{if } t \ge 5 \end{cases} = \begin{cases} 1 & t < 0.5 \\ \frac{5}{6} & \text{if } t \in [\frac{1}{2}, 3), \\ \frac{3}{6}\frac{1}{2} & \text{if } t \in [3, 5) \\ \frac{3}{6}\frac{1}{2} & \text{if } t \in [3, 5) \\ \frac{3}{6}\frac{1}{2} & \text{if } t \ge 5 \end{cases}$$
$$\hat{S}_{Z}(t) = \prod_{i:Z_{(i)} \le t} (1 - \frac{1}{n - i + 1}) = S_{Y}(t)S_{X}(t)$$

$$= \begin{cases} 1 & \text{if } t < 0.5 \\ (1 - \frac{1}{6}) & \text{if } t \in [0.5, 2) \\ (1 - \frac{1}{6})(1 - \frac{1}{5}) & \text{if } t \in [2, 3) \\ (1 - \frac{1}{6})(1 - \frac{1}{5})(1 - \frac{1}{4})(1 - \frac{1}{3})(1 - \frac{1}{2}) & \text{if } t \in [3, 5) \\ (1 - \frac{1}{6})(1 - \frac{1}{5})(1 - \frac{1}{4})(1 - \frac{1}{3})(1 - \frac{1}{2})(1 - \frac{1}{1}) & \text{if } t \ge 5 \end{cases}$$

$$= \begin{cases} 1 & t < 0.5\\ \frac{5}{6} & \text{if } t \in [0.5, 2)\\ \frac{3}{6} \frac{4}{5} & \text{if } t \in [2, 3)\\ \frac{5}{6} \frac{4}{5} \frac{3}{4} \frac{2}{3} \frac{1}{2} & \text{if } t \in [3, 5)\\ 0 & \text{if } t \ge 5 \end{cases} = \begin{cases} 1 & t < 0.5\\ \frac{5}{6} & \text{if } t \in [0.5, 2)\\ \frac{4}{6} & \text{if } t \in [2, 3)\\ \frac{1}{6} & \text{if } t \in [3, 5)\\ 0 & \text{if } t \ge 5 \end{cases}$$

A Simulation Study. Compare two estimators of $S_X(t)$:

(1) the PLE $\hat{S}_{pl}(t)$ and

(2) the incorrect estimate $\tilde{S}(t) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}_{(Z_i > t)}$ (esf).



x = sort(m)

 $\begin{aligned} &\lim_{x \to \infty} (x, 1-\operatorname{ecdf}(x)(x), type="l", lty=3) \\ &\# \operatorname{plot}(m, \operatorname{survfit}(\operatorname{Surv}(m)\sim 1) \\ &(zz=\operatorname{survreg}(\operatorname{Surv}(m,d)\sim 1)) \\ &\lim_{x \to \infty} (x, 1-\operatorname{pweibull}(x, 1/zz \\ scal, \exp(zz \\ coef)), type="l", lty=2) \\ &\lim_{x \to \infty} (x, 1-\operatorname{pweibull}(x, g, b), type="l", lty=4) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pweibull}(x, g, b), type=(1, 2, 3, 4), cex=1.0) \\ &\operatorname{leg.names} (x, 1-\operatorname{pwe$

§4.2.2.2. Homework:

- 1. Suppose that $X \sim Bin(3, 1/3)$, $Y \sim Bin(1, 0.4)$. Under the RC model, what is the limit of $\hat{S}_{pl}(2)$ when $n \to \infty$? What is the value of σ_t^2 with t = 0.5 in Theorem 2? What is the value of $\sigma_{\hat{S}_{pl}}^2(0.5)$ when n = 2? Is it closed to σ_t^2/n above? What does it tell you?
- Suppose X ~ Exp(ρ), where E(X) = 1/ρ, and Y ~ U(0,4). Select a ρ_o. Generate 100 (and then 1000) RC data, plot using these data on the same figure, (1) the PLE Ŝ_{pl}, (2) the true survival function S(·, ρ_o), (3) the S(·, ρ̂), where ρ̂ is the parametric MLE of ρ, (4) the parametric MLE of the survival function of a normal distribution N(μ, 1) using Monte Carlo method (or survreg(,dist="gaussian")) and (5) the parametric MLE of the survival function of a U(θ, 5). (Skip U(θ, 5)) Make comments on their deviations from the true S(t, ρ) at each time point t (separately for n = 100 and n = 1000).
- 3. Using the data in problem #2, we can derive the confidence intervals for $S_X(t)$ based on $\hat{S}_{pl}(t)$ and based on the two MLE estimates of S_X ($Exp(\rho)$ and $N(\mu, 1)$). They are

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

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 SE 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** (or two) confidence bands and the true survival function $S(, \rho)$.

R program for putting the two graphs on one paper.

par(mfrow=c(3,1))x=(-35:45)/10y=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")

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, ..., 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 $(Z_i, \delta_i), i = 1, ..., n$, 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 for a t you select, but $t \in (0, 2)$ with the significance level you chose and the 3 data generated. Notice that your CI is a statistic, depends

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 $a_1 < \cdots < a_m$ be all the **possible** values of $\underline{X \leq \tau}$ and $\tau = a_m$.

Assume that Y can only take values among a_j s.

only on the data.

Let d_k , c_k and r_k be defined as before (but these a_k 's are defined differently from before). ? The likelihood function becomes

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

Then the problem reduces to a parametric problem of multinomial distribution

$$(W_1, ..., W_{2m-1}) \sim M(n, c\theta_1, ..., c\theta_{2m-1}),$$

where for i = 1, ..., m, $W_i = d_i$, $W_{m+i} = c_i$, $\theta_i = p_i$, $\theta_{m+i} = \sum_{j>i} p_j$, $p_i = P\{X = a_i\}$ and $c = 1/\sum_{i=1}^{2m} \theta_i$. Here θ_i 's are function of $\vec{p} = (p_1, ..., p_{m-1})$. The multinomial distribution belongs to the exponential family and its consistency and asymptotic normality can be proved 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.

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

The PLE can be written as

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

Since m is finite,

$$\frac{d_k}{n} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Z_i = X_i = a_k\}}
\rightarrow E(\mathbf{1}_{\{Z = X = a_k\}}) \quad \text{a.s. by SLLN.}$$

$$= P(X = a_k \le Y) = P(X = a_k)P(Y \ge a_k),$$

$$\frac{r_k}{n} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Z_i \ge a_k\}}
\rightarrow E(\mathbf{1}_{\{Z \ge a_k\}}) \text{ a.s. by SLLN.}
= P(Y \ge a_k, X \ge a_k) = P(X \ge a_k)P(Y \ge a_k).$$
(2.4)
$$(2.4)$$

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

by induction on k = 1, ..., m. (Case k = 1). $S_X(t) = 1$ and $\hat{S}_{pl}(t) = 1$ for $t < a_1$. Thus $\hat{S}_{pl}(t) \rightarrow S_X(t)$ for $t \in [a_0, a_1)$ a.s.. (Case k < m). Assume that (2.6) holds. (Case k + 1). For $t \in [a_k, a_{k+1})$,

$$\hat{S}_{pl}(t) = \hat{S}_{pl}(a_k) = \prod_{j: a_k \ge a_j} (1 - \frac{d_j}{r_j})$$

$$= \prod_{j: j < k} (1 - \frac{d_j}{r_j})(1 - \frac{d_k}{r_k})$$

$$= \hat{S}_{pl}(a_{k-1})(1 - \frac{d_k}{r_k})$$

$$= \hat{S}_{pl}(a_{k-1})(1 - \frac{d_k/n}{r_k/n})$$

$$\to S_X(a_{k-1})(1 - \frac{P(X = a_k)P(Y \ge a_k)}{P(X \ge a_k)P(Y \ge a_k)}) \ a.s.$$

(by induction assumption on k and by (2.4) and (2.5))

$$=S_{X}(a_{k-1})(1 - \frac{P(X = a_{k})}{P(X \ge a_{k})})$$

= $P(X > a_{k-1})\frac{P(X > a_{k})}{P(X \ge a_{k})}$
= $P(X > a_{k})$??
= $S_{X}(a_{k}) = S_{X}(t).$ (2.7)

Thus (2.6) holds for k + 1 as well. This completes the induction proof.

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

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

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 $f(a_k) = f_X(a_k) > 0$, k = 1, ..., m + 1, where $f(a_{m+1}) = P(X > a_m)$. Then $\hat{S}_{pl}(t) = \prod_{k: a_k \leq t} (1 - \frac{d_k/n}{r_k/n})$ is a function of $d_1, ..., d_m$ and $r_1, ..., r_m$. Under these assumptions, Theorem 2 becomes:

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

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

$$nCov(\hat{S}_{pl}(t), \hat{S}_{pl}(s)) \approx S_X(t)S_X(s) \int_0^{t \wedge s} \frac{1}{S_X(x-)S_Y(x-)S_X(x)} dF_X(x), \quad t, s \le \tau. \quad \sigma^2(t) = ?$$

In particular, for $t = a_j$, j < m, we can write $\ln \hat{S}_{pl}(t) = g(d_1/n, r_1/n, \dots, d_j/n, r_j/n) = \sum_{i=1}^j \ln(1 - \frac{d_i/n}{r_i/n}).$ $g(\overline{\mathbf{w}}) = \ln(1 - \frac{\overline{w}_1}{\overline{w}_2}) + \dots + \ln(1 - \frac{\overline{w}_{2j-1}}{\overline{w}_{2j}}),$ That is, where $\mathbf{w}^t = (w_1, w_3, w_4, w_5, w_6, ..., w_{2j}),$ as $\overline{w}_2 = r_1/n = 1$ is a constant (note that $\overline{w}_{2m-1}/\overline{w}_{2m} = 1$ if $f(a_{m+1}) = 0$). Write $\overline{\mathbf{W}} = (\overline{W}_1, \overline{W}_3, \overline{W}_4, \dots, \overline{W}_{2i-1}, \overline{W}_{2i})^t = (d_1/n, d_2/n, r_2/n, \dots, d_i/n, r_i/n)^t, \quad as$ $\overline{W}_1 = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(X_i = a_1 \le Y_i),$ $\overline{W}_3 = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(X_i = a_2 \le Y_i),$ $\overline{W}_4 = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(Z_i \ge a_2), \dots$

Then

$$g(\overline{\mathbf{W}}) = \ln \hat{S}_X(t) \text{ and } g(E(\overline{\mathbf{W}})) = \ln S_X(t) \text{ (by (2.7))}.$$

Then $\overline{\mathbf{W}}$ is asymptotical normally distributed and it follows from a corollary of Slutsky's theorem that so is $g(\overline{\mathbf{W}}) = \ln \hat{S}_{pl}(t)$ with an asymptotic variance

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

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

$$\begin{split} \sqrt{n}(g(\overline{\mathbf{W}}) - g(E(\overline{\mathbf{W}}))) & \stackrel{D}{\longrightarrow} N(0, \sigma_{1,t}^2).\\ \sqrt{n}(\ln \hat{S}_{pl}(t) - \ln S_X(t)) & \stackrel{D}{\longrightarrow} N(0, \sigma_{1,t}^2).\\ \sqrt{n}(h(\ln \hat{S}_{pl}(t)) - h(\ln S_X(t))) & \stackrel{D}{\longrightarrow} N(0, \sigma^2(t)) \text{ where } h(t) = e^t.\\ \sqrt{n}(\hat{S}_{pl}(t) - S_X(t)) & \stackrel{D}{\longrightarrow} N(0, \sigma^2(t)). \end{split}$$

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

$$\sigma_2^2/n \approx \sigma_{\hat{S}_{pl}(t)}^2 \approx \frac{\partial e^x}{\partial x} \big|_{x = \ln S_X(t)} \sigma_{\ln \hat{S}_{pl}(t)}^2 \frac{\partial e^x}{\partial x} \big|_{x = \ln S_X(t)} \\ \approx (S_X(t))^2 \sigma_{\ln \hat{S}_{pl}(t)}^2,$$

as $\hat{S}_{pl}(t) = e^{\ln(\hat{S}_{pl}(t))}$. For simplicity, for now, we let j = 2. Then $\mathbf{w}^{t} = (w_{1}, w_{3}, w_{4}),$ $g(\mathbf{w}) = \ln(1 - \frac{w_{1}}{w_{2}}) + \ln(1 - \frac{w_{3}}{w_{4}}), \text{ where } w_{2} = 1.$

$$\begin{aligned} \frac{\partial g}{\partial \mathbf{w}} &= \left(\frac{1}{-w_2(1 - \frac{w_1}{w_2})} \quad \frac{1}{-w_{2j}(1 - \frac{w_{2j-1}}{w_{2j}})} \quad \frac{\frac{w_{2j-1}}{w_{2j}^2}}{(1 - \frac{w_{2j-1}}{w_{2j}})} \right)^t = \left(\frac{\frac{-1}{w_2 - w_1}}{\frac{w_{2j-1}}{w_{2j-1}}} \right) \\ \frac{\partial g}{\partial \mathbf{w}} \Big|_{\mathbf{w} = E(\overline{\mathbf{W}})} &= \left(\frac{-1}{S_X(a_1)S_Y(a_1 -)} \quad \frac{-1}{S_X(a_j)S_Y(a_j -)} \quad \frac{f_o(a_j)}{S_X(a_j)S_Y(a_j -)} \right)^t \text{ (by (2.4) \& (2.5)),} \end{aligned}$$

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

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

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

Note that $s_{ih} = s_{hi}$.

$$s_{11} = E([\mathbf{1}_{(X=Z=a_1)}]^2) - [E(\mathbf{1}_{(X=Z=a_1)})]^2 = f_o(a_1)S_Y(a_1-)(1 - f_o(a_1)S_Y(a_1-))$$

$$\begin{pmatrix} s_{12} & s_{13} \\ s_{22} & s_{23} \\ \vdots & s_{33} \end{pmatrix}$$

$$= \begin{pmatrix} -E(\mathbf{1}_{(X=Z=a_1)})E(\mathbf{1}_{(X=Z=a_2)}) & -E(\mathbf{1}_{(X=Z=a_1)})E(\mathbf{1}_{(Z\geq a_2)}) \\ E(\mathbf{1}_{(X=Z=a_2)}) - [E(\mathbf{1}_{(X=Z=a_2)})]^2 & E(\mathbf{1}_{(X=Z=a_2)}) - E(\mathbf{1}_{(X=Z=a_2)})E(\mathbf{1}_{(Z\geq a_2)}) \\ \vdots & E(\mathbf{1}_{(Z\geq a_2)}) - [E(\mathbf{1}_{(Z\geq a_2)})]^2 \end{pmatrix}$$

$$= \begin{pmatrix} -f_o(a_1)S_Y(a_1-)f_o(a_2)S_Y(a_2-) & -f_o(a_1)S_Y(a_1-)S_X(a_2-)S_Y(a_2-) \\ \vdots & S_X(a_2-)S_Y(a_2-)(1 - S_X(a_2-)S_Y(a_2-)) \end{pmatrix}$$

Then $(p_1 \ p_2 \ p_3) = (1,1) \begin{pmatrix} p_1 & 0 & 0 \\ 0 & p_2 & p_3 \end{pmatrix}$,

$$n\sigma_{\ln\hat{S}_{pl}(t)}^{2} \approx n\mathbf{p}^{t} \sum \mathbf{p} = (1,1) \begin{pmatrix} p_{1} & 0 & 0\\ 0 & p_{2} & p_{3} \end{pmatrix} \begin{pmatrix} s_{11} & s_{12} & s_{13}\\ s_{12} & s_{22} & s_{23}\\ s_{13} & s_{23} & s_{33} \end{pmatrix} \begin{pmatrix} p_{1} & 0\\ 0 & p_{2}\\ 0 & p_{3} \end{pmatrix} \begin{pmatrix} 1\\ 1 \end{pmatrix}.$$

Verify that

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

By symmetry, A' = A = 0 and

$$n\sigma_{\ln\hat{S}_{pl}(t)}^{2} \approx (1,1) \begin{pmatrix} p_{1} & 0 & 0 \\ 0 & p_{2} & p_{3} \end{pmatrix} \begin{pmatrix} s_{11} & 0 & 0 \\ 0 & s_{22} & s_{23} \\ 0 & s_{32} & s_{33} \end{pmatrix} \begin{pmatrix} p_{1} & 0 \\ 0 & p_{2} \\ 0 & p_{3} \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$
$$= (1,1) \begin{pmatrix} p_{1}p_{1}s_{11} & 0 \\ 0 & p_{2}p_{2}s_{22} + p_{2}p_{3}s_{23} + p_{3}p_{2}s_{23} + p_{3}p_{3}s_{33} \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$
$$= p_{1}p_{1}s_{11} + p_{2}p_{2}s_{22} + p_{2}p_{3}s_{23} + p_{3}p_{2}s_{23} + p_{3}p_{3}s_{33} \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$
(2.9)

Then it yields

$$\sigma_{\ln \hat{S}_{pl}(t)}^2 \approx \frac{1}{n} \sum_{a_k \le t} \frac{f_o(a_k)}{S_X(a_k) - S_Y(a_k) - S_X(a_k)}.$$
 (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 7 84-104.

$\S4.3.$ 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

$$(Y, \delta)$$
, where $\delta = \mathbf{1}_{(X \le Y)}$. $(vs. (X \land Y, \delta))$

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

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

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

(m < n or m > n ?)

Let
$$N_{-}(a_j) = \sum_{i=1}^{n} \mathbf{1}_{(X_i \le Y_i = a_j)}$$

$$N_{+}(a_{j}) = \sum_{i=1}^{n} \mathbf{1}_{(X_{i} > Y_{i} = a_{j})},$$
$$N(a_{j}) = \sum_{i=1}^{n} \mathbf{1}_{(Y_{i} = a_{j})}.$$

The likelihood function is

$$\begin{aligned} \mathbf{L}(F) &= \prod_{i=1}^{n} \mu_{F}(I_{i}) \\ &= \prod_{j=1}^{m} [(F(X_{j}))^{N_{-}(a_{j})}(S(A_{j}))^{N_{+}(a_{j})}], \ F \in \Theta, \ where \end{aligned}$$

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

$$0 \le F(a_1) \le F(a_2) \le \dots \le F(a_m) \le 1.$$

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

$$\mathcal{L}(F) = \sum_{j=1}^{m} [N_{-}(a_{j}) \log s_{j} + N_{+}(a_{j}) \log(1 - s_{j})], \quad 0 \le s_{1} \le \dots \le s_{m} \le 1.$$

$$\frac{\partial \mathcal{L}}{\partial s_{j}} = \frac{N_{-}(a_{j})}{s_{j}} - \frac{N_{+}(a_{j})}{1 - s_{j}} = 0$$

$$\Rightarrow s_{j} = \tilde{F}(a_{j}) = \frac{N_{-}(a_{j})}{N(a_{j})}.$$

$$(3.2)$$

If $\tilde{F} \in \Theta$, or $0 \le s_1 \le \cdots \le s_m \le 1$, it is the GMLE (in fact,

each summand looks like the log likelihood of a binomial distribution $Bin(N(a_j), s_j))$. Otherwise, it is not and the GMLE will be on the boundaries $0 \le s_i = s_j \le 1, i \ne j$. **Example 1**. The observations (Y_i, δ_i) s are (1, 1), (2, 0), (2, 1), (3, 1). n = 4. $a_1 = 1, a_2 = 2$ and $a_3 = 3, m = 3$.

$$\begin{aligned}
u_1 &= 1, \ u_2 &= 2 \text{ and } u_3 &= 5, \ m = 6, \\
N_-(a_1) &= 1, \ N_+(a_1) &= 0, \\
N_-(a_2) &= 1, \ N_+(a_2) &= 1, \\
N_-(a_3) &= 1, \ N_+(a_3) &= 0, \\
\\
\text{Then } \tilde{F}(t) &= \begin{cases}
0 & \text{if } t < 1 \\
1/1 & \text{if } t \in [1, 2) \\
\frac{1}{2} & \text{if } t \in [2, 3) \\
1/1 & \text{if } t \geq 3.
\end{aligned}$$

 $F \notin \Theta$. Thus \tilde{F} is not a GMLE of F. GMLE = ? (homework). **Theorem 1.** A GMLE of F_X under C1 model is

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

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

$$\hat{F}(Y_{(j)}) = \max_{i \le j} \min_{k \ge j} \frac{\sum_{i \le h \le k} \delta_{(h)}}{k - i + 1}, \ j = 1, ..., n,$$
(3.3)

and $\delta_{(h)}$ is the δ_i 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 $i \ge 0$, unless $\hat{F}(a_i) = \hat{F}(a_{i+1})$.
- 3. It can be shown that (homework) (3.3) is the same as

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

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

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

 $\max_{i < 4} m_i$

Example 2. Find the GMLE based on 4 C1 observations, (1,0), (2,1), (3,0) and (3,1). **Sol.** Ordered data: $\begin{pmatrix} 1,0 \\ Y_{(1)} \\ Y_{(2)} \\ Y_{(3)} \\ Y_{(3)} \\ Y_{(4)} \\ Y_{(4)} \\ Y_{(3)}, Y_{(4)} \end{pmatrix} = [3,3) = ?$ It is more convenient to use formula (3.5) (or (3.6)) which results in the following matrix. Let $A_{ik} = \frac{\delta_{(i)} + \dots + \delta_{(k)}}{k - i + 1}$.

	δ :	0	1	1	0		Eq.(3.5)
	iackslash k	1	2	3	4	$\min_k A_{ik}$	$\hat{F}(Y_{(i)})$
	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
	$M_k = \max_i A_{ik}$	0	1	1	2/3		
Eq. (3.6)	$\hat{F}(Y_{(k)}) = \min_{j > k} M_j$	0	2/3	2/3	2/3		

The GMLE is uniquely defined at Y_i 's *i.e.*, 1, 2 and 3, as 0, 2/3 and 2/3, and thus the GMLE is uniquely determined on the set $(-\infty, 1]$ and [2, 3], $\hat{F}(t) = 0$ and 2/3, respectively, and it arbitrary on $(1, 2) \cup (3, \infty)$.

The GMLE of
$$F(t) = \begin{cases} 1 & \text{if } t = \infty \\ \uparrow & \text{if } t \in (3, \infty) \\ 2/3 & \text{if } t \in [2, 3] \\ \uparrow & \text{if } t \in (1, 2) \\ 0 & \text{if } t \leq 1 \end{cases}$$
 A GMLE is $\hat{F}(t) = \begin{cases} 1 & \text{if } t = \infty \\ 2/3 & \text{if } t \in [2, \infty) \\ 0 & \text{if } t < 2 \end{cases}$.

Skip this page.

$$\hat{F}(Y_{(j)}) = \max_{i \le j} \min_{k \ge j} \frac{\sum_{i \le h \le k} \delta_{(h)}}{k - i + 1}, \ j = 1, ..., n,$$
(3.3)

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))
 $r[i] = max(apply(mat[1:i, i:L], 1, min))$
 $cbind(pt, r)$
}

 $\int c1(x)$

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 $p_1, ..., p_4$. Then the likelihood function

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

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

Eq.(3.5) yields $\hat{F}(t) = 0$, 2/3 and 2/3 at 1, 2, 3, respectively. The distinct observations Y_i 's, 1,2,3, partition $(-\infty,\infty)$ as 4 disjoint intervals. $(-\infty, (1, 2])((3, \infty))$ The innermost intervals are (1, 2] and $(3, \infty)$. The GMLE \hat{F} assigns weights 2/3 and 1/3 to them.

Obtain the GMLE directly from L in (3.1):

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

 $\mathcal{L} = \ln(1 - s_1) + \ln s_2 + \ln s_3 + \ln(1 - s_3), \ 0 \le s_1 \le s_2 \le s_3 \le 1$ Skip next page.

Remark: 9 Boundaries, but they can simplified:

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

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

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

$$\mathbf{s}$$
: critical pt $(0, 1, 1/2)$ s

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

(2)
$$s_1 = 0$$
: \mathcal{L} : violating the constraint no need to check

(3)
$$s_2 = s_3$$
:

$$\begin{aligned} \mathbf{s}: \quad critical \ point \ (0, 2/3, 2/3) \quad s_1 = 0 \quad s_3 = 1 \quad s_1 = s_2 = s_3 \\ \mathcal{L}: \quad \ln(\frac{2^2}{3^3}) \quad \ln(\frac{2^2}{3^3}) \quad -\infty \quad \ln\frac{1}{2^4} \\ (4) \ s_1 = s_2: \quad \mathbf{s}: \quad (1/2, 1/2, 1/2) \quad s_1 = 0 \quad s_3 = 1 \\ \mathcal{L}: \quad \ln\frac{1}{2^4} \quad -\infty \quad -\infty \end{aligned}$$

(5) $s_3 = 1$: ...

Other boundaries can be skipped (see the above remark). $\mathbf{GMLE} = ?$

§4.3.1.2. Homework.

- 1. Derive the GMLE of F using data from Example 1. Use three approaches: (1) use formula (3.5); (2) derive directly from the likelihood function in (3.1) (see Remark: 9) boundaries). (3) Use R codes.
- 2. Assuming Y is discrete with finitely many values, show that the estimator \tilde{F} in (3.2) is a consistent estimator of $F_X(a_i)$ and find its asymptotic variance.
- 3. Show that the definitions (3.4) of $\hat{F}(Y_{(i)})$ is equivalent to (3.3). Hint: Inspect the matrix $(A_{ik})_{j \times (n-j+1)}$, where $A_{ik} = \frac{\sum_{i \le h \le k} \delta_{(h)}}{k-i+1}$. Try first n = 3 or 4 and use induction argument on j.

Consistency of the GMLE is only relevant for $t \leq \tau$, where $\tau_G = \sup\{t : G(t) < 1\}$, $Y \sim G$ and $\tau = \tau_G$. 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).

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

 $\sup_{x < \tau} |\hat{F}(x) - F_X(x)| \xrightarrow{a.s.} 0$ in the C1 model, if F_X is continuous and G'(x) > 0 on $[0, \tau]$. Example 3 after the next theorem shows that this is not true.

Theorem 2. Under the C1 model, the GMLE F satisfies

(1) $F(a) \to F_X(a)$ a.s. for each $a \in \mathcal{A}$, where $\mathcal{A} = \{a : P(Y = a) > 0\}$ (Yu et al. (1998)); (2) $\lim_{n\to\infty} \int |\hat{F}(t) - F_X(t)| dG(t) = 0$ a.s. (Schick and Yu (2000));

(3) If F_X is continuous in $(0,\tau]$, $P\{Y = \tau\} > 0$ or $F_X(\tau) = 1$, and the range of Y is dense in $[0,\tau]$, then $\sup_{x \le \tau} |\hat{F}(x) - F_X(x)| \xrightarrow{a.s.} 0$, i.e., \hat{F} is uniformly strongly consistent on $(-\infty, \tau]$ (Schick and $\overline{Y}u$ (2000)).

Example 3. Consider C1 data $(Y_1, \mathbf{1}_{(X_1 \leq Y_1)}), \ldots, (Y_n, \mathbf{1}_{(X_n \leq Y_n)})$, where the survival times X_1, \ldots, X_n are i.i.d. $\sim U(0,3)$ and the inspection times Y_1, \ldots, Y_n are i.i.d. $\sim U(0, 2)$. $\tau = ?$ \vdash : $\hat{F}_n(2) = 1 > 2/3 = F_X(2)$ on an event *B* with *P*(*B*) > 1/6. Let $B = \bigcup_{j=1}^{n} B_j$, where $B_j = \{X_j \le 1 \le Y_j, Y_j > Y_i, i = 1, \dots, n, i \ne j\}$ (with $Y_j = Y_{(n)}$) and $\delta_j = 1$). On the event $\bigcup_{i=1}^n B_i$, we have $\hat{F}_n(2) = \hat{F}_n(2-) = \hat{F}(Y_{(n)}) = 1$ (as $Y_{(n)} < 2$,

$$\hat{F}(Y_{(n)}) = \max_{i \le n} \min_{k \ge n} \frac{\sum_{i \le h \le k} \delta_{(h)}}{k - i + 1} = \max_{i \le n} \frac{\sum_{i \le h \le n} \delta_{(h)}}{n - i + 1} = \frac{\sum_{n \le h \le n} \delta_{(h)}}{n - n + 1} = 1 \ (see \ (3.3))$$

(as $Y_j = Y_{(n)}$ and $\delta_j = 1$). The event has probability $1/3 - \frac{1}{3 \cdot 2^n} \ge \frac{1}{6}$, as it equals

$$P\{\bigcup_{j=1}^{n} B_j\} = nP\{X_1 \le 1 \le Y_1, Y_1 > Y_j, \ j = 2, ..., n\}$$
(*B*:s are disjoint events)

 $(B_i s \text{ are disjoint events with the same probability})$

$$=nP\{X_{1} \leq 1\}P\{1 \leq Y_{1}, Y_{1} > Y_{j}, \ j = 2, ..., n\} ?$$

$$=nP\{X_{1} \leq 1\}E(\mathbf{1}(1 \leq Y_{1}, Y_{1} > Y_{j}, \ j = 2, ..., n))$$

$$=n(1/3)E(E(\mathbf{1}(1 \leq Y_{1}, Y_{1} > Y_{j}, \ j = 2, ..., n)|Y_{1}))$$

$$=n(1/3)\int_{1}^{2}E(\mathbf{1}(1 \leq y_{1}, y_{1} > Y_{j}, \ j = 2, ..., n)|Y_{1} = y_{1})\frac{1}{2}dy_{1}$$

$$=\frac{n}{3}\int_{1}^{2}\frac{1}{2}P\{Y_{j} < y_{1}, j = 2, ..., n\}dy_{1} = \frac{n}{3}\int_{1}^{2}\frac{1}{2}(\frac{y_{1}}{2})^{n-1}dy_{1}$$

$$=\frac{n}{3}\frac{1}{n}(\frac{y_{1}}{2})^{n}|_{1}^{2}$$

$$=\frac{1}{3}(1-2^{-n}).$$

That is $P(\hat{F}(2-) = 1) \ge P(\bigcup_j B_j) = \frac{1}{3} - \frac{1}{3 \times 2^n}.$

Since $F_X(2) = F_X(2-) = 2/3$, we see that the following two statements are false: $\hat{F}(2-)$ converges to $F_X(2-)$ a.s. $(P(\hat{F}(2-) \rightarrow F_X(2-)) = 1)$

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

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

Remark. In Example 3 the GMLE $\hat{F}(t)$ is consistent for all $t < \tau$ (due to Proposition 3.2) in Schick and Yu (2000)), but not at τ or τ -.

Question: Is $\hat{F}(t)$ consistent for all $t < \tau$ under the C1 model ? **Example 4.** Suppose that $X \sim Exp(1)$ and Y has a Poisson distribution with mean 1. Then $\tau = ?$

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

 $\lim_{n\to\infty} \hat{S}(1.5) = S_X(1.5)$ a.s. $\hat{P}(S_X(1.5) = e^{-1.5})$

Answer: $\lim_{n \to \infty} \int |\hat{F}(t) - F_X(t)| dG(t) = 0 \text{ a.s. (Theorem 2).}$

 $\lim_{n \to \infty} \hat{S}(1) = S_X(1)$ a.s. !

 $\lim_{n \to \infty} \hat{S}(1.5) \neq S_X(1.5)$ a.s..

 $\lim_{n \to \infty} \hat{S}(1.5) = S_X(1)$ 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_X 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 even if F_X and G satisfy certain smoothness assumptions. Indeed, their Theorem 3 below establishes that under differentiability assumptions on F_X 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 $0 < F_X(t_o), G(t_o) < 1$, and let F_X and G be differentiable at t_o , with strictly positive derivatives $f_o(t_o)$ and $g(t_o)$, respectively. Then

$$n^{1/3} \frac{\hat{F}(t_o) - F_X(t_o)}{\{\frac{1}{2}F_X(t_o)S_X(t_o)f_o(t_o)/g(t_o)\}^{1/3}} \xrightarrow{D} 2Z, \ as \ n \to \infty$$

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 **Guassian process** if each finite-dimensional vector $(W(t_1), ..., W(t_m)) \sim N(\vec{\mu}(\mathbf{t}), \Sigma(\mathbf{t}))$, where $\Sigma(\mathbf{t})$ can be singular and $\mathbf{t} = (t_1, ..., t_m)$, $m \geq 1$; it is called a **Wiener process**, or a **Brownian Motion**, if W(t) has independent increments and $W(s+t) - W(s) \sim N(0, \sigma^2 t)$, $\forall t > 0$, **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 a given set $\mathcal{A} (= \{a : P(Y = a) > 0\})$. For $x \in (-\infty, +\infty)$, let

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

We say x is a regular point, if

- (1) x belongs to \mathcal{A} ,
- (2) x_{-} and x_{+} belong to \mathcal{A}_{*} ,
- (3) $x_{-} < x < x_{+}$ and
- (4) $F_X(x_-) < F_X(x) < F_X(x_+).$

Example 5. What are the regular points in these sets:

5.1. \mathcal{A}_1 = the support set of the Poisson distribution. 5.2. \mathcal{A}_2 = the set of all positive fractions and $F_{22}(t) = 1$

5.2. \mathcal{A}_2 = the set of all positive fractions and $F_X(t) = 1 - e^{-t}, t > 0.$

5.3. $\mathcal{A}_3 = \{0, 1, 2, 3\} \cup [4, 5] \cup (6, 8], F_X \text{ is } U(2, 10) \text{ distribution.}$

Theorem 4. (Yu et al. (1998)). Let x be a regular point. Then $n^{1/2}(\hat{F}(x) - F_X(x))$ is asymptotically normal with mean 0 and variance $F_X(x)(1 - F_X(x))/g(x)$. This asymptotic variance can be consistently estimated by $\hat{F}(x)(1 - \hat{F}(x))/N(x)$. Also, if $x_1 < \ldots < x_m$ are regular points, then $n^{1/2}(\hat{F}(x_1) - F_X(x_1), \ldots, \hat{F}(x_m) - F_X(x_m))$ 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_X is strictly monotone and Y takes on finitely many values, say at $a_1 < \cdots < a_m$, with $F(a_1) > 0$ and $F(a_m) < 1$. For n large enough,

$$\hat{F}(a_i) = \tilde{F}(a_i) \text{ (see (3.2))}, \ i = 1, ..., m - 1.$$

Reason: Let $\epsilon = \min\{F_X(a_1), F_X(a_i) - F_X(a_{i-1}), i = 2, ..., m, S_X(a_m)\}$. Then $\epsilon > 0$. Let Ω be the event that $N_-(a_j)/n$ and $N(a_j)/n$ (see (3.2)) converge for j = 1, ..., m. Then for each $\omega \in \Omega$, for n large enough, $|\tilde{F}(a_j) - F_X(a_j)| < \epsilon/3$ for all j. It follows that

$$0 \le \tilde{F}(a_1) < F_X(a_1) + \epsilon/3 < F_X(a_2) - \epsilon/3 < \tilde{F}(a_2) < \cdots \tilde{F}(a_m) \le 1.$$
(3.7)

As a consequence, \tilde{F} is a GMLE.

In fact, the asymptotic covariance matrix of $(\hat{F}(a_1), ..., \hat{F}(a_m))$ can be estimated by the expression

$$\hat{\Sigma} = \left(-\frac{\partial^2 \mathcal{L}}{\partial \mathbf{s} \partial \mathbf{s}^t} \Big|_{\mathbf{s}^t = (\hat{F}(a_1), \dots, \hat{F}(a_m))} \right)^{-1}.$$
(3.8)

Now $\mathcal{L} = \ln \prod_{i=1}^{m} s_i^{N_-(a_i)} (1 - s_i)^{N_+(a_i)}.$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{s}^{t}} &= \left(\frac{N_{-}(a_{1})}{s_{1}} - \frac{N_{+}(a_{1})}{1 - s_{1}}, \dots, \frac{N_{-}(a_{m})}{s_{m}} - \frac{N_{+}(a_{m})}{1 - s_{m}}\right).\\ \frac{\partial^{2}\mathcal{L}}{\partial \mathbf{s}\partial \mathbf{s}^{t}} &= -\left(\frac{\frac{N_{-}(a_{1})}{s_{1}^{2}} + \frac{N_{+}(a_{1})}{(1 - s_{1})^{2}} & 0 & \cdot & 0}{0 & \cdots & \frac{N_{-}(a_{m})}{s_{m}^{2}} + \frac{N_{+}(a_{m})}{(1 - s_{m})^{2}}}\right).\\ \frac{1}{n} \left(-\frac{\partial^{2}\mathcal{L}}{\partial \mathbf{s}\partial \mathbf{s}^{t}}\Big|_{\mathbf{s}^{t} = (\hat{F}(a_{1}), \dots, \hat{F}(a_{m}))}\right) & (N_{-}(a_{j}) = \sum_{i=1}^{n} \mathbf{1}(X_{i} \le a_{j} = Y_{i}))\\ \rightarrow \left(\frac{F_{X}(a_{1})g(a_{1})}{F_{X}^{2}(a_{1})} + \frac{S_{X}(a_{1})g(a_{1})}{(S_{X}(a_{1}))^{2}} & 0 & \cdot & 0\\ 0 & 0 & \cdots & \frac{F_{X}(a_{m})g(a_{m})}{F_{X}^{2}(a_{m})} + \frac{S_{X}(a_{m})g(a_{m})}{(S_{X}(a_{m}))^{2}}\right) \text{ a.s.}\\ = \left(\frac{g(a_{1})}{F_{X}(a_{1})} + \frac{g(a_{1})}{S_{X}(a_{1})} & 0 & \cdot & 0\\ 0 & 0 & \cdots & \frac{g(a_{m})}{F_{X}(a_{m})} + \frac{g(a_{m})}{S_{X}(a_{m})}\right)\\ = \left(\frac{g(a_{1})}{F_{X}(a_{1})S_{X}(a_{1})} & 0 & \cdot & 0\\ 0 & 0 & \cdots & \frac{g(a_{m})}{F_{X}(a_{m})S_{X}(a_{m})}\right)
\end{aligned}$$

$$n\hat{\Sigma} \to \begin{pmatrix} \frac{F_X(a_1)S_X(a_1)}{g(a_1)} & 0 & \cdot & 0\\ \cdot & \cdot & \cdot & \cdot\\ 0 & 0 & \cdots & \frac{F_X(a_m)S_X(a_m)}{g(a_m)} \end{pmatrix} a.s.$$

Theorem 4 is an extension of the above observation.

However, if F_X is not strictly increasing on $\{a_1, ..., a_m\}$, it is not true that

$$F(a_i) = F(a_i), i = 2, ..., m - 1$$
, if n is large enough,

as (3.7) does not hold. In such a case, (3.8) 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),

$$n^{1/3}(\hat{F}(t_o) - F_X(t_o)) \xrightarrow{D} hZ,$$

where

$$h = \{4F_X(t_o)S_X(t_o)f_o(t_o)/g(t_o)\}^{1/3}$$

Thus, a 95% CI for $F_X(t_o)$ 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.,

$$\hat{h} = \{4\hat{F}(t_o)\hat{S}(t_o)\hat{f}(t_o)/\hat{g}(t_o)\}^{1/3}$$
$$\hat{f}(t) = \int \frac{1}{w_n} K(\frac{x-t}{w_n}) d\hat{F}(x),$$

 w_n is the window width, $w_n \to 0$, $K(\cdot)$ is a kernel $(e.g., K(x) = \frac{1}{2} \mathbf{1}_{(-1 < x \leq 1)})$ and thus

$$\hat{f}(t) = \frac{\hat{F}(t+w_n) - \hat{F}(t-w_n)}{2w_n},$$

$$\begin{split} \hat{g}(t) &= \int \frac{1}{w_n} K(\frac{x-t}{w_n}) d\hat{G}(x), \end{split} \tag{Lebesgue Stieltjes} \\ \hat{G}(x) &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(Y_i \leq x)}. \end{split}$$

§4.3.1.3. Homework.

4. Let F_X 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 $\hat{F}(1/\rho)$ 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 $\hat{F}(3)$ 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 $\text{Exp}(\rho)$. Plot the survival curves of $S(t; \rho_o)$, the MLE and the GMLE on the same graph. Now pretend the data (Y_i, δ_i) 's are 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 intervalcensored 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.

\S **4.4. Self-consistent (SC) algorithm.**

Under the RC model, C1 model and LC model, the GMLE of F_X 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 $I_1, ..., I_n$ denote the observed intervals.

Define innermost intervals A_j , j = 1, 2, ..., m, induced by $I_1, ..., I_n$ 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

$$A_j \cap I_i = \emptyset \text{ or } A_j \ \forall i \text{ and } j.$$

Let the endpoints of the innermost intervals be a_j and b_j , j = 1, ..., m, where

 $a_1 \le b_1 \le a_2 \le b_2 \le \dots \le a_m \le b_m.$

The following example illustrates the procedure for finding innermost intervals.

Example 2.1. Suppose that there are five observed intervals I_i 's:

 $\begin{array}{l} (1,4], [2,2], (2,6], [5,5], \text{ and } (1,6]. \\ \text{Then there are two exact observations,} \\ I_2 = [2,2] \text{ and } I_4 = [5,5], \text{ and} \\ \text{three censored intervals,} \\ I_1 = (1,4], I_3 = (2,6] \text{ and } I_5 = (1,6]. \\ \text{Furthermore, there are three innermost intervals,} \\ A_1 = [2,2], A_3 = [5,5], \text{ and } A_2 = (2,4]. \\ \begin{array}{c} 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ ((& \cdot(&] & \cdot &]] \\ ((& [](&] & [] &]] \\ s_1 & s_2 & s_3 \end{array} \qquad (\delta_{ij})_{5\times 3} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \end{pmatrix}$

Peto (1973) shows that the GMLE of F_X only

assigns weights, say $s_1, ..., s_m$, to the corresponding innermost intervals $A_1, ..., A_m$. The generalized likelihood function L can be simplified as

$$\mathbf{L} = \mathbf{L}(s_1, ..., s_m) = \prod_{i=1}^{n} [\sum_j \delta_{ij} s_j], \qquad (4.1)$$

where $\delta_{ij} = \mathbf{1}_{(A_j \subset I_i)}$, $\mathbf{s}^t = (s_1, \dots, s_{m-1}) \in D_s$ is the transpose of \mathbf{s} , $D_s = \{\mathbf{s}: s_i \ge 0, s_1 + \dots + s_{m-1} \le 1\}$ 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 $j = 1, ..., m$.
At step $h, s_j^{(h)} = \sum_{i=1}^n \frac{1}{n} \frac{\delta_{ij} s_j^{(h-1)}}{\sum_{k=1}^m \delta_{ik} s_k^{(h-1)}}, \quad j = 1, ..., m, \ h \ge 2.$

Stop when $\mathbf{s}^{(h)}$ converges, *i.e.*, $||\mathbf{s}^{(h)} - \mathbf{s}^{(h-1)}||$ is small enough. He shows that, as $h \to \infty$, $s_j^{(h)}$ converges to the GMLE, \hat{s}_j , which maximizes L and satisfies

the system of self-consistent equations

$$s_j = \sum_{i=1}^n \frac{1}{n} \frac{\delta_{ij} s_j}{\sum_{k=1}^m \delta_{ik} s_k}, \quad j = 1, ..., m, \mathbf{s} \in D_s.$$
(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

 $\hat{F}(b_i) = \hat{F}(a_{i+1}-) = \hat{s}_1 + \dots + \hat{s}_i,$

but is not uniquely defined for t being in a non-singleton innermost interval

(Peto, 1973; Turnbull, 1976).

To avoid this ambiguity we define

$$\hat{F}(t) = \sum_{j: A_j \subset (-\infty, t]} \hat{s}_j = \sum_{b_j \le t} \hat{f}(b_j).$$
(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 2.1** (continued).

$$\begin{split} s_1^{(1)} &= s_2^{(1)} = s_3^{(1)} = 1/3. \\ s_1^{(h)} &= \frac{1}{5} [\frac{s_1^{(h-1)}}{s_1^{(h-1)} + s_2^{(h-1)}} + \frac{s_1^{(h-1)}}{s_1^{(h-1)}} + \frac{s_1^{(h-1)}}{s_1^{(h-1)} + s_2^{(h-1)} + s_3^{((h-1))}}], \\ s_2^{(h)} &= \frac{1}{5} [\frac{s_2^{(h-1)}}{s_1^{(h-1)} + s_2^{(h-1)}} + \frac{s_2^{(h-1)}}{s_2^{(h-1)} + s_3^{(h-1)}} + \frac{s_2^{(h-1)}}{s_1^{(h-1)} + s_2^{(h-1)} + s_3^{((h-1))}}], \\ s_3^{(h)} &= 1 - s_1^{(h)} - s_2^{(h)}, \ h \ge 2 \end{split}$$

Remark. An SCE of **s** can be viewed as

a critical point of L(s), subject to the constraint on s.

The reason is as follows.

Given $j \in \{1, ..., m\}$, let $\mathbf{s}(\epsilon) = (s_1(\epsilon), ..., s_{m-1}(\epsilon))$ be defined by

$$s_k(\epsilon) = \begin{cases} \frac{s_k}{1+\epsilon} & \text{if } k \neq j \\ \frac{s_j+\epsilon}{1+\epsilon} & \text{if } k = j, \end{cases} \qquad \qquad = \begin{cases} \frac{s_k}{1+\epsilon} & \text{if } k \neq j \\ \frac{s_j-1}{1+\epsilon} + 1 & \text{if } k = j, \end{cases}$$

Write

$$\Lambda_{j}(\epsilon) = \ln L(\mathbf{s}(\epsilon)) = \sum_{i=1}^{n} \ln \sum_{k=1}^{m} \delta_{ik} s_{k}(\epsilon) \qquad =>$$

$$\frac{\partial \Lambda_{j}(\epsilon)}{\partial \epsilon} \Big|_{\epsilon=0} = -n + \sum_{i=1}^{n} \frac{\delta_{ij}}{\sum_{k=1}^{m} \delta_{ik} s_{k}}. \qquad \left(= \sum_{i=1}^{n} \frac{\sum_{k=1}^{m} \delta_{ik} \frac{-s_{k} \mathbf{1}(k \neq j) - (s_{j} - 1) \mathbf{1}(k = j)}{(1 + \epsilon)^{2}} \Big|_{\epsilon=0}\right)$$

If s is the GMLE, then the value $\epsilon = 0$ maximizes $\Lambda_j(\epsilon)$ for each j subject to $\epsilon \ge 0$. In other words, we have

$$\left. \frac{\partial \Lambda_j(\epsilon)}{\partial \epsilon} \right|_{\epsilon=0} = -n(1 - \frac{1}{n} \sum_{i=1}^n \frac{\delta_{ij}}{\sum_{k=1}^m \delta_{ik} s_k}) \le 0$$

with equality

$$1 = \frac{1}{n} \sum_{i=1}^{n} \frac{\delta_{ij}}{\sum_{k=1}^{m} \delta_{ik} s_k}$$
 (compare to (4.2)) $s_j = \sum_{i=1}^{n} \frac{1}{n} \frac{\delta_{ij} s_j}{\sum_{k=1}^{m} \delta_{ik} s_k}$

unless $s_j = 0$. Consequently, (4.2) holds. That is, the GMLE satisfies (4.2). Of course the fact that (4.2) holds for s does not implies that s is the GMLE, but it implies that s is a critical point of the likelihood.

```
library(survival)
rm(list=ls())
p = 1/2
n = 100
gmle1=rep(0,10)
for(i in 1:10){
    x = rexp(n,p) \# survival time
    y_1 = rexp(n,p) \# the next 3 lines are for Case 2 censoring pattern.
    z = rexp(n,p)
    y2 = y1+z
    status = 2^{as.numeric}(x \le y1) + 3^{as.numeric}(y1 \le x)^{as.numeric}(x \le y2)
    +0^{*}as.numeric(v2<x) #
                                  3 lines for (L,R,status) used in Surv()
    t1 = 0^{*}(status = 2) + y1^{*}(status = 3) + y2^{*}(status = 0) \# L
    t2 = y1^{*}(status = 2) + y2^{*}(status = 3) + 100^{*}(status = 0) \# R
\# next two lines are SCE codes
    my.surv = Surv(time=t1,time2=t2,event=status,type="interval")
                                          \# works for PLE and self-consistent estimator
    temp = survfit(my.surv \sim 1)
    tt=summary(temp)
                                           \#SCE of the i-th loop.
                                                        # = \hat{S}(2)? or \hat{S}(2-)?
    gmle1[i]=ttsurv[tttime>=2][1]
}
                    #SCE of last loop.
\operatorname{tt}
                                                       lower95%CI
     time
            n.risk
                      n.event
                                  survival
                                              std.err
                                                                       upper 95\% CI
     1.04
            44.000
                     4.69e + 00
                                   0.4645
                                              0.0508
                                                           0.3750
                                                                            0.575
     1.25
                     1.75e - 01
            35.306
                                   0.4622
                                              0.0508
                                                           0.3726
                                                                            0.573
     1.50
            33.131
                     1.43e - 03
                                   0.4622
                                              0.0508
                                                           0.3726
                                                                           0.573
     1.72
            33.130
                     3.47e + 00
                                              0.0517
                                                           0.3239
                                   0.4138
                                                                           0.529
     2.18
            24.660
                     7.36e + 00
                                   0.2903
                                              0.0526
                                                           0.2035
                                                                           0.414
     2.32
            16.299
                     9.54e - 02
                                   0.2886
                                              0.0526
                                                           0.2019
                                                                           0.413
     3.04
                     1.18e + 01
                                                           0.0233
            15.203
                                   0.0639
                                              0.0329
                                                                           0.175
     4.21
                     3.74e - 08
             2.369
                                   0.0639
                                              0.0329
                                                           0.0233
                                                                            0.175
     8.60
                     3.69e - 01
                                              NaN
                                                            NA
                                                                            NA
             0.369
                                   0.0000
                                                            dis
                                                                          normal
(s100=sd(gmle1))
```

[1] 0.07321893

Remark. In the above codes,

$$(t1, t2, event) = \begin{cases} (0, t2, 2) & \text{if LC at } t_2 \\ (t1, t2, 3) & \text{if SIC in } (t1, t2] \\ (t1, 100, 0) & \text{if RC at } t1 \text{ (assuming } t1 < 100, \text{ OW replace } 100 \text{ by } 10^9 \end{cases}$$

Another way:

> l = c(-Inf,2,4,3,9)> r = c(5,2,7,8,Inf) > require(interval) # SC algorithm codes

> fit1=icfit(l,r)
> summary(fit1)

Interval Probability

 $\begin{array}{ccccccc} 1 & [2,2] & 0.2667 \\ 2 & (4,5] & 0.5333 \\ 3 & (9,Inf) & 0.2000 \end{array}$

§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 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

$$H(x) = \sum_{i=1}^{n} \frac{1}{n} \frac{\mu_H(I_i \cap (-\infty, x]))}{\mu_H(I_i)} \qquad (s_j = \frac{1}{n} \sum_{i=1}^{n} \frac{\delta_{ij} s_j}{\sum_{k=1}^{m} \delta_{ik} s_k})$$
$$= \sum_{i=1}^{n} \frac{1}{n} \mathbf{1}_{(l_i \le x < r_i)} \frac{H(x) - H(l_i)}{H(r_i) - H(l_i)} + \sum_{i=1}^{n} \frac{\mathbf{1}_{(r_i \le x, r_i \in \mathcal{S}_H)}}{n}, \ H \in \Theta,$$
(4.4)

where S_H is the support set of H, namely, $x \in S_H$ if $|H(x + \epsilon) - H(x - \epsilon)| > 0 \forall \epsilon > 0$ (or its density > 0). Eq. (4.4) is equivalent to

$$H(x) = \int_{\substack{l \le x < r \\ r}} \frac{H(x) - H(l)}{H(r) - H(l)} d\hat{Q}(l, r) + \int_{r \le x} \mathbf{1}_{(r \in \mathcal{S}_H)} d\hat{Q}(l, r), \quad H \in \Theta, \quad (4.5)$$

where
$$\hat{Q}(l,r) = \sum_{i=1}^{n} \frac{\mathbf{1}_{(l_i \le l, r_i \le r)}}{n}$$
 (the edf of Q) and $Q(l,r) = P(L \le l, R \le r).$

Definition A solution H (or s) to (4.5) (or (4.2)) is called a self-consistent estimator (SCE) of F_X (or s).

Remark GMLE \Rightarrow SCE, but SCE \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].

 $\left(\begin{array}{c} -\infty & (1 \quad 3] \\ \end{array} \right) \left(\begin{array}{c} 5 \end{bmatrix} \left(\begin{array}{c} 7 \\ \end{array} \right) \right)$

Solving equation (4.2) (homework) 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 **s** yield two SCEs of F_X :

$$\hat{F}(t) = \begin{cases} 0 & \text{if } t < 3\\ 1/3 & \text{if } 3 \le t < 5\\ 2/3 & \text{if } 5 \le t < 7\\ 1 & \text{if } t \ge 7; \end{cases} \text{ and } \tilde{F}(t) = \begin{cases} 0 & \text{if } t < 3\\ 1/2 & \text{if } 3 \le t < 7\\ 1 & \text{if } t \ge 7. \end{cases}$$

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 (in hw) that $L(\hat{F}) > L(\tilde{F})$.

Thus \hat{F} is a GMLE but \tilde{F} is not a GMLE.

 \hat{F} and \tilde{F} both satisfy the summation equation in (4.5).

Thus \hat{F} and \tilde{F} are SCEs.

Remark Gu and Zhang (1993) defined an SCE H of F_X to be a solution to the equation in (4.5) without the restriction $H \in \Theta$. They thought that the solution of (4.5) 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 (4.5). **Example 2** Six DC data.

 $N_{1} = \sum_{i=1}^{n} \mathbf{1}_{(L_{i}=R_{i}=2)} = 1 \ ([2,2]), \qquad N_{2} = \sum_{i=1}^{n} \mathbf{1}_{(L_{i}=R_{i}=5)} = 3 \ ([5,5]), \\ N_{3} = \sum_{i=1}^{n} \mathbf{1}_{((L_{i},R_{i}]=(-\infty,4])} = 1 \text{ and } N_{4} = \sum_{i=1}^{n} \mathbf{1}_{((L_{i},R_{i})=(2,+\infty))} = 1.$ A solution *H* to the SC equation in (4.5) 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 (4.5) at $-\infty$, 2, 4 and 5. Verify as follows. $H(x) = \int_{l \le x < r} \frac{H(x) - H(l)}{H(r) - H(l)} d\hat{Q}(l, r) + \int_{r \le x} \mathbf{1}_{(r \in S_H)} d\hat{Q}(l, r).$ (L-Stieltjes)

$$\begin{split} H(-\infty) &= 0, \quad \text{RHD of } (4.5) = \frac{N_1}{n}0 + \frac{N_2}{n}0 + \frac{N_3}{n}0 + \frac{N_4}{n}0 = 0 \text{ trivially.} \\ H(2) &= 1/2, \quad \text{RHD of } (4.5) = \frac{1}{6} + \frac{3}{6} \cdot 0 + \frac{1}{6}\frac{1/2 - 0}{1/4 - 0} + \frac{1}{6} \cdot 0 = 1/2. \\ H(4) &= 1/4, \quad \text{RHD of } (4.5) = \frac{1}{6} + \frac{3}{6} \cdot 0 + \frac{1}{6} \cdot 1 + \frac{1}{6}\frac{1/4 - 1/2}{1 - 1/2} = 1/4. \\ H(5) &= 1, \quad \text{RHD of } (4.5) = \frac{1}{6} + \frac{3}{6} + \frac{1}{6} + \frac{1}{6}\frac{1 - 1/2}{1 - 1/2} = 1. \end{split}$$

Remark. The above examples illustrate that

- A solution to the SC equation is a critical point under the constraint that ∑_{i=1}^m s_i = 1.
 Moreover, the solution to the SC algorithm with initial point s_i ≥ 0 and ∑_{i=1}^m s_i = 1 is a critical point under the constraint that s_i ≥ 0 and ∑_{i=1}^m s_i = 1.
 Finally, the solution to the SC algorithm with initial point s_i ≥ 0 and ∑_{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_X 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
$$(L, R) = (-\infty, Z)\mathbf{1}_{(X \le Z)} + (X, X)\mathbf{1}_{(Z \le X \le Y)} + (Y, +\infty)\mathbf{1}_{(X > Y)}$$
.

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 $F_X(\tau_l) > 0$ (or $F_X(\tau_r) < 1$), $F_X(x)$ is not identifiable for $x < \tau_l$ (or $x > \tau_r$). Thus in such cases, we only consider the estimation of $F_X(x)$ for $x \in [\tau_l, \tau_r]$. Denote $P_c(x) = P\{X \text{ is not censored} | X = x\}$ and

 $K(x) = P\{Z < x \le Y\}.$

$$K(x) = P_c(x)$$
 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)
$$P\{L = \tau_r\} > 0 \text{ if } F_X(\tau_r) < 1; \text{ and } P\{R = \tau_l\} > 0 \text{ if } F_X(\tau_l) > 0$$

(AS3)
$$\int \frac{d[F_Y(u) + F_Z(u)]}{E_Y(u) - E_Y(u)} < \infty;$$

 $\int_{\tau_l < u < \tau_r} F_Z(u) - F_Y(u)$

(AS4) $K(x) = P_c(x) > 0 \text{ for all } x \in (\tau_l, \tau_r);$

(AS5) Z and Y take on finitely many values, say $b_1 < \cdots < b_N$,

and
$$0 < F_X(b_1) < \cdots < F_X(b_N) < 1$$
.

(AS6) There are at most m IIs for each sample size nand $\mu_{F_X}(A_j) > 0$ for each II A_j .

Theorem 1 (Yu and Li (2001)) Suppose that (AS1) and (AS2) hold, the SCE \hat{F} satisfies

$$\lim_{n \to \infty} \sup_{x \in [\tau_l, \tau_r]} |\hat{F}(x) - F_X(x)| = 0 \text{ almost surely,}$$
$$\lim_{n \to \infty} \sup_x |\hat{F}(x) - F_X(x)| = 0 \text{ almost surely if } F_X(\tau_l) = 0 \text{ and } F_X(\tau_r) = 1.$$

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 $\sqrt{n}(\hat{F} - F_X)$ 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 $b_1 < \cdots < b_{k+1}$ be the distinct right endpoints of all the innermost intervals induced by the observed intervals. Then the covariance matrix of the SCE $(\hat{F}(b_1), ..., \hat{F}(b_k))$ 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 $\left(\begin{array}{ccc} c_1 & d_1 & 0 & 0 & \cdots & 0 & 0 \\ c_1 & c_1 & c_2 & c_3 & c_4 \end{array} \right)$

 $form \ J^{-1} = \begin{pmatrix} c_1 & d_1 & 0 & 0 & \cdots & 0 & 0 \\ d_1 & c_2 & d_2 & 0 & \cdots & 0 & 0 \\ 0 & d_2 & c_3 & d_3 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & d_{k-1} & c_k \end{pmatrix}.$

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 $b_1 < \cdots < b_{k+1}$ 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 $(E(J(F_X)))^{-1}$ 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 probabilities g_1, g_2 and g_3 , respectively, where $g_1 + g_2 + g_3 = 1$, and let $F_X(x) = p_1 \mathbf{1}(x \ge 1) + p_2 \mathbf{1}(x \ge 5)$, where $p_1 + p_2 = 1$. Derive all possible SCEs of F_X based on Eq. (4.2) in §4.4. Derive the GMLE of F_X . Are they consistent and asymptotically normal ? (if $g_1 = g_2 = g_3$ and $p_1 = p_2$). **Solution** We first solve for SCE of **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], $(4, +\infty)$, $(-\infty, 4]$, $(2, +\infty)$. -- 1-2----4--5------

Possible IIs: [1,1], (2,4], [5,5] (when n is large enough), with weights s_1, s_2, s_3 . $s_j = \sum_{i=1}^n \frac{1}{n} \frac{\delta_{ij} s_j}{\sum_{k=1}^m \delta_{ik} s_k}, \quad j = 1, ..., m \ (=3).$ Simple case: 1 of each 6 I_i 's.

$$s_1 = \frac{1}{6} \left[\frac{s_1}{s_1} + \frac{s_1}{s_1} + \frac{0}{s_3} + \frac{0}{s_3} + \frac{s_1}{s_1 + s_2} + \frac{0}{s_2 + s_3} \right]$$

$$s_2 = \frac{1}{6} \left[\frac{0}{s_1} + \frac{0}{s_1} + \frac{0}{s_3} + \frac{0}{s_3} + \frac{s_2}{s_1 + s_2} + \frac{s_2}{s_2 + s_3} \right]$$

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

$$s_{1} = \frac{N_{1}}{n} \frac{s_{1}}{s_{1}} + \frac{N_{2}}{n} \cdot 0 + \frac{N_{3}}{n} \cdot 0 + \frac{N_{4}}{n} \frac{s_{1}}{s_{1} + s_{2}}$$
 Can $s_{1} = 0$?

$$s_{2} = \frac{N_{1}}{n} \cdot 0 + \frac{N_{2}}{n} \cdot 0 + \frac{N_{3}}{n} \cdot \frac{s_{2}}{s_{2} + s_{3}} + \frac{N_{4}}{n} \frac{s_{2}}{s_{1} + s_{2}}$$
 Can $s_{2} = 0$?

(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 **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 **s** as $s_i \ge 0$ and $s_1 + s_2 \le 1$. The second one is an SCE if and

only if $U_n \ge 0$. Let 1. $H_1 = \frac{N_1 + N_4}{n} \mathbf{1}_{(x \ge 1)} + \frac{N_2 + N_3}{n} \mathbf{1}_{(x \ge 5)}$. 2. $H_2 = \frac{N_1}{N_1 + N_3} \mathbf{1}_{(x \ge 1)} + U_n \mathbf{1}_{(x \ge 4)} + \frac{N_2}{N_2 + N_4} \mathbf{1}_{(x \ge 5)}$. Thus there are possible two SCEs of F_X : H_1 or $H_2 \text{ if } U_n \ge 0.$ The GMLE is $\hat{F} = \begin{cases} H_1 & \text{if } U_n < 0 \\ H_2 & \text{if } U_n \ge 0. \end{cases}$

The reason is as follows.

Theorem 4 Given interval-censored data (C1, C2, DC or MIC data), let A_i , i = 1, ..., mbe all the distinct IIs. Suppose \hat{s}_i , i = 1, ..., 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 $\tilde{s}_i, i = 1, ..., m$, then $\hat{\mathbf{s}}$ is a GMLE.

The proof of the theorem utilizes the convexity of $-\ln L$ in $(s_1, ..., s_k)$. Note that GMLE is also an SCE.

If $U_n < 0$, there is only one SCE of s, *i.e.* H_1 . Thus H_1 is a GMLE if $U_n < 0$. If $U_n \ge 0$, there are two SCEs: H_1 and H_2 .

Then one can use two ways to show H_2 is the GMLE:

(1) Theorem 4.

(2) $L(H_2) > L(H_1)$

To prove the consistency of the GMLE, check (AS1) and (AS2) by Th. 1:

(AS1) $P\{X \in (\tau_l, \tau_r] \text{ and } K(X) = 0\} = 0$, where $K(x) = P\{Z < x \le Y\}$.

(AS2) (1) $P\{L = \tau_r\} > 0$ if $F_X(\tau_r -) < 1$; and (2) $P\{R = \tau_l\} > 0$ if $F_X(\tau_l) > 0$. (Z,Y) $\in \{(0.5,2), (2,4), (4,8)\}$ and X = 1 or 5.

Now verify them (when $g_i = 1/3$ and $p_i = 1/2$) as follows:

AS1: $K(x) = 1/3 \ \forall \ x \in (0.5, 2] \cup (2, 4] \cup (4, 8] \supset [1, 5] \ (= [\tau_l, \tau_r]).$

Thus AS1 holds *i.e.* $P(X \in [\tau_l, \tau_r] \text{ and } K(X) = 0) = 0.$

AS2: (1) $P\{L = \tau_r\} = P\{X = 5, (Z, Y) = (4, 8)\} = 0.5/3 > 0$ and $F(\tau_r -) = 1/2 < 1$. (2) $P\{R = \tau_l\} = P\{X = 1, (Z, Y) = (0.5, 2)\} = 0.5/3 > 0$ and $F(\tau_l) = 1/2 > 0$. Thus AS2 holds.

Thus the SCEs are consistent.

To discuss the asymptotic normality (see Theorem 2), verify that

(1) AS5 and AS6 fail, as $F_X(2) = p_1 \not < F_X(4) = p_1$ and

(2) AS4 fails, as $K(3) = P\{Z < 3 \le Y\} = g_2 > 0$, but $P_c(3)$ is not defined. Thus Theorem 2 is not valid. Does it mean that the GMLE is not asymptotically normal? Th. 3 is not valid. Is J^{-1} given in Th.3 the asymptotic covariance matrix of \hat{F} ? We now establish the asymptotic normality directly. Note that

$$E(N_1/n) = P\{X = 1, (Z, Y) = (0.5, 2) \text{ or } (2, 4)\} = p_1(g_1 + g_2)$$

$$E(N_2/n) = P\{X = 5, (Z, Y) = (4, 8) \text{ or } (2, 4)\} = p_2(g_2 + g_3)$$

$$E(N_3/n) = P\{X = 5, (Z, Y) = (0.5, 2)\} = p_2g_1$$

$$E(N_4/n) = P\{X = 1, (Z, Y) = (4, 8)\} = p_1g_3.$$

(5.1)

$$U_n = \frac{N_4}{N_2 + N_4} - \frac{N_1}{N_1 + N_4} \xrightarrow{a.s.} \frac{p_1 g_3}{p_1 g_3 + p_2 (g_2 + g_3)} - \frac{p_1 (g_1 + g_2)}{p_1 (g_1 + g_2) + p_2 g_1} < 0$$

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 .

$$\vdash: \qquad \frac{H_1(t) - F_X(t)}{\sigma_{H_1(t)}} \xrightarrow{D} N(0, 1), \quad t \in [1, 5) \qquad \text{(by the CLT)},$$

 $\begin{aligned} \mathbf{Reason:} \ H_1(t) &= \begin{cases} (N_1 + N_4)/n & \text{if } t \in [1,5) \\ 1 & \text{if } t \ge 5 \end{cases} \\ N_1 &= \sum_{i=1}^n \mathbf{1}((L_i, R_i) = (1,1) \text{ or } (-\infty, 2)) \text{ and} \qquad N_4 = \sum_{i=1}^n \mathbf{1}((L_i, R_i) = (-\infty, 4)) \end{aligned} \\ \text{If } t \in [1,5) \text{ then } H_1(t) &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}((L_i, R_i) \in \{(1,1), (-\infty, 2), (-\infty, 4)\}) \\ &= \overline{\mathbf{1}(X = 1, (Z, Y) \in \{(0.5, 2), (2, 4), (4, 8)\})} \\ &= \overline{\mathbf{1}(X = 1)} \\ \sigma_{H_1(t)}^2 &= \begin{cases} Var(\sum_{i=1}^n \mathbf{1}_{(X_i=1)}/n) = p_1 p_2/n & \text{if } t \in [1,5) \\ 2 & \text{if } t \ge 5 \end{cases}. \end{aligned} \end{aligned}$ (5.2)

By (5.1), $\mathcal{L} = N_1 \ln F(1) + N_2 \ln(1 - F(4)) + N_3 \ln(1 - F(1)) + N_4 \ln F(4)$, where F(5) = 1. Verify the Fisher information matrix is

$$J/n = -\frac{1}{n} \frac{\partial^2 \mathcal{L}}{\partial (F(1), F(4)) \partial (F(1), F(4))^t} = \begin{pmatrix} \frac{N_1}{nF^2(1)} + \frac{N_3}{n(1-F(1))^2} & 0\\ 0 & \frac{N_4}{nF^2(4)} + \frac{N_2}{n(1-F(4))^2} \end{pmatrix}$$

$$\xrightarrow{a.s.} \begin{pmatrix} \frac{g_1+g_2}{p_1} + \frac{g_1}{p_2} & 0\\ 0 & \frac{g_3}{p_1} + \frac{g_2+g_3}{p_2} \end{pmatrix} = \begin{pmatrix} 2 & 0\\ 0 & 2 \end{pmatrix}$$

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 $F_X(x) = p_1 \mathbf{1}_{(x \ge 1)} + p_2 \mathbf{1}_{(x \ge 3)} + p_3 \mathbf{1}_{(x \ge 5)}$, 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_X based on Eq. (4.2) in §4.4.
 - 1.b. Find the limits of the SCEs directly.
 - 1.c. Derive the GMLE of F_X .
 - 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 $t \in \mathbb{R}^1$.
 - 2.b. Are they consistent estimators of F_X ?
 - 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 $-H = -E\left(\frac{\partial^2 \mathcal{L}}{\partial F(b_i)\partial F(b_j)}\Big|_{F=F_X}\right)_{2\times 2}$, where $(b_1, b_2) = (1, 4)$ and $\mathcal{L} = \ln \mathbb{L}$.

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, $H = F_X$ satisfies the self-consistent equation

$$H(x) = \int_{l \le x < r} \frac{H(x) - H(l)}{H(r) - H(l)} dQ(l, r) + P\{R \le x\}, \ H \in \Theta.$$

4. Under double censoring with Y = -X and Z = Y - 1, derive the function $E(\mathbf{1}(X \le t, L \le t < R | (L, R) = (l, r))$ of (l, r, t) for an F_X you select, and show that $H = F_X$ 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) (mixture of RC model and mixed IC model). Assume:

- (1) N is a random integer;
- (2) $T, Y_1 < Y_2 < \cdots < Y_k < \cdots$ are inspection times;
- (3) X and $(N, T, Y_i, i \ge 1)$ 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, ..., 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_X . Define

- $\tau = \sup\{x : F_X(x) < 1 \text{ and } F_T(x) < 1\},\\ \tau_Y = \sup\{t : F_{Y_N}(t) < 1\}.$
- $au_N = \sup\{i: f_N(i) > 0\}.$

A point x is called a support point of a cdf F if $|F(x + \epsilon) - F(x - \epsilon)| > 0 \ \forall \ \epsilon > 0$. Denote S_F the set of all support points of F. If S(t) = exp(-t), t > 0, then $S_F = (0, \infty)$ or $[0, \infty)$? People make use of the assumptions as follows:

(AS1) $\tau_Y \leq \tau$;

(AS2) P{T or $Y_N = \tau$ } > 0 if $F_X(\tau -) < 1$.

- (AS3) H_n is right continuous and $\mathcal{S}_{H_n} \subset \{R_1, ..., R_n\};$
- (AS4) $F_X(\tau) > 0$ and $\bigcup_{i \leq \tau_N} \mathcal{S}_{F_{Y_i}} \subset \mathcal{S}_{F_X}$;
- (AS5) $\cup_{i \leq \tau_N} \mathcal{S}_{F_{Y_i}}$ is a finite set and F_X 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 $\mu_{F_X}(A_j) > 0$ 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

$$\lim_{n \to \infty} \sup_{x} |H_n(x) - F_X(x)| = 0 \text{ a.s. if } F_X(\tau) = 1 \text{ and } \lim_{n \to \infty} \sup_{x \le \tau} |H_n(x) - F_X(x)| = 0 \text{ a.s.}.$$

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

$$\sqrt{n}(H_n(x) - F_X(x)) \xrightarrow{D} Z(x) \text{ for } x \leq \tau,$$

where Z is a Gaussian process on $[0, \tau]$. Let $b_1 < \cdots < b_{k+1}$ be the right endpoints of all the distinct IIs induced by the observed intervals. The covariance matrix of $(H_n(b_1), ..., H_n(b_k))$ can be estimated by $-\left(\frac{\partial^2 \mathcal{L}}{\partial F(b_i)\partial F(b_j)}\right)_{k \times k}^{-1}|_{F=\hat{F}}$ and $\mathcal{L} = \ln L$.

Remark. The assumptions AS1, AS4, AS5 and AS6 eliminate the case in which $\mu_{F_X}(A_j) = 0$ 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 $\mu_{F_X}((4, 5]) = 0$.

Remark. For IC data, $-\ln \mathbf{L}$ is strictly convex in the interior of $\mathbf{s} \in D_s$, where

 $D_s = \{ \mathbf{s} : \vec{s} = (s_2, ..., s_m), s_i \ge 0, \sum_{i=2}^m s_i \le 1 \}, s_1 = 1 - s_2 - \dots - s_m.$

Thus there is a unique GMLE of **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 $x \in \mathbb{R}^1$, as $\frac{d^2(-\ln x)}{dx^2} = \frac{1}{x^2} > 0$, $-\ln \mathbf{L}$ is strictly convex. $\delta \mathbf{4} \in \mathbf{2}$ Homoverk

§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_X . (**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.)

$\S4.7.$ Case 2 data and the GMLE

The mixed case IC model is used to formulate case 2 data (without exact observations). Let $\mathbf{Y} = \{Y_j : j \ge 1\}$ be a sequence of random variables such that $Y_1 < Y_2 < \cdots : Y_0 = -\infty$. K be a positive random integer such that $(K, \mathbf{Y}) \perp X$;

On the event $\{K = k\}$, let $(L, R) = \begin{cases} (Y_{i-1}, Y_i) & \text{if } Y_{i-1} < X \le Y_i, i \in \{1, 2, ..., k\} \\ (Y_k, \infty) & \text{if } X > Y_k \end{cases}$,

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.

 $\int_{-\infty}^{x} d\nu(t) = \int_{-\infty}^{x} (dP(L \le t) + dP(R \le t)). \qquad \qquad \int_{-\infty}^{\infty} d\nu(t) = ?$ **Theorem 1.** (Schick and Yu (2000)) Suppose $E(K) < \infty$. Let $A = \bigcup_{i \le \tau_K} \mathcal{S}_{F_{Y_i}}$, where $\tau_K = \sup_{i \le \tau_K} \{i: P\{K = i\} > 0\}$, and \hat{F} be a GMLE of F_X . Then

1. $\int |\hat{F} - F_X| d\nu \to 0$ almost surely;

2. If $P(Y_i = a) > 0$ for some $i \leq \tau_K$, then $\hat{F}(a) \to F_X(a)$ a.s.;

3. If A is dense in $[0,\infty)$ and F_X is continuous, then $\sup_x |\hat{F}(x) - F_X(x)| \to 0$ a.s..

Remark The assumption that $E(K) < \infty$ can be removed.

Question: Relation between the 3 statements in the theorem ?

Question: Why not SCE ? Recall Th 1 in §4.5 says that

the SCE is consistent under DC model with AS1 and AS2.

Example 1 (inconsistent SCE with C2 data). Let $K \equiv 2$ (Y = (Y₁, Y₂)),

 $P(X = 2) = P(X = 4) = P(X = 6) = 1/3, P(\mathbf{Y} = (1, 5)) = P(\mathbf{Y} = (3, 7)) = 1/2.$ If n is large enough, the observations are

 N_1 (- ∞ , 3], N_2 (5, ∞), N_3 (3,7] and N_4 (1,5].

IIs: (1,3], (3,5] 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 for n large enough.

$$H_1(t) = \frac{N_1 + N_4}{n} \mathbf{1}_{(t \ge 3)} + \frac{N_2 + N_3}{n} \mathbf{1}_{(t \ge 7)},$$
$$\hat{F}(t) = \frac{N_1}{N_1 + N_3} \mathbf{1}_{(t \ge 3)} + U_n \mathbf{1}_{(t \ge 5)} + \frac{N_2}{N_2 + N_4} \mathbf{1}_{(t \ge 7)}.$$

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 that Y_i 's are discrete, thus we expect a GMLE is consistent at 1, 3, 5, 7. The limit of $H_1(3)$ is

$$\lim_{n \to \infty} \frac{N_1 + N_4}{n} = P\{X = 2, \mathbf{Y} = (3, 7)\} + P\{X \le 4, \mathbf{Y} = (1, 5)\} = 1/2 \neq 1/3 = F_X(2).$$

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 II's with right endpoints b_i for each sample size, F_X is strictly increasing on b_i 's, $\begin{pmatrix} \hat{F}(b_1) - F_X(b_1) \end{pmatrix}$

then
$$\sqrt{n} \begin{pmatrix} \Gamma(0_1) & \Gamma_X(0_1) \\ & \ddots & \\ \hat{F}(b_k) - F_X(b_k) \end{pmatrix} \xrightarrow{\mathcal{D}} N(0_k, \Sigma_k) \text{ as } n \to \infty, \text{ where}$$

$$\Sigma_k = -n \left(E(\frac{\partial^2 \mathcal{L}}{\partial F(b_i) \partial F(b_j)} \Big|_{F=F_X}) \right)_{k \times k}^{-1}$$

and the asymptotic covariance matrix of \hat{F} can be estimated by $-\left(\frac{\partial^2 \mathcal{L}}{\partial F(b_i)\partial F(b_j)}\right)_{k\times k}^{-1}|_{F=\hat{F}}$. **Theorem 3.** (Groeneboom (1996)). Let F_X 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 $0 < \epsilon_o \le 1/(2M)$, so g does not have mass close to the diagonal.

Then we have at each point $t_o \in (0, M)$

$$n^{1/3} \{ 2a(t_o) / f_o(t_o) \}^{1/3} \{ \hat{F}(t_o) - F_X(t_o) \} \xrightarrow{\mathcal{D}} 2Z^*,$$

where Z^* is defined as in Theorem 2, and

$$a(t_o) = \frac{g_1(t_o)}{F_X(t_o)} + k_1(t_o) + k_2(t_o) + \frac{g_2(t_o)}{1 - F_X(t_o)},$$

$$k_1(u) = \int_u^M \frac{g(u, v)}{F_X(v) - F_X(u)} dv \text{ and } k_2(v) = \int_0^v \frac{g(u, v)}{F_X(v) - F_X(u)} du.$$

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_X and G have continuous derivatives,

with their densities $f_o(x_o) > 0$ and $g(x_o, x_o) > 0$, and

let F be the estimator of F_X , obtained at the first step of the iterative convex minorant algorithm, starting the iterations with F_X . Then the statistic

 $(n \ln n)^{1/3} \frac{\tilde{F}(x_0) - F_X(x_0)}{\{\frac{3}{4}(f_o(x_0))^2/g(x_0, x_0)\}^{1/3}} \xrightarrow{\mathcal{D}} 2Z^*,$ 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 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_X is unknown and thus it is impossible to start the iterations from F_X except in simulation.

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, ..., K, takes on finitely many values. Then the convergence rate is $n^{1/2}$.
- 2. 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 $n^{1/3}$.
- 3. 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 $(n \ln n)^{1/3}$.

§4.7.2. Homework

- 1. a. Derive the limits of the SCEs in Example 1 and compare to F_X .
 - 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 $\{1, 2, ..., 9\}$ and 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 intervalcensored 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.

$\S4.8.$ Application to model justification.

It is desirable to use parametric models instead of nonparametric models. Given a RC data set of n = 100, one can try various MLEs (see Figure 4.8.1). They can be quite different. How can we decide which to choose ?



Fig. 4.8.1. MLEs of 4 parametric survival functions based on 100 observations That is, 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 themselves.

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 $S_X(t)$ and its 95% CI along $t < \tau$ based on the IC data.

2. Plot the MLE of $S_X(t)$ for the potential parametric model based on the IC data.

3. If the MLE curve lies within the confidence band, the parametric model may fit the data. $\begin{cases} edf & \text{if the data is complete} \\ BLE & \text{if the data is complete} \end{cases}$

Here, GMLE is given by
$$\begin{cases} PLE & \text{if data are RC} \\ \text{the max-min form} & \text{if data are C1}, \\ SCE & \text{otherwise.} \end{cases}$$
 Just survfit()?

In the next figure, we add the GMLE curve together with the pointwise CI of the PLE. It suggests to eliminate two distributions which two ?

One may further draw the pointwise CI curves of the other two and check which one of them contains the GMLE better than the other.

The survival function of the Weibull distribution is

$$S(t;\theta,\kappa) = e^{-(t/e^{\hat{\theta}})^{\kappa}}, t > 0, \qquad \qquad \theta = \text{intercept}, \tau = e^{\theta},$$

A CI is $e^{-(t/e^{\hat{\theta} \pm 1.96SE})^{\kappa}}.$

The survival function of the lognormal distribution is

$$S(t:\mu,\sigma) = \int_{\frac{\ln t - \mu}{\sigma}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

A CI is $S(t:\mu,\sigma) = \int_{\frac{\ln t - \hat{\mu} \pm 1.96SD}{\sigma}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$



$$\begin{split} & \text{lines}(x,1\text{-plnorm}(x,zz\$coef+1.96*w,zz\$scale), type="l", lty=5) \\ & (zz=\text{survreg}(\text{Surv}(m,d)\sim 1)) \\ & w= \text{summary}(zz)\$table[1,2] \\ & \text{SD of intercept} \\ & \text{lines}(x,1\text{-pweibull}(x,1/zz\$scal,exp(zz\$coef-1.96*w)), type="l", lty=3) \ \# \text{CI} \\ & \text{lines}(x,1\text{-pweibull}(x,1/zz\$scal,exp(zz\$coef+1.96*w)), type="l", lty=3) \\ & \text{leg.names} < -c("ple", "weib", "lognormal") \\ & \text{legend}(300, \ 0.88, \ \text{leg.names}, \ \text{lty}=c(1,3,5), \text{cex}=1.0) \\ & \text{The CI plots suggest a proper parametric distribution.} \end{split}$$



Fig. 4.8.3. MLEs of 2 parametric survival functions based on 100 observations We shall introduce several other methods here.

§4.8.1. Q-Q plot

For a complete data set

 $X_1, ..., X_n$ from a cdf F_X .

If we suspect F_X belongs to a parametric distribution family, say

 $H_0: F_X = F(\cdot, \theta),$

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 $X_{(1)} \leq \cdots \leq X_{(n)}$ be order statistics;

They are $100\frac{1}{n+1}$, ..., $100\frac{n}{n+1}$ -th sample percentiles (quantiles) of the sample. Let $y_i = \sup\{u: F(u, \hat{\theta})\} < \frac{i}{n+1}\}, i = 1, ..., n.$ — estimated population percentiles (quantiles). Plot $(y_i, X_{(i)}), i = 1, ..., n.$

89

If the plot is close to the line y = x, then $F \approx F_X$. qqplot() is a function in R. It only needs the sample and F_X . Recall the the 100*p*-th quantile of a cdf F is a

$$q = F^{-1}(p)$$
, where $F^{-1}(p) = \sup\{q : F(q) < p\}$.

We first explain the reasoning of QQ-plot with complete data.

Complete data. Let $a_1 < \cdots < a_m$ be distinct points among X_1, \ldots, X_n . Since the edf $\hat{F} \to F_X$ a.s. on $(0, \infty)$, we expect the quantile functions $(\hat{F})^{-1}$ and F_X^{-1} are close. These $(a_i, F_X^{-1}(\hat{F}(a_i))), i = 1, \ldots, m$ 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 $X \sim F_X$ and (L_i, R_i) 's are iid IC observations. If we suspect that F_X is some of given form, say $F_X = F_*(\cdot, \theta)$, we can also do QQ-plot as follows.

- 1. Obtain the IIs based on observations (L_i, R_i) , i = 1, ..., n, denoted by A_j , j = 1, ..., m.
- 2. Obtain the GMLE of the distribution function F_X , 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 $F_X(t) = F_*(\frac{t-\mu}{\sigma}, \theta)$.

$\S4.8.2.$ Hazard plot for RC data

 $(Z_i, \delta_i), i = 1, ..., n.$

Plot sample integrated hazard v.s. integrated hazards of $F_X(\cdot, \theta)$. To get sample integrated hazards:

a. Order Z_i 's as $Z_{[n]} \leq \cdots \leq Z_{[1]}$ (reverse order);

- b. Denote $\delta_{[i]}$, i = 1, ..., 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 ${\cal Z}_{[k]}$ is

$$\hat{H}(Z_{[k]}) = -\sum_{i: \ Z_{[i]} \le Z_{[k]}, \delta_{[i]} = 1} \log[1 - \hat{h}(Z_{[i]})].$$

 $(\hat{S}(t) = \prod_{i:Z_{(i)} < t} (1 - \frac{\delta_{(i)}}{n+1-i}))$

f. Plot $(\hat{H}(Z_{[k]}), H_o(Z_{[k]}))$, where $\delta_{[k]} = 1$ and $H_o(x) = -\log S_X(x, \hat{\theta})$). **Example**. Hazard plot. Use Sample 0 in Table 1.1 v.s. exponential distribution.

Table 1.2. Calculation of sample hazards						
Patient	Remission	Reverse	h(x)	Integrated		
Number	Time	Order (K)	= f(x)/S(x-)	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.

Note that the H(6) appears three times as there are three exact observations at x = 6. However, they are obviously the same. Thus we plot

 $(0.15, 6\rho), (0.21, 7\rho), (0.28, 10\rho), (0.37, 13\rho), (0.47, 16\rho), (0.62, 22\rho), (0.80, 23\rho),$ 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 $a_1 < \cdots < a_m$ be distinct point among all exact observations. Since $F_{pl} \to F_*$ a.s., we expect sample cumulative hazard $\hat{H} \to H_o$ 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_X completely.

3 ways for model checking:

- (1) plot the cdf's together with the CI band of one cdf (compare two cdf's curves),
- (2) QQ-plot (check for linearity),
- (3) Hazard plot (check for linearity).

$\S4.8.3$. Diagnostic plot under regression set-up with IC data

In parametric regression analysis, in addition to survival time X, we also observe covariate Z and the survival function is a function of $\beta' \mathbf{Z}$, such as $X = \beta' \mathbf{Z} + W$ in linear regression or $S(t|\mathbf{z}) = (S_o(t))^{\exp(\beta' \mathbf{Z})}$ in Lehmann Model (or Cox model), among others. Typically, under the linear (or log-linear) regression there are 3 ways for model checking:

(1) plot the MLE of the parametric cdf of W together with its CI band and the GMLE,

(2) QQ-plot (check for linearity),

(3) The marginal distribution (MD) plot.

The MD plot: Instead plot the MLE \tilde{F}_W and the GMLE \hat{F}_W , the MD approach plots the GMLE $\hat{S}_X(t)$ based on the IC data (L_i, R_i) 's, and the MLE of S_X under the parametric assumption, denoted by \tilde{S}_X , which is obtained as follows.

1. Let $W_1, ..., W_m$ be i.i.d. from the baseline distribution, where m is large, say 100.

2. Let $Z_1^*, ..., Z_m^*$ be i.i.d. from the empirical distribution of Z_i 's.

3. Let $X_i^* = W_i + \hat{\beta}' \mathbf{Z}_i^*, i = 1, ..., m$.

4. Finally, plot the edf based on X_i^* 's against the GMLE of $S_X(t)$.

For instance, under the log linear regression model, $\ln X = \beta' \mathbf{Z} + W$, where the baseline survival function S_W is either known (or known up to a parameter), the vectors

 $(\ln L_i - \beta' \mathbf{z}_i, \ln R_i - \beta' \mathbf{z}_i), i = 1, ..., n, \text{ are i.i.d. from } S_W,$

where either $(\mu_W, \sigma_W) = (0, 1)$, or μ_W or σ_W is a parameter. One can plot S_W against the

GMLE of S(t) based on modified IC data $(L_i - \hat{\beta}' \mathbf{z}_i, R_i - \hat{\beta}' \mathbf{z}_i)$'s, where $\hat{\beta}$ is the parametric MLE of β .

On the other hand, we can also use the following QQ-plot to check whether the IC regression data satisfy a certain parametric distribution F_* . We now explain the QQ-plot 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 $(L_i, R_i, \mathbf{z}_i), i = 1, ..., n$.

- 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 $(\ln L_i \hat{\beta}^t \mathbf{z}_i, \ln R_i \hat{\beta}^t \mathbf{z}_i), i = 1, ..., n$, denoted by $A_j, j = 1, ..., m$.
- 3. Obtain the GMLE of the distribution function based on $(\ln L_i \hat{\beta}^t \mathbf{z}_i, \ln R_i \hat{\beta}^t \mathbf{z}_i), i = 1, ..., n$, 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*, where m_i is either the midpoint or end points of the finite A_i 's.

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.

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 $\{1, 2, ..., 28\}$. Conditional on K = k, the follow-up time Y_i , i = 1, ..., 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 Tis 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 based on pivotal function 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 $(T(L_i, \mathbf{z}_i, \hat{\beta}), T(R_i, \mathbf{z}_i, \hat{\beta})), i = 1, ..., n,$ denoted by $A_j, j = 1, ..., m$.
- 4. Obtain the GMLE of the distribution function based on $(T(L_i, \mathbf{z}_i, \hat{\beta}), T(R_i, \mathbf{z}_i, \hat{\beta})), i = 1, ..., n$, 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,

$$S_X(x) = e^{-x^{\kappa} e^{\beta^t \mathbf{Z}}}, \ x > 0.$$
 not a nice form

One of such pivoting functions is

$$T = X^{\kappa} e^{\beta^t \mathbf{Z}} = \left(\frac{X}{e^{\frac{-\beta^t \mathbf{Z}}{\kappa}}}\right)^{\kappa},$$

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 $(T(L_i, \mathbf{z}_i, (\beta, \kappa)))$, $T(R_i, \mathbf{z}_i, (\beta, \kappa))] = (L_i^{\kappa} e^{\beta^t \mathbf{Z}_i}, R_i^{\kappa} e^{\beta^t \mathbf{Z}_i}] = (L_{i,\beta,\kappa}, R_{i,\beta,\kappa}]$. 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, ..., n. Finally we plot \hat{F}_* against F_T , the exponential distribution Exp(1).

A nicer form is $T = (\frac{X}{exp(\beta'\mathbf{Z})})^{\kappa}$ with $S_X(x) = exp(-(\frac{x}{exp(\beta'\mathbf{Z})})^{\kappa}).$

Remark. It is worth mentioning that under the Cox model, the first two approach, *i.e.*, the parametric and GMLE of S_W plots and the QQ-plot may not be feasible. The reason is that there may not exists a pivotal function.

§4.8.3.2. Homework:

In the following, make comments on whether the plots suggest $F_X(x) = F_*(\frac{x-\mu}{\sigma})$.

- 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 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}) .

The semi-parametric analysis deals with regression data. It assumes that $X|\mathbf{Z}$ satisfies a certain model, *i.e.*, Cox's model or the LR model, but the baseline distribution $F_{X|\mathbf{Z}}(\cdot|0)$ is unknown.

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,

 $\mathbf{Z} = \mathbf{1}_{\text{(patient is from sample 1)}}$.

Two typical models are considered, among other models.

1. Proportional hazards (PH) model: Conditional on $\mathbf{Z} = \mathbf{z}$, the hazard function

$$h(t|\mathbf{z}) = \psi(\beta, \mathbf{z})h_o(t), \ t < \tau, \tag{1.1}$$

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

$$S_{X|\mathbf{Z}}(t|\mathbf{z}) = (S_o(t))^{\psi(\beta,\mathbf{Z})} \quad (H_{X|\mathbf{Z}}(t|\mathbf{z}) = \psi(\beta,\mathbf{z})H_o(t))$$
(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.

Otherwise, (1.1) and (1.2) are different models. None of them is a special case of the other model.

 $(F(x) - F(a)) = \int_{a}^{x} F'(t)dt \ \forall \ x)$

2. Accelerated lifetime model: Conditional on $\mathbf{Z} = \mathbf{z}$,

$$X = X_o/\psi(\beta, \mathbf{z}) \quad (\ln X = \ln X_o - \ln \psi(\beta, \mathbf{z})). \tag{1.3}$$
$$S_{X|\mathbf{Z}}(t|\mathbf{z}) = S_o(\psi(\beta, \mathbf{z})t) \qquad (S_o = S_{X_o}).$$

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
$$\psi$$
:
$$\begin{cases} e^{\beta^{t} \mathbf{Z}} - \log \text{ linear}; \\ log(1 + e^{\beta^{t} \mathbf{Z}}) - \log \text{ listic.} \end{cases} \quad \psi \ge 0. \tag{1.4}$$

$\S5.2.$ 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 $(M_i, \delta_i, \mathbf{z}_i), i = 1, ..., n$ be i.i.d. copies of $(M, \delta, \mathbf{Z}).$ The log likelihood function is

$$\mathcal{L}(\beta) = \ln \prod_{i=1}^{n} (f(M_i | \mathbf{z}_i))^{\mathbf{1}_{e,i}} (S(M_i | \mathbf{z}_i))^{\mathbf{1}_{r,i}}$$

= $\ln [\prod_{i: ex} h(M_i | \mathbf{z}_i) \prod_{i=1}^{n} S(M_i | \mathbf{z}_i)]$
= $\sum_{i: ex, M_i < \tau} \ln \psi(\beta, \mathbf{z}_i) + \sum_{i: ex, M_i < \tau} \ln h_o(M_i) + \sum_{i=1}^{n} \psi(\beta, \mathbf{z}_i) \ln S_o(M_i).$

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:

 $a_1 < \cdots < a_g$ — all the distinct exact observations. By rearranging the index, assume $X_i = a_i, i = 1, ..., g$. $\mathcal{R}_j = \mathcal{R}(a_j) = \{i : M_i \ge a_j\},\ \phi(i) = \psi(\beta, \mathbf{z}_i), \qquad (\text{see } (1.4))$ $\mathcal{D} = \{i : \delta_i = 1, i = 1, ..., n\}$ (note that all exact observations are distinct), Define a modified likelihood function

$$lik = \prod_{i \in \mathcal{D}} \frac{\phi(i)}{\sum_{k \in \mathcal{R}_i} \phi(k)}.$$
(2.1)

The log likelihood

$$l(\beta) = \ln lik = \sum_{i \in \mathcal{D}} [\ln \phi(i) - \ln \sum_{k \in \mathcal{R}_i} \phi(k)].$$
(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 $\S5.2$, let

$$\psi(\beta, \mathbf{z}) = e^{\beta^t \mathbf{Z}}.$$
(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

 $\begin{array}{l} (1,1,1), (2,0,5), (2.5,1,2), (4,0,1), (7,1,2),.\\ \text{Rearrange (according to exact obs.) } (M_i, \delta_i, \mathbf{z}_i)\text{'s as}\\ (1,1,1), (2.5,1,2), (7,1,2), (2,0,5), (4,0,1),.\\ \text{We use this order from now on.}\\ 3 \text{ exact observations: } (a_1, a_2, a_3) = (1, 2.5, 7).\\ \mathcal{D} = \{1,2,3\},\\ \mathcal{R}_1 = \mathcal{R}(1) = \{1,2,3,4,5\}, \ \mathcal{R}_2 = \mathcal{R}(2.5) = \{2,3,5\}, \ \mathcal{R}_3 = \mathcal{R}(7) = \{3\};\\ \phi(i)\text{'s are } e^{\beta}, \ e^{2\beta}, \ e^{2\beta}, \ e^{\beta}, \ e^{\beta}, \end{array}$

$$lik = \frac{e^{\beta}}{2e^{2\beta} + e^{5\beta} + 2e^{\beta}} \cdot \frac{e^{2\beta}}{2e^{2\beta} + e^{\beta}} \cdot \frac{e^{2\beta}}{e^{2\beta}} = \frac{1}{(2 + e^{-\beta})(2e^{\beta} + e^{4\beta} + 2)}$$

§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.

\S **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:

 $a_1 < \cdots < a_g$ are all distinct exact observations;

 d_j is the # of deaths at a_j ;

 S_j is the collection of all the combinations of selecting d_j elements out of those in $\mathcal{R}(a_j)$; $r_j = |\mathcal{R}(a_j)|$;

$$\mathcal{R}_j = \mathcal{R}(a_j) = \{i : M_i \ge a_j\},\\ \mathcal{D}_i = \{i : \delta_i = 1, M_i = a_i\}:$$

$$\phi(i) = \psi(\beta, \mathbf{z}_i) \text{ (often } = e^{\beta \mathbf{Z}_i}).$$

A log likelihood is defined as

$$l(\beta) = \ln \prod_{j=1}^{g} \frac{\prod_{i \in \mathcal{D}_{j}} \phi(i)}{\sum_{(i_{1},\dots,i_{d_{j}}) \in \mathcal{S}_{j}} \phi(i_{1}) \cdots \phi(i_{d_{j}})}.$$
 $(lik = \prod_{i \in \mathcal{D}} \frac{\phi(i)}{\sum_{k \in \mathcal{R}_{i}} \phi(k)})$

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 4/1 is log linear and there

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 $\mathcal{R}_{i_j} = \mathcal{R}(a_j), j = 1, ..., g,$ where $g = \sum_{i=1}^n \delta_i.$

Cox made the following assumptions:

1. there is only one death at a_i and thus we can pretend that

there is only one death in $[a_i, a_i + \epsilon)$ for a small ϵ ;

2. death only occurs in $[a_j, a_j + \epsilon), j = 1, ..., g$. Then

$$P\{\text{the } i\text{-th person died in } [a_j, a_j + \epsilon) | \text{the } i\text{-th person survived time } a_j)$$
(0)
= $P(X_i \in [a_j, a_j + \epsilon) | X_i \ge a_j)$
= $P(X_i \in [a_j, a_j + \epsilon)) / P(X_i \ge a_j)$
 $\approx h(a_i | \mathbf{z}_i) \epsilon,$

as $h(t)\epsilon = \frac{f(t)\epsilon}{S(t-)}$ and P{the k-th death in $[a_j, a_j + \epsilon)} \approx f(a_j | \mathbf{z}_i)\epsilon$ },

$$P\{ \text{ a death in } [a_j, a_j + \epsilon) | \text{ each in } \mathcal{R}_{i_j} \text{ survives time } a_j \}$$
(1)

$$\left(= P\{\bigcup_{k \in \mathcal{R}_{i_j}} A_k | \cap_{h \in \mathcal{R}_{i_j}} C_h\} \quad (A_k = \{X_k \in [a_j, a_j + \epsilon)\} \text{ and } C_h = \{X_h \ge a_j\})$$

$$(\text{as any one in the risk set may die}) \right)$$

$$= \sum_{k \in \mathcal{R}_{i_j}} P\{ \text{ the } k\text{-th died in } [a_j, a_j + \epsilon)| \text{ each in } \mathcal{R}_{i_j} \text{ survives time } a_j \}$$

$$= \sum_{k \in \mathcal{R}_{i_j}} P\{A_k | \cap_{h \in \mathcal{R}_{i_j}} C_h\}$$

$$= \sum_{k \in \mathcal{R}_{i_j}} P\{A_k | C_k \cap \cap_{h \in \mathcal{R}_{i_j}, h \neq k} C_h\}$$

$$= \sum_{k \in \mathcal{R}_{i_j}} P\{A_k | C_k\}$$

$$(by independence of observations)$$

$$= \sum_{k \in \mathcal{R}_{i_j}} P\{ \text{ the } k\text{-th died in } [a_j, a_j + \epsilon)| \text{ the } k\text{-th survived time } a_j\}$$

$$\approx \sum_{k \in \mathcal{R}_{i_j}} h(a_j | \mathbf{z}_k) \epsilon$$

$$(by (0))$$

$$i.e. \sum_{k \in \mathcal{R}_{i_j}} P\{A_k | \cap_{h \in \mathcal{R}_{i_j}} C_h\} \approx \sum_{k \in \mathcal{R}_{i_j}} h(a_j | \mathbf{z}_k) \epsilon$$

$$(2)$$

Let $B_i = \{ \text{one death in } [a_j, a_j + \epsilon) \text{ from } \mathcal{R}_{i_j}, \text{ and each in } \mathcal{R}_{i_j} \text{ survived } a_j \},\$

$$(as \ P(A|BC) = \frac{1}{P(BC)} = \frac{1}{P(BC)} = \frac{1}{P(BC)} if \ A \subset B)$$

$$\approx \frac{h(a_j | \mathbf{z}_i) \epsilon}{\sum_{k \in \mathcal{R}_{i_j}} h(a_j | \mathbf{z}_k) \epsilon} \qquad (by \quad (1) \text{ and } (2))$$

$$= \frac{\phi(i)}{\sum_{k \in \mathcal{R}_{i_j}} \phi(k)} \text{ (since } h(a_j | \mathbf{z}_i) = \phi(i)h_o(a_j)). \qquad (3)$$

Cox defines the conditional likelihood to be

$$lik = \prod_{i \in \mathcal{D}} P\{A_i | B_i\} = \prod_{i \in \mathcal{D}} \frac{\phi(i)}{\sum_{k \in \mathcal{R}_{i_j}} \phi(k)}.$$

§5.2.4. Partial likelihood approach in continuous case

Let that $a_j = M_{i_j}$, j = 1, ..., g. Let $A_{i_j} = \{\text{the } i_j \text{-th person died in } [a_j, a_j + \epsilon)\}$ and $B_{i_j} = \{\text{one death in } [a_j, a_j + \epsilon) \text{ from } \mathcal{R}_{i_j}, \text{ and each in } \mathcal{R}_{i_j} \text{ survived } a_j\}, i = 1, ..., g.$ Then define a full likelihood under assumptions 1 and 2 as $L = P\{A_{i_1}B_{i_1}\cdots A_{i_g}B_{i_g}\}.$ Since P(ABC) = P(A)P(B|A)P(C|AB),

$$\begin{split} L &= P\{A_{i_1}B_{i_1}\cdots A_{i_g}B_{i_g}\} = P\{B_{i_1}A_{i_1}\cdots B_{i_g}A_{i_g}\} \\ &= P(B_{i_1})P(A_{i_1}|B_{i_1})P(B_{i_2}|B_{i_1}A_{i_1})P(A_{i_2}|B_{i_1}A_{i_1}B_{i_2})\times\cdots \\ &\times P(B_{i_g}|B_{i_1}A_{i_1}\cdots B_{i_{g-1}}A_{i_{g-1}})P(A_{i_g}|B_{i_1}A_{i_1}\cdots B_{i_{g-1}}A_{i_{g-1}}B_{i_g}) \\ &= P(A_{i_1}|B_{i_1})P(A_{i_2}|A_{i_1}B_{i_1}B_{i_2})\cdots P(A_{i_g}|A_{i_1}\cdots A_{i_{g-1}}B_{i_1}\cdots B_{i_g}) \\ &\times P(B_{i_1})P(B_{i_2}|A_{i_1}B_{i_1})\cdots P(B_{i_g}|A_{i_1}\cdots A_{i_{g-1}}B_{i_1}\cdots B_{i_{g-1}}) \\ &= P(A_{i_1}|B_{i_1})P(A_{i_2}|B_{i_2})\cdots P(A_{i_g}|B_{i_g}) \\ &\times P(B_{i_1})P(B_{i_2}|A_{i_1}B_{i_1})\cdots P(B_{i_g}|A_{i_1}\cdots A_{i_{g-1}}B_{i_1}\cdots B_{i_{g-1}}) \end{split}$$

by independence.

$$lik = P(A_{i_1}|B_{i_1})P(A_{i_2}|B_{i_2})\cdots P(A_{i_g}|B_{i_g}).$$
 ((see (3) in §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:

$$H_o(t) = \int_0^t h_o(u) du \text{ (cts)}$$

or $\sum_{u \le t} \ln(1 - h_o(u)) \approx \sum_{u \le t} h_o(u) \text{ (discrete)}.$
 $S_o(t) = exp(-H_o(t)).$

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

$$\hat{S}_o(t) = exp(-\sum_{a_j \le t} \frac{d_j}{\sum_{l \in \mathcal{R}(a_j)} \hat{\phi}(l)}).$$
(4)

Note that if $\phi(i) = 1$,

$$\hat{S}_o(t) = exp(-\sum_{a_j \le t} \frac{d_j}{r_j}) = exp(-\sum_{a_j \le t} \hat{h}(a_j)).$$

Breslow (1972, JRSS,B) also proposes another estimator. They can be computed by R codes: library(MASS)

 $library(survival) \\ u=coxph(Surv(m,d) \sim x) \\ y=survfit(u) \\ plot(y)$

> library(MASS)

> library(splines)

> library(survival)

> attach(gehan)

> gehan

pairtime cens treat1 1 1 1 control $\mathbf{2}$ 1 10 6 - MP1 3 2221 control 4 27 6 - MP1 53 3 1 control 3 32 6 - MP6 0 $> coxph(Surv(time,cens) \sim treat,gehan,method="exact") # discrete$ exp(coef)se(coef)coef \boldsymbol{z} ptreatcontrol 1.6282 5.09490.4331 $3.759 \quad 0.00017$ Likelihood ratio test=16.25 on 1 df, p=5.544e-05n=42, number of events= 30 $> \operatorname{coxph}(\operatorname{Surv}(\operatorname{time, cens}) \sim \operatorname{treat, gehan, method} = "breslow")$ exp(coef)coef se(coef)zptreatcontrol 1.5092 4.52310.4096 $3.685 \quad 0.000229$ Likelihood ratio test=15.21 on 1 df, p=9.615e-05n=42, number of events= 30 > (x=coxph(Surv(time,cens)~treat,gehan,method="efron")) coef exp(coef)se(coef)zp1.57214.8169 0.4124 3.812 0.000138 treatcontrol Likelihood ratio test=16.35 on 1 df, p=5.261e-05n = 42, number of events = 30 > (x=coxph(Surv(time,cens)~treat)) exp(coef)coef se(coef)zptreatcontrol 1.5721 4.8169 0.4124 3.812 0.000138 Likelihood ratio test=16.35 on 1 df, p=5.261e-05n=42, number of events= 30 > summary(x) n=42, number of events= 30 coef exp(coef)se(coef)Pr(>|z|)z0.4124 treatcontrol 1.57214.81693.812 0.000138 * * * exp(coef)exp(-coef)lower.95 upper.95 4.817 2.147treatcontrol 0.207610.81Concordance = 0.69 (se = 0.041) Likelihood ratio test= 16.35 on 1 df, p=5e-05 Wald test = 14.53 on 1 df, p=1e-04 Score (logrank) test = 17.25 on 1 df, p=3e-05> (y=survfit(x, conf.type="log-log")) # baseline survival function median 0.95LCL 0.95UCL events n42 1330 8 22> summary(y)

i	time	n.risk	n.event	survival	std.err	lower 95% CI	upper95% CI
	1	42	2	0.964	0.0254	0.8604	0.991
	2	40	2	0.926	0.0367	0.8098	0.973
	3	38	1	0.907	0.0414	0.7834	0.962
	÷						
	23	7	2	0.169	0.0784	0.0516	0.344
> plc	$\operatorname{ot}(y)$						



Fig. 1. Plot of S_o under the Cox model

Cox proposes $\hat{S}_o(t) = exp(-\sum_{a_j \le t} \frac{d_j}{\sum_{l \in \mathcal{R}(a_j)} \hat{\phi}(l)}).$

Breslow (1972, JRSS,B) also proposes another estimator of S_o (computed by R codes). 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

$$\mathcal{L}(F) = \prod_{i: \ ex} [(S_o(M_i -))^{\psi(\beta, \mathbf{Z}_i)} - (S_o(M_i))^{\psi(\beta, \mathbf{Z}_i)}] \prod_{i: \ rc} (S_o(M_i))^{\psi(\beta, \mathbf{Z}_i)}.$$
 (5)

which corresponding to the Lehmann model. They are the same if X is continuous.

- §5.2.5.2. Homework. 5 $(M_i, \delta_i, \mathbf{z}_i)$'s: (2.5, 1, 2), (2, 0, 5), (4, 0, 1), (1, 1, 1), (7, 1, 2).
 - 1 Which of the three MPLEs of β through coxph() is the solution in §5.2.1.2.
 - 2 Derive S_o with three methods: (2.a) R program, (2.b) with Eq. (4) in this section, (2.c) with Eq. (5) and $\beta = \hat{\beta}$ derived in §5.2.1.2.
- \S 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 extension of the PH model to IC data. She defines the SMLE of β to be the one that maximizes the likelihood function

$$L(\beta, S_o) =$$

$$\prod_{i: \ L_i < R_i} [(S_o(L_i))^{\psi(\beta, \mathbf{Z}_i)} - (S_o(R_i))^{\psi(\beta, \mathbf{Z}_i)}] \prod_{i: \ L_i = R_i} [(S_o(L_i-))^{\psi(\beta, \mathbf{Z}_i)} - (S_o(R_i))^{\psi(\beta, \mathbf{Z}_i)}]$$
(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 $S(t) = \prod_{x_i \le t, x_i \in D_f} (1 - h(x_i)),$

 $f(t|\mathbf{z}) = h(t|\mathbf{z})S(t-|\mathbf{z})$, 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

$$\mathcal{L}(\beta, h_o) = \prod_{i: \ L_i < R_i} [S(L_i | \mathbf{z}_i) - S(R_i | \mathbf{z}_i)] \times \prod_{i: \ L_i = R_i} [h(L_i | \mathbf{z}_i) S(L_i - | \mathbf{z}_i)]$$
(5.3.2)

where
$$h(t|\mathbf{z}) = \begin{cases} \exp(\beta \mathbf{z})h_o(t) & \text{if } t < \tau \text{ or } t < M_{(n)} \\ 1 & \text{if } t = \tau \text{ or } t = M_{(n)} \end{cases}$$
,
 $S(t|\mathbf{z}) = \prod_{x_i \le t, x_i \in D_{f_o}} (1 - h(x_i|\mathbf{z})), \text{ and}$
 $S(t - |\mathbf{z}) = \prod_{x_i < t, x_i \in D_{f_o}} (1 - h(x_i|\mathbf{z})).$

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 $A_1, ..., A_m$ be the II's induced by I_i 's, the observed intervals, and $p_j = \mu_F(A_j)$. Consider S_o of form

$$S_o(t) = \sum_{j: A_j \cap (t,\infty] \neq \emptyset} p_j \cdot$$

The variance of the SMLE $\hat{\beta}$ can be estimated by a $p \times p$ matrix \hat{V}_{11}/n , where

$$\begin{pmatrix} \hat{V}_{11} & \hat{V}_{12} \\ \hat{V}_{12}^t & \hat{V}_{22} \end{pmatrix} = \hat{J}^{-1}, \qquad (J = -\frac{\partial^2 (\ln L)/n}{\partial (\beta^t, p_1, \dots, p_{m-1})^{\prime}})$$

 \hat{J} is the empirical Fisher information matrix, i.e., minus the second partial derivatives matrix of the $\frac{1}{n}$ log likelihood function with respect to $(\beta^t, p_1, ..., p_{m-1})$, with all parameters replaced by their SMLE. Note that $\ln L = O(n)$ and $(\ln L)/n = O(1)$.

$$\mathbf{Q:}\ \Sigma_{\hat{\beta}} = V_{11}\ ?\ \Sigma_{\hat{\beta}} \neq V_{11}\ ?\ \Sigma_{\hat{\beta}} \approx V_{11}\ ?\ \Sigma_{\hat{\beta}} = V_{11}/n\ ?\ \Sigma_{\hat{\beta}} \neq V_{11}/n\ ?\ \Sigma_{\hat{\beta}} \approx V_{11}/n\ ?$$

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_i 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 from two groups, corresponding to $Z_i = 0$ or 1.

Note that the baseline survival function S satisfies

S(4-) = 1, $S(4) = S(6-) = p_2 + p_3,$ $S(6) = S(8-) = p_3$ and S(8) = 0.

The likelihood is

$$\begin{split} \mathbf{L} = & [S(2) - S(5)][S(3) - S(4)][S(5)^{e^{\beta}} - S(9)^{e^{\beta}}][S(1)^{e^{\beta}} - S(6)^{e^{\beta}}][S(7) - S(8)] \\ = & p_1 p_1 (p_2 + p_3)^{e^{\beta}} (1 - p_3^{e^{\beta}}) p_3 & (simpler) \\ \mathcal{L} = & \ln \mathbf{L} = \log[p_1^2 p_3 (1 - p_1)^{e^{\beta}} (1 - p_3^{e^{\beta}})] & \text{as } p_1 + p_2 + p_3 = 1. \end{split}$$

The parameter space is $\Omega = \{(\beta, p_1, p_3) : \beta \in (-\infty, \infty), p_1 \ge 0, p_3 \ge 0, p_1 + p_3 \le 1\}$ 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).

The NR method points to $p_2 < 0$, which maximizes L without constraint (see the figure below). So it stops at the boundary.

But the SMLE needs to be searched on the boundaries: what are they ?) It reduces to 2 dimensions (α, p_i) .



Fig 2. Illustration of the drawback of NR method.

A feasible algorithm for the GMLE.

It can be illustrated by the next figure.

Abusing notations, we identify S with a vector $(S_1, ..., S_m)$. Similarly, we identify $S^{(i)}$ with $(S_1^{(i)}, ..., S_m^{(i)})$.



Fig. 3. Illustration of the S-step in the feasible algorithm

Step 0. Let $b^{(0)} = 0$ be the initial estimate of β and the GMLE of a survival function with observations $(L_j, R_j), j = 1, ..., n$ be the initial estimate of $S^{(0)}$.

- Step i + 1 $(i \ge 0)$. 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 $S (= (S_1, ..., S_m))$ so that the log likelihood function $\mathcal{L}(\cdot, b^{(i+1)})$ is maximized (or increases). Since $S_j = p_{j+1} + \dots + p_m$ for some j, let $p^{(i+1),0} = p^{(i)}$. At Sub-step j (j = 1, ..., m), update $(p_1, ..., p_m)$ by $(p_1^{(i+1),j}, ..., p_m^{(i+1),j})$, where $p_h^{(i+1),j} = p_{h,u_o}$ and $p_{h,u} = \begin{cases} \frac{p_h^{(i+1),j-1} + u}{1+u} & \text{if } h = j, \\ \frac{p_h^{(i+1),j-1}}{1+u} & \text{if } h \neq j, \end{cases}$

 $h = 1, ..., m, u_o > 0$ is a number maximizing $L(b^{(i+1)}, S_{\cdot,u})$ where $S_{\cdot,u} = (S_{1,u}, ..., S_{m,u})$ and $S_{i,u} = p_{i+1,u} + \cdots + p_{m,u}$. Note: If such u_o is difficult to choose, one may choose a u_o satisfying

$$L(b^{(i+1)}, S^{(i+1),j}) > L(b^{(i+1)}, S^{(i+1),j-1}).$$
(3.1)

In particular, if $\frac{\partial}{\partial u} \ln L(b^{(i+1)}, S_{\cdot,u}) \Big|_{u=0} > 0$, $u_o = c^k \frac{\partial}{\partial u} \ln L(b^{(i+1)}, S_{\cdot,u}) \Big|_{u=0}$, where $S_{\cdot,u} = (S_{1,u}, ..., S_{m,u}), c \in (0.1)$, 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)	simplification
(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	$(1-e^eta p_1)-0$	$1 - e^{eta} p_1$
(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 (AL) model and regression analysis with IC data The AL model (or the log linear regression model) assumes 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,

$$\ln X = \ln T - \beta z \qquad (often written as \ln Y = \beta X + \alpha + \epsilon)$$

lnT has a distribution which does not depend on β and Z, $E(\ln T) = \alpha$ and $\sigma_{\ln T}^2 = \sigma^2$. For simplicity, we shall replace β , lnX and lnT by $-\beta$, X and $\alpha + \epsilon$, respectively. That is, the ordinary linear regression set-up:

$$X = \beta z + \alpha + \epsilon$$
, where $E(\epsilon) = 0$ and $Var(\epsilon) = \sigma^2$

Here $\beta \mathbf{z}$ is interpreted as $\beta' \mathbf{z}$. For simplicity, we only consider the univariate case. Question: $(\alpha, \beta) = ?$ If one has complete data, $(X_i, z_i), i = 1, ..., n$,

the least squares estimate (LSE) 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 $\mathbf{X}^t = (X_1, ..., X_n)$.

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 $W_i(\alpha,\beta) = X_i - \alpha - \beta z_i, i = 1, ..., n$. Moreover, the LSE is also the MLE under the normal distribution.

Thus it is the solution to the equation $\frac{\partial \ln \mathbf{L}}{\partial \beta} = 0$ and $\frac{\partial \ln \mathbf{L}}{\partial \alpha} = 0$ under $N(\mu, \sigma)$. **Definition.** $S_o(t) = S_{X|\mathbf{Z}}(t|0)$ is called the baseline survival function. Let $W = \alpha + \epsilon$, $\tau_o = \tau_{C_o}$, where $C_o = Y - \beta \mathbf{Z}$ and $\tau_{C_o} = \sup\{t : F_{C_o}(t) < 1\}$.

Under the standard right censorship model, we say that the parameter is identifiable if $S(t - \mathbf{bz}) = S_o(t - \beta \mathbf{z}) \forall \mathbf{z}$ and $\forall t \in \mathcal{D}$, where

$$\mathcal{D} = \begin{cases} \{t : t < \tau_o\} & \text{if } P(Y = \tau_Y) = 0 < S_W(\tau_o -) \\ \{t : t \le \tau_o\} & \text{if } P(Y = \tau_Y)S_W(\tau_o) > 0 \\ (-\infty, \infty) & \text{if } S_W(\tau_o -) = 0 \text{ or } P(Y = \tau_Y)S_W(\tau_o -) > 0 = S_W(\tau_o) \\ => (S(t), \mathbf{b}) = (S_o(t), \beta) \ \forall \ t \in \mathcal{D}, \text{ where } \tau_M = \sup\{t : S_M(t) > 0\}. \end{cases}$$

=> $(S(t), \mathbf{b}) = (S_o(t), \beta) \forall t \in \mathcal{D}$, where $\tau_M = \sup\{t : S_M(t) > 0\}$. **Example 1.** Suppose $E(X_i) = \mu + \alpha_i$, i = 1, 2. Then $(\mu, \alpha_1, \alpha_2)$ is not identifiable. *e.g*, $E(X_1) = 2$ and $E(X_2) = 5$ yield $(\mu, \alpha_1, \alpha_2) = (0, 2, 5)$ or (2, 0, 3), not uniquely determined. **Theorem 1.** (Yu and Dong (2019)) Suppose that $\tau_o < \infty$. Then

(a) The survival function $S_W(t)$ is identifiable iff $t \in \mathcal{D}$.

(b) The parameter β is identifiable iff $\mathcal{B}_{\mathbf{Z}_0} \neq \emptyset$, where $\mathbf{z}_0 \in \mathcal{D}_{\mathbf{Z}}$ such that $\tau_Y - \beta \mathbf{z}_o = \tau_o$ and

$$\mathcal{B}_{\mathbf{Z}_{0}} = \left\{ (w_{1}, \mathbf{z}_{1}, ..., w_{p}, \mathbf{z}_{p}) : \mathbf{z}_{1} - \mathbf{z}_{0}, ..., \mathbf{z}_{p} - \mathbf{z}_{0} \text{ are linearly independent,} \\ \mathbf{z}_{i} \in \mathcal{D}_{\mathbf{Z}}, w_{i} \in \mathcal{D}_{W} \text{ and } w_{i} + \beta' \mathbf{z}_{i} \left\{ \leq \tau_{Y} & \text{if } P(Y = \tau_{Y}) > 0 \\ < \tau_{Y} & \text{otherwise} \end{array} \right\}$$

Theorem 2. (Yu and Dong (2019)) Suppose that $\tau_o = \infty$.

(a) The survival function $S_W(t)$ is identifiable for each t.

(b) The parameter β is identifiable iff $\exists \mathbf{z}_0 \in \mathcal{D}_{\mathbf{Z}}$ such that $\mathcal{B}_{\mathbf{Z}_0} \neq \emptyset$, where

 $\mathcal{B}_{\mathbf{Z}_0} = \{ (\mathbf{z}_1, ..., \mathbf{z}_p): \ \mathbf{z}_1 - \mathbf{z}_0, ..., \mathbf{z}_p - \mathbf{z}_0 \text{ are linearly independent, and } \mathbf{z}_i \in \mathcal{D}_{\mathbf{Z}} \}. \text{ Here } \mathbf{z}_0 = \mathbf{0} \text{ if } \mathbf{0} \in \mathcal{D}_{\mathbf{Z}}, \text{ otherwise } \mathbf{z}_0, ..., \mathbf{z}_p \text{ are linearly independent vectors belonging to } \mathcal{D}_{\mathbf{Z}}.$

Remark. Under the standard right censorship model, a sufficient condition for identifiability is $S_o(\tau_M) = 0$ and \mathbf{Z} takes on p + 1 distinct values $\mathbf{0}, \mathbf{z}_1, ..., \mathbf{z}_p$ and $\mathbf{z}_1, ..., \mathbf{z}_p$ are linearly independent.

Under right censoring, there are several estimators:

- (1) Miller's estimator. (1976, Biometrika).
- (2) Buckley-James estimator (1979, Biometrika).
- (3) M-estimator approach (Zhang and Li (1996, Annals of Statistics)

(4) SMLE (Yu and Wong (2003, JSCS))

(5) Modified SMLE (Yu and Wong (2005, Technometics))

There are three issues for each estimator:

- (1) How to justify the estimator ?
- (2) How to derive the estimate ?
- (3) Is it consistent or efficient ?

§5.4.1. Miller Estimator (1976, Biometrika). The estimator of (α, β) with RC data is $(\hat{\alpha}, \hat{\beta}) = argmin_{a,\mathbf{b}}r(a, \mathbf{b})$ (that minimizes $r(a, \mathbf{b})$ over all (a, \mathbf{b})), where

$$r(a, \mathbf{b}) = n \int t^2 d\hat{F}_{a, \mathbf{b}}(t), \ \hat{F} = \hat{F}_{a, \mathbf{b}}$$
 is the GMLE of F_{ϵ} based on $(W_i(a, \mathbf{b}), \delta_i)$'s,

and $W_i(a,b) = M_i - a - b\mathbf{z}_i$.

If $F_X(\tau_Y) < 1$, then Miller's estimator is not consistent, as $\alpha = E(\epsilon)$ is not identifiable.

Let $s_i, ..., s_m$ 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}_o 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 $\hat{F}_\alpha(t) = \hat{F}_o(t + \alpha)$. Let $\eta_1, ..., \eta_m$ 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, ..., 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

$$\alpha = \hat{\alpha}(\beta) = \sum_{i=1}^{n} w_i(\beta) (M_i - \beta z_i).$$
(1.0)

Thus it suffices to search $\hat{\beta}$ that minimizes

$$H(\beta) = \sum_{j=1}^{n} w_j(\beta) [M_j - \hat{\alpha}(\beta) - \beta z_i]^2.$$
 (1.1)

Miller suggested the following iterative procedure:

1. Assign an initial value $\beta = \frac{\sum_{i: ex} X_i(z_i - \overline{z}_u)}{\sum_{j: ex} (z_j - \overline{z}_u)^2}$, where $\overline{z}_u (\overline{X}_u)$ 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 $\frac{\sum_{i=1}^{n} w_i(\beta) X_i(z_i \overline{z}_w)}{\sum_{j=1}^{n} w_j(\beta) (z_j \overline{z}_w)^2}$, where $\overline{z}_w = \sum_{j=1}^{n} w_j z_j$.
- 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

$$\hat{\sigma}_{\hat{\beta}}^2 = \frac{\sum_{j=1}^n w_j(\hat{\beta})(M_i - \hat{\alpha} - \hat{\beta} z_i)^2}{n \sum_{j=1}^n w_j(\hat{\beta})(z_i - \overline{z}_w)^2}, \text{ comparing to } \sigma^2 / \sum_{i=1}^n (z_i - \overline{z})^2$$

The consistency and asymptotic normality were considered under the assumption that

$$P\{X \text{ is not censored} | X = t\} > 0 \text{ for all possible } t$$
(1.2)

and the censoring distribution is of form

$$G(y|z) = G_o(y - \beta z)$$
, where G_o 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).

First review the LR model with complete data. The LR model:

 $X_i = \alpha + \beta' z_i + \epsilon_i, i = 1, ..., n, \beta \in \mathbb{R}^p$, with i.i.d. observations (X_i, z_i) 's. The common estimator is the LSE, which is the solution to

$$\operatorname{argmin}_{a,\mathbf{b}} \sum_{i=1}^{n} (X_i - a - bz_i)^2 \quad (= \operatorname{argmin}_{a,\mathbf{b}} n \int t^2 d\hat{F}_{a,\mathbf{b}}(t)),$$

where $\hat{F} = \hat{F}_{a,\mathbf{b}}$ is the GMLE based on $X_i - a - b\mathbf{z}_i$'s. Or the LSE of β is the root of

$$H(b) = \sum_{i=1}^{n} (X_i - \overline{X} - b(\mathbf{z}_i - \overline{\mathbf{z}}))(\mathbf{z}_i - \overline{\mathbf{z}}) \quad \text{with } a = \overline{X} - b\overline{\mathbf{z}}$$

$$H(b) = \sum_{i=1}^{n} (X_i - b\mathbf{z}_i)(\mathbf{z}_i - \overline{\mathbf{z}}) \quad (as \quad \sum_{i=1}^{n} (\overline{X} - b\overline{\mathbf{z}})(\mathbf{z}_i - \overline{\mathbf{z}}) = 0).$$

Given RC data, $(M_i, \delta_i, \mathbf{z}_i)$ s, Buckley and James propose to estimate β by the "root" of

$$H(\mathbf{b}) = \sum_{i=1}^{n} (\hat{X}_{i}^{*} - b\mathbf{z}_{i})(\mathbf{z}_{i} - \overline{\mathbf{z}}), \text{ where}$$
$$\hat{X}_{i}^{*} = \hat{X}_{i}^{*}(\mathbf{b}) = M_{i}\delta_{i} + (1 - \delta_{i})[\mathbf{b}'\mathbf{z}_{i} + \frac{\sum_{t \in \mathcal{A}_{i}} t\hat{f}_{\mathbf{b}}(t)}{\hat{S}_{\mathbf{b}}(T_{i})}], \qquad (2.1)$$
$$\mathcal{A}_{i} = \{t : t > T_{i}, \hat{f}_{\mathbf{b}}(t) > 0\},$$

 $\hat{S}_{\mathbf{b}}$ is the PLE of $S_{T(\mathbf{b})}$, the survival function based on $(T_i(\mathbf{b}), \delta_i), i = 1, ..., n,$ $T_i = T_i(\mathbf{b}) = M_i - \mathbf{b'}\mathbf{z}_i$, with $\mathbf{b} \in \mathcal{R}^p$, $(T(\beta) = X - \beta \mathbf{Z} = \alpha + \epsilon \text{ if } \delta = 1)$, $\hat{f}_{\mathbf{b}}$ and $\hat{F}_{\mathbf{b}}$ the PLE's of the density function and the cdf, respectively.

It is worth mentioning that $\hat{S}_{\mathbf{b}}$ depends on **b** as T_i 's depend on **b**.

If the largest observation $T_{(n)}$ among T_i 's is right censored, then treat $T_{(n)}$ as uncensored (suggested by Buckley and James), or set $\hat{S}_{\mathbf{b}}(T_{(n)} + 1) = 0$ if

 $\delta_{(n)} = 0.$

The root of H(b) is called the Buckley and James estimator (BJE), denoted by $\hat{\beta}$ and $\hat{\alpha} = \overline{\hat{X}^*} - \hat{\beta}' \overline{\mathbf{z}}$, where $\overline{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{z}_i$ and $\hat{\mathbf{X}}^* = \hat{\mathbf{X}}^*(\mathbf{b}) = (\hat{X}_1^*, ..., \hat{X}_n^*)'$. **Remark.** Buckley and James point out that the solution may not exist. In the latter case,

Remark. Buckley and James point out that the solution may not exist. In the latter case, the BJE of β is modified as the zero-crossing point of H(b) if p = 1 (Lai and Ying (1991)). One says that \hat{b} is a zero-crossing point (zcp) of a function H if H(b) changes its sign at $b = \hat{b}$. Then the BJE is well defined, though it may not be unique. Question: If $H(t) = \mathbf{1}(t \leq 1)$, what is a zcp of H?

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

$$E(M|Z) \neq \alpha + \beta Z.$$

$$Define \qquad X^* = X\delta + (1-\delta)E(X|X > Y) \qquad (= X\delta + (1-\delta)E(X|\delta = 0) \quad (2.2)$$
$$\neq \hat{X}_i^* = M_i\delta_i + (1-\delta_i)[\mathbf{b}'\mathbf{z}_i + \frac{\sum_{t \in \mathcal{A}_i} t\hat{f}_{\mathbf{b}}(t)}{\hat{S}_{\mathbf{b}}(T_i)}]).$$

Then

$$E(X^*|Z) = \alpha + \beta Z.$$

Reason: Abusing notation, write $E_z(W) = E(W|Z)$, where W is a random variable. Then

$$E_z(X^*) = E_z(E_z(X^*|\delta))) = P(\delta = 1)E_z(X^*|\delta = 1) + P(\delta = 0)E_z(X^*|\delta = 0)$$

$$=P(\delta = 1)E_{z}(X|\delta = 1) + P(\delta = 0)E_{z}(E_{z}(X|X > Y)|\delta = 0)$$
(2.3)
$$=P(\delta = 1)E_{z}(X|\delta = 1) + P(\delta = 0)E_{z}(E_{z}(X|\delta = 0)|\delta = 0)$$
(by (2.2))
$$=P(\delta = 1)E_{z}(X|\delta = 1) + P(\delta = 0)E_{z}(X|\delta = 0)$$

$$=E_{z}(E_{z}(X|\delta))$$
$$=E_{z}(X|\delta) = E_{z}(X|\delta) = E_{z}(X|\delta)$$

If we could observe $X_1^*, ..., X_n^*$, then $H(b) = \sum_{i=1}^n (X_i^* - bz_i)(z_i - \overline{z}) = 0$ leads to

the "LSE"
$$\hat{\beta} = \frac{\sum_{i=1}^{n} X_i^* (z_i - \overline{z})}{\sum_{j=1}^{n} (z_j - \overline{z})^2}$$
 and $\hat{\alpha} = \overline{X^*} - \hat{\beta}\overline{z}.$ (2.4)

Since we cannot observe all X_i^* 's, in view of (2.3) we replace X_i^* 's by their predictors

$$\hat{X}_{i}^{*} = X_{i}\delta_{i} + \hat{E}(X_{i}|X_{i} > Y_{i})(1 - \delta_{i}) \qquad \text{(which is (2.1)), where}$$

$$\hat{E}(X_{i}|X_{i} > Y_{i}) = \hat{\beta}z_{i} + \frac{\sum_{t > M_{i} - \hat{\beta}z_{i}} t\hat{f}_{\hat{\beta}}(t)}{\hat{S}(M_{i} - \hat{\beta}z_{i})}, \qquad (2.5)$$

$$E(\beta \mathbf{Z} + W|W > Y_{i} - \beta z_{i}) \qquad (W = \alpha + \epsilon)$$

and $\hat{S}_{\hat{\beta}}$ is the PLE of S_W based on observations $(M_i - \hat{\beta} z_i, \delta_i)$'s, though $\hat{\beta}$ is an estimate. (2) An M-estimator based on $N(\mu, \sigma)$.

An M-estimator is the solution to $\frac{\partial \ln \underline{\mathbf{L}}}{\partial \beta} = 0$. The likelihood function for given RC data is

$$\mathbf{L} = \prod_{i=1}^{n} (f(T_i))^{\delta_i} (S(T_i))^{1-\delta_i}.$$
$$\frac{\partial \mathrm{ln}\mathbf{L}}{\partial \beta} = \sum_{i=1}^{n} \{\delta_i \frac{f'}{f}(T_i) + (1-\delta_i) \frac{-f}{S}(T_i)\}(-\mathbf{z}_i). \tag{A}$$

Under $N(\mu, \sigma^2)$ assumption, $f(t) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(t-\mu)^2}{2\sigma^2}),$ $\frac{f'}{\epsilon}(t) = -\frac{t-\mu}{2}.$

$$\frac{f(t)}{f(t)} = -\frac{1}{\sigma^2},$$

$$f(t) = -\int_t^\infty f'(x)dx = -\int_t^\infty \frac{f'}{f}(x)f(x)dx = \int_t^\infty \frac{x-\mu}{\sigma^2}f(x)dx.$$
 Then Eq. (A) yields

$$\frac{\partial \ln \mathcal{L}}{\partial \beta} = \sum_{i=1}^{n} \{ \delta_i \frac{(T_i - \mu)}{\sigma^2} + (1 - \delta_i) \frac{\int_{T_i}^{\infty} (x - \mu) f(x) dx}{\sigma^2 S(T_i)} \} \mathbf{z}_i.$$
(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 \overline{T}^* and multiplying σ^2 .

In other words, (2.1) is the same as

$$H(\mathbf{b}) = \sum_{i=1}^{n} (T_i^* - \overline{T^*}) \mathbf{z}_i, \text{ where } T_i^* = \hat{X}_i^* - \mathbf{b} \mathbf{z}_i.$$

Notice

$$H(b) = \sum_{i=1}^{n} (X_i - \overline{X} - b(\mathbf{z}_i - \overline{\mathbf{z}}))(\mathbf{z}_i - \overline{\mathbf{z}})$$

=
$$\sum_{i=1}^{n} (X_i - b\mathbf{z}_i))(\mathbf{z}_i - \overline{z})$$

=
$$\sum_{i=1}^{n} (X_i - \overline{X} - b(\mathbf{z}_i - \overline{\mathbf{z}}))\mathbf{z}_i \qquad (\sum_{i=1}^{n} (X_i - \overline{X} - b(\mathbf{z}_i - \overline{\mathbf{z}}))\overline{\mathbf{z}} = 0)$$

The consistency and the asymptotic properties of the BJE have been established under certain continuous assumptions by Lai and Ying (1991) and

under certain discrete assumptions by Kong and Yu (2006).

In addition, if $\epsilon \sim N(\mu, \sigma^2)$, then the BJE is asymptotically efficient, just like the LSE. Otherwise, it is not efficient.

On the other hand, under certain discontinuous assumptions,

the BJE may not have asymptotic normal distribution (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 $\epsilon_1, ..., \epsilon_n$ are i.i.d. with variance σ^2 , given U (or \mathbf{z}_i s), the covariance matrix of the LSE is

$$\Sigma_{\hat{\beta}} = \sigma^2 (U'U)^{-1}.$$
 (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 (2006) present 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

$$I(\hat{\beta}) = \sum_{i=1}^{n} (T_i^* - \overline{T^*}) (\mathbf{z}_i - \overline{\mathbf{z}}) \{ (T_i^* - \overline{T^*}) (\mathbf{z}_i - \overline{\mathbf{z}}) \}' / \hat{\sigma}^4$$

and $\hat{\sigma}^2$ is an estimator of σ^2 .

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

$$\hat{\sigma} = \frac{Tq - T \cdot q}{\overline{q^2} - (\overline{q})^2},$$

where $\overline{T} = \frac{1}{m} \sum_{i=1}^n \delta_i T_i,$
 $\overline{q} = \frac{1}{m} \sum_{i=1}^n \delta_i q_i,$
 $\overline{q^2} = \frac{1}{m} \sum_{i=1}^n \delta_i q_i^2$ and
 $\overline{Tq} = \frac{1}{m} \sum_{i=1}^n \delta_i T_i q_i.$

An estimator of the variance of $\hat{\beta}$ given by Buckley and James is

$$\hat{\sigma}_{\hat{\beta}}^{2} = \frac{\hat{\sigma}_{u}^{2}}{\sum_{i: ex} (z_{i} - \overline{z}_{u})^{2}},$$
where $\hat{\sigma}_{u}^{2} = \frac{1}{\sum_{i=1}^{n} \delta_{i} - 2} \sum_{i: ex} [X_{i} - \overline{X}_{u} - \hat{\beta}(z_{i} - \overline{z}_{u})]^{2},$
 $\overline{X}_{u} = \frac{\sum_{i: ex} X_{i}}{\sum_{i=1}^{n} \delta_{i}} \text{ and } \overline{z}_{u} = \frac{\sum_{i: ex} z_{i}}{\sum_{i=1}^{n} \delta_{i}}.$

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 $\tilde{\sigma}_{\epsilon} = \int (t - \mu_{\epsilon,u})^2 d\hat{F}_{\hat{\beta}}(t),$

$$\begin{split} \mu_{\epsilon,u} &= \int t dF_{\hat{\beta}}(t),\\ \tilde{\Sigma}_z^2 &= \int_{t < \infty} (z(t) - \mu_{z,u}) (z(t) - \mu_{z,u})' d\hat{F}_{\hat{\beta}}(t),\\ \mu_{z,u} &= \int z(t) d\hat{F}_{\hat{\beta}}(t), \end{split}$$

 $\hat{F}_{\hat{\beta}}$ is the modified PLE that moves the weight from $+\infty$ to the largest observation, and

 $z(t) = \text{average of } z \text{ in } \{z_i : T_i(\hat{\beta}) = t, i = 1, ..., n\}.$

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: $\hat{\alpha} = \overline{\hat{X}^*} \hat{\beta}\overline{z}$ and $\hat{\beta} = \frac{\sum_{i=1}^n X_i^*(z_i \overline{z})}{\sum_{j=1}^n (z_j \overline{z})^2}$.
- 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 BJE exits (see Example 1 below).

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(MASS), library(survival), library(rms) > set.seed(1000) > fun=function() { # prepare a function for computing the BJE $y=b^*x+w$ d = ifelse(y > c, 0, 1) $m = y^*d + c^*(1-d)$ $f=bj(Surv(m, d) \sim x, link="identity", control=list(iter.max=50))$ return(f(scoef[2]))> n=20 # Set sample size. > b=1 # Set β value. > c = 1 \S A.1. Example 1. The algorithm yields a point which is the BJE. > w = rbinom(n, 1, 0.5)> x = rbinom(n, 1, 0.5)> fun()CoefS.E.Wald ZPr(>|Z|)Intercept 0.4444 0.15922.790.00520.9556 0.26643.590.0003 xSA.2. Example 2. The algorithm cannot give an answer to a BJE.) > w=rbinom(n,1,0.5) > x = rbinom(n, 1, 0.5)> fun()No convergence in 80 steps Failure in bj.fit \$fail [1] TRUE **Note:** The BJE of (a, b) is a zero crossing point of b: BJE=1. (the data set is [1,] 1 1 0[2,] 1 1 0[3,] 1 1 1[4,] 1 0 1[5,] 1 1 0[6,] 1 0 1[7,] 1 1 0[8,] 1 0 1[9,] 0 1 0[10,] 1 0 1[11,] 1 0 1[12,] 1 0 1

[13,] 0 1 0[14,] 1 1 0

[15,] 1 0 1

 $\begin{array}{c} [16,] \ 1 \ 1 \ 0 \\ [17,] \ 1 \ 1 \ 0 \\ [18,] \ 1 \ 1 \ 0 \\ [19,] \ 1 \ 1 \ 0 \\ [20,] \ 0 \ 1 \ 0 \end{array}$

§A.3. Example 3. The algorithm oscillates and yields a point which is not the BJE. # Then we generate another set of data and try again.

> n=20> w = runif(n, 0, 1)> x = rbinom(n, 1, 0.5)*0.5> c = 0.9> fun()Cycle period = 2No convergence in 52 steps, but cycle found - average beta returned > (bj(Surv(m, d) ~ x, link="identity", control=list(iter.max=50))) CoefS.E.Wald ZPr(>|Z|)Intercept 0.46550.11853.93< 0.0001x1.17300.39312.980.0028 which is not the BJE, the BJE of (a,b) is a zero-crossing point (0.191187, 1.158995),

which is based on

the data (stored in qyu/data/reg/dis/new/simubin20)

0.9000000 0 0.5 0.22477293 1 0.0 0.90000000 0 0.50.42949612 1 0.0 $0.83872464 \ 1 \ 0.5$ 0.9000000 0 0.50.9000000 0 0.50.9000000 0 0.50.87571891 1 0.0 0.90000000 0 0.5 $0.09953812 \ 1 \ 0.0$ $0.52228869 \ 1 \ 0.5$ 0.80450739 1 0.0 $0.32050270\ 1\ 0.0$ $0.84096155 \ 1 \ 0.5$ 0.9000000 0 0.50.90000000 0 0.50.61437741 1 0.5

 $0.50425090 \ 1 \ 0.0 \\ 0.9000000 \ 0 \ 0.5$

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 $q_1 < \cdots < q_m$ 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, ..., m, first compute the PLE \hat{S}_b for a $b \in (q_h, q_{h+1})$. e.g., (the midpoint of the interval if 0 < i < m

$$let b = \begin{cases} q_1 - 1 & \text{if } i = 0\\ q_m + 1 & \text{if } i = m. \end{cases}$$

Then compute
$$(M_i^*(b), z_i^*(b))'s$$
 and $\hat{b}_h = \frac{\sum_{j=1}^n (z_j - \overline{z})M_j^*}{\sum_{k=1}^n (z_k - \overline{z})z_k^*}$, where
 $M_i^* = M_i\delta_i + (1 - \delta_i) \frac{\sum_{t>T_i(b)} \hat{f}_b(t) \frac{\sum_{j=1}^n M_j \mathbf{1}_{(T_j(b)=t,\delta_j=1)}}{\sum_{k=1}^n \mathbf{1}_{(T_k(b)=t,\delta_k=1)}}}{\hat{S}_b(T_i(b))}$
 $z_i^* = \delta_i z_i + (1 - \delta_i) [\frac{\sum_{t>T_i(b)} \hat{f}_b(t) \frac{\sum_{j=1}^n z_j \mathbf{1}_{(T_j(b)=t,\delta_j=1)}}{\sum_{k=1}^n \mathbf{1}_{(T_k(b)=t,\delta_k=1)}}}{\hat{S}_b(T_i(b))}]$

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 $H(q_i+)$ (use \hat{S}_b for $b \in (q_i, q_{i+1})$, $H(q_i)$ and $H(q_i-)$ (use \hat{S}_b for $b \in (q_{i-1}, q_i)$, i = 1, ..., m.

$$H(b) = \sum_{i=1}^{n} (M_i^* - bz_i^*)(z_i - \overline{z})$$
(2.8)

If $H(q_i)H(q_i+) \leq 0$, or $H(q_i) = 0$, 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 Eq. (2.8) rather than

$$H(b) = \sum_{i=1}^{n} (\hat{X}_{i}^{*}(b) - bz_{i})(z_{i} - \overline{z}) \text{ where } \hat{X}_{i}^{*}(b) = M_{i}\delta_{i} + (1 - \delta_{i})[bz_{i} + \frac{\sum_{t \in \mathcal{A}_{i}} t\hat{f}_{b}(t)}{\hat{S}_{b}(T_{i})}]$$
(2.9)

because even though (M_i^*, z_i^*) depends on b, it is constant in b in each interval (q_j, q_{j+1}) . **Remark.** We shall make a correction on a mistake in the statement in Yu and Wong (2003). The NPMLE $\hat{S}_b(T_i(b))$ is constant in (b_{i-1}, b_i) , where b_1, b_2, \ldots are the ordered solutions to $T_k(b) = T_j(b)$ with $(\delta_k, \delta_j) \neq (0, 0)$, not only $\delta_k \neq \delta_j$ as stated in Yu and Wong (2003) (see Example 4). However, the above algorithm still works. The reason is as follows.

(1) (M_i^*, z_i^*) only need to be modified when $\delta_i = 0$, and the modification depends on $\hat{S}_b(T_i(b))$ and $\hat{f}_b(T_j(b))$ where $T_j(b) > T_i(b)$.

(2) If $\delta_i = \delta_j = 1$ and $T_i(q_k) = T_j(q_k) \ (\neq T_h(q_k) \ \forall h \notin \{i, j\})$, then $\hat{S}_b(T_i(b))$ is constant if $\delta_i = 0$, and $\hat{f}_b(T_i(q_k+)) + \hat{f}_b(T_j(q_k-)) = \hat{f}_b(T_i(q_k))$.

Example. Let $(M_1, ..., M_4) = (1, 2, 3, 4), (z_1, ..., z_4) = (0, 0, 1, 1), (\delta_1, ..., \delta_4) = (1, 0, 0, 1),$ $\mathbf{T}(b) = (T_1(b), ..., T_4(b)) \text{ and } \hat{S}_b(\mathbf{T}(b)) = (\hat{S}_b(T_1), ..., \hat{S}_b(T_4)).$ $M_i - bz_i = M_j - bz_j$ leads to $b \in \{3 - 1, 4 - 1, 4 - 2\} = \{2, 3\}$ if $(\delta_i, \delta_j) \neq (0, 0)$; but it leads only to b = 2 if $\delta_i \neq \delta_j$. Since $T_1(b) = 1, T_2(b) = 2, T_3(b) = 3 - b$ and $T_4(b) = 4 - b$,

$$\mathbf{T}(b) = \begin{cases} (1, 2+, 2+, 3) & \text{if } b = 1\\ (1, 2+, 1+, 2) & \text{if } b = 2\\ (1, 2+, 0.5+, 1.5) & \text{if } b = 2.5\\ (1, 2+, 0+, 1) & \text{if } b = 3\\ (1, 2+, -1+, 0) & \text{if } b = 4 \end{cases} \text{ and } \hat{S}_b(\mathbf{T}(b)) = \begin{cases} (\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, 0) & \text{if } b < 2\\ (\frac{1}{3}, \frac{3}{8}, \frac{3}{4}, \frac{3}{8}) & \text{if } b = 2\\ (\frac{2}{3}, \frac{1}{3}, 1, \frac{1}{3}) & \text{if } b \in (2, 3)\\ (\frac{1}{3}, \frac{1}{3}, 1, \frac{3}{3}) & \text{if } b = 3\\ (\frac{1}{3}, \frac{1}{3}, 1, \frac{3}{3}) & \text{if } b > 3 \end{cases}$$

which is not constant in $(2, \infty)$ (as claimed in Step 1 in Yu and Wong (2002a)). To find the BJE, just check H(b) for $b \in \{2, 3, ... ??\}$ instead of checking all $b \in (-\infty, \infty)$ (see (2.9))!

Remark. If data are complete, both the BJE $\hat{\beta}_{BJ}$ and the Miller estimator $\tilde{\beta}_M$ reduce to the LSE. In fact, $\tilde{\beta}_M$ = argmin_b $\int t^2 d\hat{F}_b(t)$, where \hat{F}_b is the PLE of F_o based on $X_i - \overline{X} - b(z_i - \overline{z})$'s is an extension of $\operatorname{argmin}_{a,b} \sum_{i=1}^n (X_i - a - bz_i)^2/n$, and $b = \hat{\beta}_{BJ}$ is the zero-crossing point of H(b) and H(b) is the extension of $\frac{d}{db} \sum_{i=1}^n (X_i - a - bz_i)^2$.

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 $^{\circ}$ 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.



Fig. 5.1. MSMLE vs. BJE For Insulation Data

If one sets $\hat{f}_b(M_{(n)}(b)) = 0$ when $\delta_{(n)} = 0$, rather than $\hat{f}_b(M_{(n)}(b)) = \hat{S}_b(M_{(n)}(b))$, 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).

If one sets $\hat{f}_b(M_{(n)}(b)) = \hat{S}_b(M_{(n)}(b))$ (when $\delta_{(n)} = 0$), then there is just unique BJE $\hat{\beta} = -0.0193$ and it is a root of H. The current R program yields the same value.

Table 4.1 presents two data sets, of sample size 30 each, generated from simulation. Both set $\hat{f}_b(M_{(n)}(b)) = \hat{S}_b(M_{(n)}(b))$ if $\delta_{(n)} = 0$.

The first data set has no zero-crossing point of the sum of least squares

and has exactly one solution to Eq. (2.3) with (a, b) = (5.750711, 0.807222)The current R-program converges, and yields (6.7148, 0.7999)

The second one does not have a solution to Eq. (2.3),

but has a unique zero-crossing point at (a, b) = (7.237446, 1.078106). The current R-program cannot converges, and yields (9.7784, 1.0734).

	0	\		/ 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/
		Tabl	~	1 1 Cimerala	1:00	. Freemale

 Table 4.1. Simulation Examples

§5.4.2.2. Homework

- 1. There are 4 observations (M_i, δ_i, 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

 $X_i^* = E(X_i | X_i \in I_i)$, where I_i is the *i*-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).

3. Derive the most possible "observations" and most possible BJE's under the assumptions in the simulation example A.2 when n = 1000.

$\S5.4.3.$ An M-estimation approach

Huber (1964) proposed an M-estimator

which is a zero point of a score function $\sum_{i=1}^{n} \psi(\theta, T_i)$ (e.g., $= \frac{\partial}{\partial \theta} \ln L$, but not always) where θ is the parameter of interest and T_i 's are observations. Modifying Huber's M-estimation,

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 $\frac{\partial}{\partial \mathbf{b}} \ln \mathbf{L}$ in **b**. We shall first illustrate via RC data.

Note that the likelihood function can be written as

$$\mathcal{L}(b, S, f) = \prod_{i=1}^{n} (f(M_i - b'z_i))^{\delta_i} (S(M_i - b'z_i))^{1-\delta_i}.$$

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 $\frac{\partial \mathbf{L}}{\partial b} = 0$.

Moreover, to eliminate the effect of α , one needs to centralize z_i in L. Thus the derivative of $-\ln L$ is

$$\Phi = \sum_{i=1}^{n} (z_i - \overline{z}) \left(\delta_i(\frac{f'}{f}) (M_i - bz_i) - (1 - \delta_i)(\frac{f}{S}) (M_i - bz_i) \right).$$
(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

$$\frac{f}{S}(t) = \frac{-\int_{x>t} \frac{f'(x)}{f(x)} f(x) dx}{S(t)} = \frac{-\int_{x>t} \frac{f'(x)}{f(x)} dF(x)}{S(t)} = \frac{\int_{x>t} \frac{f'(x)}{f(x)} dS(x)}{S(t)}$$

Thus
$$\Phi = \Phi(b, S, \frac{f'}{f}) = \sum_{i=1}^{n} (z_i - \overline{z}) \left(\delta_i(\frac{f'}{f})(T_i(b)) - (1 - \delta_i) \frac{\int_{x > T_i(b)} \frac{f'(x)}{f(x)} dS(x)}{S(T_i(b))} \right),$$

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

$$\Phi = -\sum_{i=1}^{n} (z_i - \overline{z}) \left(\delta_i(T_i(b)) - (1 - \delta_i) \frac{\int_{x > T_i(b)} x dS_b(x)}{\hat{S}_b(T_i(b))} \right)$$

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

$$\Phi(b, \hat{S}_b, \frac{\tilde{f}'_b}{\tilde{f}_b})$$

where \hat{S}_b is the PLE of S_o based on $(M_i - bz_i, \delta_i)$'s, and \tilde{f}_b is a kernel estimator

$$\tilde{f}_b(t) = \frac{1}{h} \int K(\frac{x-t}{h}) d\hat{F}_b(x), \text{ with } K \ge 0, \int K(x) dx = 1,$$

the set $\{x : K'(x) \ne 0\}$ is not a null set, $n^{-1/2}/h \rightarrow 0$

e.q., $h = cn^{-1/5}$ and c is a predetermined constant. (3.2)

Examples of such kernels are

$$K(x) = \frac{3}{4}(1 - x^2)\mathbf{1}_{(|x| \le 1)},$$

and

$$K(x) = (1 - |x|)\mathbf{1}_{(|x| \le 1)}.$$

Other examples can be found in Härdle W. (1990, p.45). It is obvious that an Mestimate 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 $\hat{\Sigma} =$

$$\left(\sum_{i=1}^{n} \left(\delta_{i}[(\frac{\hat{f}'}{\hat{f}})(M_{i}-bz_{i})]^{2}-(1-\delta_{i})(\frac{\hat{f}}{\hat{S}})^{2}(M_{i}-bz_{i})\right)(z_{i}-\overline{z})(z_{i}-\overline{z})'\right)^{-1},$$

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 two numbers b_1 and b_2 such that $\Phi(b_1)\Phi(b_2) < 0$. WLOG, assume $b_1 > 0$ 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 $i \ge 2$ 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, $\Phi = -2(\sum_{i=1}^{n} (z_i - \overline{z})(X_i - \overline{X}) - b\sum_{i=1}^{n} (z_i - \overline{z})^2).$

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 ${\bf b}$ 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 $X_i = \beta \mathbf{z}_i + W_i$, i = 1, ..., n, where $W_i = \alpha + \epsilon_i$. 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

$$\mathbf{L} = \prod_{i=1}^{n} \mu_F (I_i - \mathbf{b} \mathbf{z}_i), \tag{4.1}$$

where F is a cdf, I_i 's are the observed intervals containing X_i , and $I_i - c$ is a shift of I_i by c units. The SMLE of (β, F_o) is

$$(\hat{\beta}, \hat{F}) = \operatorname{argmax}_{\mathbf{b}, F} \{ \mathbf{L}(\mathbf{b}, F) : \mathbf{b} \in \mathcal{R}^p, F(t) \uparrow, F(-\infty) = 0, F(\infty) = 1 \}.$$

Remark. With complete data, the LSE minimizes

$$\sum_{i=1}^{n} (X_i - \alpha - \mathbf{b}z_i)^2 \text{ or } \int t^2 d\hat{F}_{\epsilon}(t).$$

The SMLE of (α, β, F_o) maximizes

$$\mathcal{L}(a, \mathbf{b}, F) = \prod_{i=1}^{n} f_{\epsilon} (X_i - a - \mathbf{b} z_i).$$

Example 1. Consider a simple example of RC data, say,

3 (M_i, δ_i, z_i) s are (1, 1, 1), (2, 0, 1), (3, 1, 0).

Let
$$T_i(\beta) = M_i - \beta z_i$$
 $(X_i - \beta Z_i \sim F_o(t) = F_{X|Z}(t|0).$
 I_i 's are $\{X_1\}, (X_2, \infty), \{X_3\}$
 $I_i - bz_i$'s are $\{T_1(b)\}, (T_2(b), \infty), \{T_3(b)\}.$
 $L(b, F) = \prod_{i=1}^n \mu_F(I_i - bz_i) = f(1-b)S(2-b)f(3)$ Why ?
 $L(b, F) \leq \hat{f}_b(1-b)\hat{S}_b(2-b)\hat{f}_b(3)$

where $\hat{S}_b(t)$ is the PLE based on $T_1(b)$, $T_2(b)$, $T_3(b)$, and $\hat{f}_b(t) = ?$ The PLE \hat{S}_b depends on the ranks of $T_i(b) = M_i - bz_i$'s, Let (r_1, r_2, r_3) be the ranks of $(T_1 \quad T_2 + \quad T_3)(b)$ $(1-b \quad (2-b) + \quad 3)$ $T_i(b)$'s change their ranks after their ties: $M_k - bz_k = M_j - bz_j$ with $z_k \neq z_j$. $1-b=3 \ (T_1(b)=T_3(b))$ and $2-b=3 \ (T_2(b)=T_3(b))$.

Their solutions are $b_1 = -2$ and $b_2 = -1$.

-2 and -1 partition $(-\infty, \infty)$ into 5 disjoint intervals. r_1, r_2, r_3 remain constant in each interval.

Each point in $\{-2\} \cup (-1, \infty)$ is an SMLE of β . For examples

One SMLE of (β, F_o) is $(\hat{\beta}, \hat{F}_{\hat{\beta}}(t)) = (-2, \frac{2}{3}\mathbf{1}(t \ge 3)).$ How about $(-2, \frac{2}{3}\mathbf{1}(t \ge 3)) + \frac{1}{3}\mathbf{1}(t \ge 5)?$ $(-2, \frac{2}{3}\mathbf{1}(t \ge 3)) + \frac{1}{3}\mathbf{1}(t \ge 4)?$

Example 2. Consider the simple linear regression (p = 1), with complete data, say $X = \beta Z + W$, where W and $Z \sim Bin(1,0.5)$ 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, ..., 4. Thus there are 4 possible values of T(b), say, $T_i(b) = X_i - bZ_i$: 0, 1, 1 - b, 2 - b. Suppose a random sample of size n contains N_1 , N_2 , N_3 and N_4 of them. One may consider

the parametric approach (say the MLE or the MME of (p, β) , assuming $W \sim Bin(1, p)$; the semi-parametric approach (the LSE of β , assuming F_W is unknown);

the SMLE approach (assuming F_W is unknown) as follows. The empirical df

$$\hat{f}_b(T_i(b)) = \begin{cases} \frac{N_i}{n} & \text{if } T_i(b) \neq T_j(b) \ \forall \ j \neq i, \\ \frac{N_i + N_j}{n} & \text{if } T_i(b) = T_j(b) \text{ for only one } j \neq i \text{ where } i, j = 1, 2, 3, 4. \end{cases}$$
(4.2)

The possible solutions b to the equations $T_i(b) = T_j(b)$ are 0, 1, 1 and 2.

They partition $(-\infty, \infty)$ into 7 intervals.

$$\begin{pmatrix} b: & 0 & 1 & 2 & OW \\ (T_1, ..., T_4): & (0, 1, 1-b, 2-b) = (0, 1, 1, 2) & 0, 1, 0, 1 & 0, 1, -1, 0 & no \ ties \\ \hat{f}_{T(b)}: & (\frac{N_1}{n})^{N_1}, (\frac{N_2+N_3}{n})^{N_2}, (\frac{N_2+N_3}{n})^{N_3}, (\frac{N_4}{n})^{N_4} \end{pmatrix}$$

For *n* large enough, $N_i \approx n/4$, and likelihood function (1.2) or (4.1) $\mathcal{L} = \prod_{i=1}^n \hat{f}_b(T_i(b)) =$

$$\begin{cases} \left(\frac{N_1+N_3}{n}\right)^{N_1+N_3} \left(\frac{N_2+N_4}{n}\right)^{N_2+N_4} \approx (0.5)^n & \text{if } b = 1, \\ \left(\frac{N_1}{n}\right)^{N_1} \left(\frac{N_2+N_3}{n}\right)^{N_2+N_3} \left(\frac{N_4}{n}\right)^{N_4} \approx (0.5)^{n/2} (0.25)^{n/2} & \text{if } b = 0, \\ \left(\frac{N_1+N_4}{n}\right)^{N_1+N_4} \left(\frac{N_2}{n}\right)^{N_2} \left(\frac{N_3}{n}\right)^{N_3} \approx (0.5)^{n/2} (0.25)^{n/2} & \text{if } b = 2, \\ \left(\frac{N_1}{n}\right)^{N_1} \left(\frac{N_2}{n}\right)^{N_2} \left(\frac{N_3}{n}\right)^{N_3} \left(\frac{N_4}{n}\right)^{N_4} \approx (0.25)^n & \text{otherwise,} \end{cases}$$
(4.3)

is maximized by b = 1. Thus the SMLE of β is $\hat{\beta} = 1$ for all large n. consistent ? efficient ? The SMLE of F_o is $\hat{F}(t) = \hat{p}\mathbf{1}(t \ge 0) + (1 - \hat{p})\mathbf{1}(t \ge 1)$, where $\hat{p} = (N_1 + N_3)/n$. The LSE $\tilde{\beta}_{LSE} = \beta + \frac{N_1N_4 - N_2N_3}{(N_1 + N_4)(N_2 + N_3)}$. Thus,

 $P(\hat{\beta}_{LSE} \neq \beta \text{ i.o.}) = 1.$ The LSE $\hat{\beta}_{LSE}$ satisfies $\sqrt{n}(\tilde{\beta}_{LSE} - \beta) \xrightarrow{D} N(0, 1)$, and $n\sigma_{\tilde{\beta}_{LSE}}^2 \to 1$ as $n \to \infty$, as the asymptotic variance $\sigma_{\tilde{\beta}_{LSE}}^2 = \frac{1}{n}$.

It can be shown that $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{D} N(0,0)$ and $n\sigma_{\hat{\beta}}^2 \to 0$.

The MLE does not have this property.

Remark: If p > 1, the solution to SMLE is little bit more complicated, but doable. Lemma 5.4.1. Suppose that $X = \beta' \mathbf{Z} + W$, where W and \mathbf{Z} are uncorrelated, $E(W) = \mu$ and $\sigma_W^2 = \sigma$, $\beta = (\beta_1, ..., \beta_p)'$, $\mathbf{Z} = (Z_1, ..., Z_p)'$ and $Cov(\mathbf{Z})$ is diagonal. Then $X = \beta_i Z_i + W^*$, where $W^* = W + \sum_{j \neq i} \beta_j Z_j$. $E(W^*) = \mu_*$ and $\sigma_{W^*}^2 = \sigma_2^2$.

In other words, under the assumptions in Lemma 5.4.1, the estimator of β can be estimated coordinate by coordinate. The assumption that Cov(Z) is diagonal is a key. If Cov(Z) is not diagonal, find a matrix B such that $B\Sigma(Z)B'$ is diagonal. In application, let $\hat{C}ov(Z) = \sum_{i=1}^{n} Z_i Z'_i - \overline{Z}'\overline{Z}$, and find a matrix B such that $B\hat{C}ov(Z)B'$ is diagonal. Then let $X^* = \beta Z + W^*$, where $X^* = BX$ and $Z^* = BZ$.

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 . The maximum of L is $L(b, \hat{F}_b)$.

To find the SMLE of (β, F_o) , it suffices to maximizes $L(b, F_b)$.

Let $a_1 < \cdots < a_m$ be all the solutions to equations of form

 $M_k - bz_k = M_j - bz_j, z_k \neq z_j$, where $M_i = L_i$ or $R_i, a_0 = -\infty$ 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, $L(b, F_b)$ is constant in b on the interval (a_i, a_{i+1}) . There are at most 2m + 1 different values of $L(b, \hat{F}_b)$'s.

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 $\mathcal{L}(b, \hat{F}_b)$ over all $b \in \mathcal{A}$ is an SMLE of β . We summarize as an algorithm:

- 1. Derive \mathcal{B} and \mathcal{A} , where $\mathcal{B} = \{b: b = \frac{M_i M_j}{z_i z_j}, z_i \neq z_j, (\delta_i, \delta_j) \neq (0, 0), 1 \le i < j \le n\}.$ Let $a_1 < \cdots < a_m$ be the distinct elements of \mathcal{B} .
- 2. Derive $\mathcal{L}(b) = \mathcal{L}(b, \hat{F}_b)$ 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 close 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 $\beta = ?$ Is β identifiable ? **Ans.** The SMLE of β is the $b \in \mathcal{B}$ that $T_{i_1}(b) = \cdots = T_{i_m}(b)$, where *m* is the largest.

In particular, if W is continuous, then each $b \in \mathcal{B}$ is an SMLE of β , as the edf gives equal weight to each $T_i(b)$ if $b \notin \mathcal{B}$; otherwise, the edf gives weight 2/n to $T_i(b)$ (= $T_j(b)$) and 1/n to the rest $T_k(b)$'s.

Properties of the SMLE with RC data under certain regularity conditions:

- 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 and there exists two distinct values of \mathbf{Z} , say \mathbf{z}_1 and \mathbf{z}_2 such that $W + \beta \mathbf{X} \leq \tau_Y$ 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 $W \sim N(\mu, \sigma^2)$, 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, $P\{\delta = 1\} \in (0, 1)$ and under certain regularity conditions, then the SMLE is consistent.

An inconsistent SMLE example. Let $X = \beta Z + W$, $\beta = 1$, $Y \equiv 1$ and $P(Z = 0) = P(Z = 0.6) = 0.5 = P(W \in (0, 0.1)) = P(W = 0.5)$ $\begin{pmatrix} type \ Z \ W \ M \ \delta \ T \ \# of \ T'_i s \\ 1 \ 0 \ (0, 0.1) \ (0, 0.1) \ 1 \ (0, 0.1) \ \approx n/4 \\ 2 \ 0 \ 0.5 \ 0.5 \ 1 \ 0.5 \ \approx n/4 \\ 3 \ 0.6 \ (0, 0.1) \ (0.6, 0.7) \ 1 \ (0.6 - 0.6b, 0.7 - 0.6b) \ \approx n/4 \\ 4 \ 0.6 \ 0.5 \ 1 \ 0 \ 1 - 0.6b \ \approx n/4 \end{pmatrix}$ if $n \approx \infty$. $\mathcal{B} \subset (1/6, 2/6) \cup \{5/6\} \cup (1 - 1/6, 1 + 1/6) \cup (9/6, 10/6)$,

$$L \approx \begin{cases} \frac{1}{n^{k}} \frac{1}{n^{k}} \frac{1}{4^{k}} \frac{1}{4^{k}} & \text{if } b \in (1/6, 2/6) \\ \frac{1}{n^{k}} \frac{(k+1)^{k+1}}{n^{k+1}} \frac{1}{n^{k-1}} \frac{1}{4^{k}} & \text{if } b = 2/6 \\ \frac{1}{n^{n}} & \text{if } b = 5/6 & (1) \\ \frac{1}{n^{k-1}} \frac{1}{n^{k-1}} \frac{2^{2}}{n^{2}} \frac{1}{2^{2k}} = \frac{1}{(4k)^{2k-2}} \frac{4}{(4k)^{2}} \frac{1}{2^{2k}} = \frac{4}{(4k)^{2k}} \frac{1}{2^{2k}} = \frac{1}{(4k)^{2k}} \frac{4}{2^{2k}} & \text{if } b \in (5/6, 7/6) \\ \frac{1}{(4k)^{k}} (3/4)^{k} (\frac{3}{4*2k})^{k} (3/8)^{k} = \frac{1}{(4k)^{2k}} (3\sqrt{3}/8)^{2k} = largest & \text{if } b \ge 10/6 \end{cases}$$

 $3^k (3/8)^{2k} = 27^k / 8^{2k} = (\sqrt{27}/8)^{2k} \ 3^k (3/8) * 2 = \sqrt{27/16} > 1.$

5. If F_o is continuous, $P\{\delta = 1\} \in (0, 1)$ and under certain regularity conditions, then the SMLE is asymptotically normal, with estimated asymptotic covariance matrix

$$\hat{\Sigma} = \left(\sum_{i=1}^{n} (1-\delta_i) \left(\frac{\tilde{f}(T_i(\hat{\beta}))}{\hat{S}_{\hat{\beta}}(T_i(\hat{\beta}))}\right)^2 z_i z_i'\right)^{-1},\tag{4.1}$$

where \tilde{f} is a kernel estimate of the df of F_o . Note that

$$\ln(f^{\delta}(T(b))S^{1-\delta}(T(b))) = \delta \ln f(T(b)) + (1-\delta)\ln S(T(b))$$
$$\frac{\partial \ln(f^{\delta}(T(b))S^{1-\delta}(T(b)))}{\partial b} = \left[-\delta \frac{f'_o(T(b))}{f_o(T(b))} + (1-\delta)\frac{f_o(T(b))}{S_o(T(b))}\right]Z$$

Thus the conjecture says that the first term is missing in Σ .

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

$$E[(\delta \frac{f'(T(b))}{f(T(b))} + (1-\delta) \frac{f(T(b))}{S(T(b))})^2 Z Z']$$

We now verify property 1.

Assume that p = 1, W and \mathbf{z} are both continuous independent random variables. Let $(X_i, z_i), i = 1, ..., n$ are observations. Then $\hat{\beta} = \frac{X_1 - X_2}{z_1 - z_2}$ Why ??

 $\begin{array}{l} X_i,\, z_i,\, \hat{\beta} \text{ are all continuous random variables.} \\ \hat{\beta} = \beta + \frac{W_1 - W_2}{z_1 - z_2}. \\ \text{Let } \hat{\beta}_n = \max\{b:\ b \in \mathcal{B}\},\, \text{then } \hat{\beta}_n \to \infty \text{ a.s., as } \min\{|Z_1 - Z_2|\} \to 0 \text{ a.s..} \\ \text{ That is, } \hat{\beta}_n \text{ is an inconsistent SMLE.} \\ \text{On the other hand, the SMLE that is closest to the LSE is consistent.} \end{array}$

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, W and Y are independent. We consider several cases in our simulation studies:

- (1) F_o is continuous (Examples 5, 6 and 7), or discrete (Example 3), or discontinuous but not discrete (Examples 4).
- (2) All the underlying distributions belong to the exponential family (Examples 5 and 7) or F_o does not belong to the exponential family (other examples). In the following examples, let $X = \beta Z + W$ and $E(W) = \alpha$.

Example 3. Suppose *W* 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. Suppose W 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)							
		(lpha,eta)	SMLE $(\hat{\alpha}, \tilde{\beta})$	BJE			
Example 3 (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)			
Exa	ample 4 (disco	ntinuous 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)			
E	xample 5 (con	tinuous 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)			
E	xample 6 (con	tinuous 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)			
E	Example 7 (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

 $f(x) = \frac{1}{\sigma} e^{-\left[\frac{x-\mu}{\sigma}+1\right]} \mathbf{1}_{(x>\mu-\sigma)}.$

Q: Does it belong to the exponential family ?

Example 5. Suppose W, Y and Z have distributions Exp(5,2), Exp(3,4) and Exp(2,2), respectively. $(\alpha, \beta) = (5, 1)$.

Example 6. Suppose $W \sim U(-1,1)$, $Z \sim Exp(0,19)$, $Y \sim Exp(0,25)$, $(\alpha,\beta) = (0,2)$. **Example 7.** Suppose $W \sim N(0,1)$, $Y \sim N(0,6)$, $Z \sim N(2,1)$, $(\alpha,\beta) = (0,1)$.

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 verse 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 8. (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)			

millions of dollars 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 W is continuous and there is no censoring, then

$$\hat{f}_{\mathbf{b}}(T_i(\mathbf{b})) = 1/n \ \forall \ i, \text{ except perhaps for two, at which } \hat{f}_{\mathbf{b}}(T_i(\mathbf{b})) = 2/n.$$
 (2.2)

Consequently,

$$\mathbf{L}(\hat{S}_{\hat{\beta}},\hat{\beta}) = 2^2/n^n \text{ and } \mathbf{L}(\hat{S}_{\beta},\mathbf{b}) = 1/n^n \text{ if } \mathbf{b} \text{ is not an SMLE}$$
$$|\frac{1}{n}\mathrm{ln}\mathbf{L}(\hat{S}_{\hat{\beta}},\hat{\beta}) - \frac{1}{n}\mathrm{ln}\mathbf{L}(\hat{S}_{\mathbf{b}},\mathbf{b})| = \begin{cases} \frac{1}{n}\mathrm{ln}4 & \text{if } \mathbf{b} \text{ is not an SMLE}\\ 0 & \text{otherwise} \end{cases}$$
(2.3)

In other words, if W is continuous, the inconsistent SMLE $\hat{\beta}$ also satisfies that

$$|\frac{1}{n} \mathrm{lnL}(\hat{S}_{\hat{\beta}}, \hat{\beta}) - \frac{1}{n} \mathrm{lnL}(\hat{S}_{\beta}, \beta)| \to 0$$

In a consistency proof, one may want to establish

$$\lim_{n \to \infty} \left[\frac{1}{n} \ln \mathcal{L}(\hat{S}_{\hat{\beta}}, \hat{\beta}) - E\{ \frac{1}{n} \ln \mathcal{L}(S_o, \beta) \} \right] = 0 \ a.s..$$

However, it does not work here, because $E\{\frac{1}{n}\ln L(S_o,\beta)\} = -\infty$ 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 **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 W, Y and Z (= **Z**) have distributions Exp(3, 1), Exp(1, 1) and Exp(0, 1), respectively. $(\alpha, \beta) = (3, 1).$

n		$\beta (\sigma_{\hat{eta}})$	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	

Table 3. Simulation Results on estimating (α, β)

Remark. Under the semiparametric model, $X = \alpha + \beta' \mathbf{Z} + \epsilon = \beta' \mathbf{Z} + T$, the location parameter $\alpha = E(T)$ is not identifiable under censoring. This is the major reason why people do not consider the model $X = \alpha + \beta' \mathbf{Z} + \epsilon$, where $E(\epsilon) = 0$. 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).
 - a. Find the SMLE of β under the linear regression model.
 - b. Find the PSMLE of β under the linear regression model.
- 2. Suppose n = 1000 right-censored observations are from the model with $Y \equiv 1$ and $X = \beta Z + W$, where $W \perp Z$ and $\beta = 1$. W and $Z \sim Bin(1, p), p \in (0, 1)$.

a. Derive the MME of (β, p) and S_W , assuming β is unknown (as no closed form

solution for the MLE).

- b. Derive the SMLE of β and S_W (without the Bin(1,p) assumption).
- c. Derive the NPMLE of S_M based on M_i 's only.
- d. Derive the NPMLE of S_X based on (M_i, δ_i) 's.
- e. What are S_W , S_M and S_X ?
- d. What are the limits of the above estimators (you can assume p = 0.5)?

$\S5.4.5.$ A modified SMLE (MSMLE) 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 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 L, one can only find a critical point of L, or a root to $\frac{\partial \ln L}{\partial b}$. 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

$$f_F(x) = \int_{-\infty}^{\infty} \frac{1}{h} K(\frac{x-t}{h}) dF(t), \ K(x) = \frac{1}{2} \mathbf{1}_{(-1 < x \le 1)}, \ h > 0, \ \lim_{n \to \infty} h = 0$$
(5.1)

(e.g., $h = O(n^{-1/5})$, 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 $T_i(\mathbf{b}) = M_i - \mathbf{b}\mathbf{z}_i, i = 1, ..., n$. Then

$$\mathcal{L}(\hat{S}_{\mathbf{b}}, \mathbf{b}) = \prod_{i=1}^{n} \left[\left(\frac{[\hat{S}_{\mathbf{b}}(T_{i}(\mathbf{b}) - h) - \hat{S}_{\mathbf{b}}(T_{i}(\mathbf{b}) + h)]}{2h} \right)^{\delta_{i}} (\hat{S}_{\mathbf{b}}(T_{i}(\mathbf{b})))^{1 - \delta_{i}} \right].$$
(5.2)

Then estimate β by $\hat{\beta}$ which maximizes $\mathcal{L}(\hat{S}_{\mathbf{b}}, \mathbf{b}, \mathbf{b}, \mathbf{b})$ over $\mathbf{b} \in \mathcal{R}^p$. $\hat{\beta}$ is called the MSMLE of β . Then $\hat{S}_{\hat{\beta}}(t)$ is an MSMLE of $S_o(t)$, where $T \sim F_o$. Let $\alpha = E(T)$ if E(T) 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

$$\hat{\alpha} = a(\hat{\beta}), \text{ where } a(\hat{\beta}) = \frac{\int_{x \le T_{(n)}(\hat{\beta})} x d\hat{F}_{\hat{\beta}}(x)}{\int_{x \le T_{(n)}(\hat{\beta})} 1 d\hat{F}_{\hat{\beta}}(x)},$$
(5.3)

where $T_{(1)} \leq \cdots \leq T_{(n)}$ 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 $\frac{1}{2h}$ in (5.2) does not depend on **b**, it suffices to maximize for $\mathbf{b} \in \mathcal{R}^p$,

$$l(\mathbf{b}) = \prod_{i=1}^{n} \left[[\hat{S}_{\mathbf{b}}(T_i(\mathbf{b}) - h) - \hat{S}_{\mathbf{b}}(T_i(\mathbf{b}) + h)]^{\delta_i} (\hat{S}_{\mathbf{b}}(T_i(\mathbf{b})))^{1 - \delta_i} \right].$$
(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 $a_1 < \cdots < a_m$ are all the distinct values of

$$b = \frac{M_i - M_j + kh}{z_i - z_j}, \ z_i \neq z_j, \ k = 0, \pm 1, \pm 2, \ \text{and} \ \begin{cases} \delta_i = \delta_n = 1 & \text{if } k \in \{0, \pm 2\}, \\ \delta_j > \delta_i & \text{if } k = \pm 1. \end{cases}$$
(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, ..., 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. If F_o is discontinuous and there exists two distinct values of \mathbf{Z} , say \mathbf{z}_1 and \mathbf{z}_2 such that
- $W + \beta \mathbf{z}_i \leq \tau_Y$, then $P\{\beta \neq \beta \text{ i.o.}\} = 0$.
 - Moreover, it is <u>conjectured</u> that * If (M, δ, π) 's <u>conside</u> from the stand
 - * If $(M_i, \delta_i, \mathbf{z}_i)$'s are i.i.d. from the standard RC model, and β is identifiable then the MSMLE $\hat{\beta}$ is consistent without any additional assumptions (under preparation). Simulation suggests that it is also efficient if the parametric MLE is.

In particular, it suggests that

if F_o is continuous, then the limit of $n\Sigma_{\hat{\beta}}$ attains the efficient lower bound of the variance, *i.e.*, $-n\left(\frac{\partial \log L}{\partial \beta \partial \beta'}\right)^{-1}$.

The following are simulation results supporting the above conjectures. In our simulation, we assume that T, \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 $T \sim U(-1, 1)$ (the uniform distribution), $Z \sim U(0, 9)$ and $(\alpha, \beta) = (0, 2)$.

Example 5.2. Suppose T 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 $T \sim N(0, 0.09), Z \sim U(0, 9)$ 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 $T \sim N(\mu, \sigma^2)$, provided $n \geq 200$. In the next 4 cases, there are right-censored data. Define $Y^c = Y - \beta' \mathbf{Z}$ and $\tau = \sup\{t : P(Y^c < t) < 1\}$.

Table 1. Simulation Results on estimating β without censoring.

	SMLE	parameter β	LSE	$\text{MSMLE } \hat{\beta}$		
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_c	<u>,</u>)			
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)		
Exa	mple 5.5 $(N(\mu, \sigma^2$	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)		
Exan	ple 5.6 (continuo	ous).		$F_o(\tau -) = 1.$		
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 (Exp(μ, σ)). $F_o(\tau -) = 1.$						
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)		

ation Results on estimating β with con

Example 5.4. Suppose T is a mixture of U(0, 0.5) and 51, w.p. 0.5 and 0.5, respectively, $Z \sim U(1,2), Y \sim U(4,4.1), \text{ and } (\alpha,\beta) = (25.625,1).$

Example 5.5. Suppose $T \sim N(0, 0.09), Z \sim U(0, 9)$ and Y equals 0.5 and 39 w.p. 0.5 and 0.5, respectively. $(\alpha, \beta) = (0, 1)$.

Example 5.6. Suppose $T \sim U(-1,1)$, $Z \sim Exp(0,19)$, $Y \sim Exp(0,25)$, $(\alpha,\beta) = (0,2)$. **Example 5.7.** Suppose T, Y and Z have distributions Exp(5,2), Exp(3,4) and 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 β_{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 = 0 in 5.4), and consider the case of complete data in Example 5.2. Then 5.4) 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 $T_i = T_i(\beta)$. Let $n_1 = \sum_{i=1}^{n} \mathbf{1}_{(T_i=0.45)}$. If n is large enough, one expects that $n_1 \ge 10$. 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 $b \ne \beta$, 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 $T \sim N(\mu, \sigma^2)$. 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 $Var(T)/(n \cdot Var(Z)) = 2.5^2/n$ (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		$\sqrt{n \cdot ELB}$
$\sqrt{n}\hat{\sigma}_{\hat{eta}}$:	(10.363)	5.707	3.875	3.132	2.503)	$\begin{pmatrix} 2.5 \end{pmatrix}$

(5) It is well known that if $T \sim N(\mu, \sigma^2)$ 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 $\frac{SE_{\hat{\beta}_{BJE}}^2}{SE_{\hat{\beta}_{BJE}}^2} = 1$ 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.

Example:	5.1	5.2	5.3	5.4	5.5	5.6	5.7
$rac{\hat{\sigma}^2_{ ilde{eta}_{BJE}}}{\hat{\sigma}^2_{\hat{eta}}}$:	(1.3	∞	1.0	3.9	?	1.8	3.7
F_o :	\bigcup unif.	mixture	normal	mixture	normal	unif.	expon.)

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 $h = n^{-1/5}$ yields $\hat{\beta} = 0$, which does not lead to a satisfactory estimate. We choose $h = 3n^{-1/5}$.

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	-0.028(0.015)	-0.036(0.017)	$0.058\ (0.023)$				
	$0.004 \ (0.017)$							
	$0.002\ (0.016)$							

Table 5. Regression analysis on the Stanford heart transplant data

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 $h = n^{-1/5}$).
- 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 Dong, J.Y. (2019). Identifiability Conditions For The Linear Regression Model Under Right Censoring *Comm. Statist. A—Theory and Methods* https://www.tandfonline.com/doi/full/10.1080/03610926.2020.1743315
- * 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 semiparametric 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. (2007). Asymptotic distribution of the Buckley-James estimator under non-standard conditions *Statistica Sinica*. 17 341-360
- * Yu, Q.Q and Wong, G.Y.C. (2006) A data set that there are two BJEs. 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
 - $H_0: F = F(\cdot; \theta)$, where the form $F(\cdot|\cdot)$ is given.
- 3. $H_0: F = F_o$, where F_o is a given cdf.
- 4. $H_0: \theta = \theta_o$, where θ is a parameter of F_X , such as the mean and the variance, or the parameter in a certain parametric distribution.

In the elementary statistics course, 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 $H_0: \hat{\beta} = \beta_o \text{ vs. } H_1: \hat{\beta} \neq \beta_o$ is $\phi = \mathbf{1}((\hat{\beta} - \beta_o)'J^{-1}(\hat{\beta} - \beta_o)|_{\beta=\beta_o} > \chi^2_{\alpha,p})$, where $J = -(\frac{\partial^2 \ln L}{\partial \beta \partial \beta'})^{-1}$, $\beta \in \mathcal{R}^p$ and $\chi^2_{\alpha,p}$ is the $(1 - \alpha)$ 100-th percentile of the χ^2 distribution with degree freedom p.

(Or
$$\phi_1 = \mathbf{1}(|\beta - \beta_o|/\hat{\sigma}_{\hat{\beta}} > z_{\alpha/2})).$$

(2) Score test approach: Under certain assumptions, a test for H_0 : $\beta = \beta_o$ is $\phi = \mathbf{1}(U(\beta)'JU(\beta)|_{\beta=\beta_o} > \chi^2_{\alpha,p})$, where $U(\beta) = \frac{\partial \ln \underline{L}}{\partial \beta}$.

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 $X = \beta Z + \theta Z^2 + W$.

\S **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_0 : $S_o(t_0) = p_0$, where p_0 is known: For RC data, under the RC model, $\phi = \mathbf{1}_{(\frac{|\hat{S}_{pl}(t_0) - p_0|}{\hat{\sigma}_{\hat{S}_{nl}(t_0)}} \ge z_{\alpha/2})}$,

where $\hat{\sigma}^2_{\hat{S}_{pl}(t_0)}$ 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

$$\phi = \mathbf{1}_{(\frac{|\hat{S}(t_0) - p_0|}{\hat{\sigma}_{\hat{S}(t_0)}} \ge z_{\alpha/2})},$$

where $\Phi(z_{\alpha}) = 1 - \alpha$, Φ is the cdf of N(0, 1), and $\hat{\sigma}^2_{\hat{S}(t_0)}$ is the estimate of $\sigma^2_{\hat{S}(t)}$ 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 $\S4.7$.

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

$$\phi = \mathbf{1}_{(|T| > z_{\alpha/2})}$$
, where $T = \frac{\overline{X} - \mu}{s/\sqrt{n}}$, and $s^2 = \frac{\sum_{i=1}^n (X_i - \overline{X})^2}{n-1}$

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

$T = \frac{U}{\hat{\sigma}_U}, \text{ where } U = \int_0^\infty W(x)(\hat{S}_{pl}(x) - S_o(x))dx, \qquad (1.1)$

and W is a weight function, which may depend on F_o and S_Y . For example, let W(t) be an estimate of $P(X \land Y \ge t)$ 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 $b_1 < \cdots < b_{m+1}$. 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

$$U = b_0 + (b_2 - b_1, ..., b_{m+1} - b_m)(\hat{S}(b_1), ..., \hat{S}(b_m))^t + \int_0^\tau S_o(t)dt.$$
$$\sigma_U^2 = (b_2 - b_1, ..., b_{m+1} - b_m)\Sigma(b_2 - b_1, ..., b_{m+1} - b_m)^t$$

where Σ is the covariance matrix of $(\hat{S}(b_1), ..., \hat{S}(b_m))$.

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 (Koriol 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, \ldots, I_n ,
- (1) resample I_1^* , ..., I_n^* , and get U_1 ;
- (2) repeat step 1 N times (including the first time), and get $U_1, ..., U_N$;
- (3) $U_1, ..., U_N$ 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} .

\S 6.2. Two-sample problem

Suppose that there are two independent random samples with sizes n_1 and n_2 , from survival function S_1 and S_2 , respectively. $H_0: S_1 = S_2$, v.s. $H_1: S_1(t) \ge S_2(t)$ for all t and $S_1 \ne S_2$.

(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 $\sigma_U^2 = Var(\int_0^{\tau} \hat{S}_1(t)dt) + Var(\int_0^{\tau} \hat{S}_2(t)dt)$, 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 $(L_i, R_i), i = 1, ..., n_1 + n_2$, 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 $100(1 - \alpha)$ 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 $(x \land y, \delta)$ 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.

$$U_{ij} = \begin{cases} 1 & \text{if } a_i > b_j \text{ or } a_i + \ge b_j \\ & (\text{we know for sure that obs-}i \text{ in sample } 1 > \text{obs-}j \text{ in sample } 2), \\ -1 & \text{if } a_i < b_j \text{ or } a_i \le b_j + \\ & (\text{we know for sure that obs-}i \text{ in sample } 1 < \text{obs-}j \text{ in sample } 2), \\ 0 & \text{if not sure.} \end{cases}$$

 $U = \sum_{j=1}^{n_2} \sum_{i=1}^{n_1} U_{ij}. \text{ A test is } \phi = \mathbf{1}_{\left(\frac{U}{\hat{\sigma}_U} > z_\alpha\right)}.$ 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 $c_i +, i = 1, ..., n_1 + n_2$. 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

 $V_{kh} = \begin{cases} 1 & \text{if we know for sure obs-} k > \text{obs-} h, \\ -1 & \text{if we know for sure obs-} k < \text{obs-} h, \\ 0 & \text{if not sure,} \end{cases}$ (2.1)

(for RC data, we know for sure obs-k > obs-h iff $c_k > c_h$ or $c_k + \ge c_h$.)

$$V_k = \sum_{h=1}^{n_1+n_2} V_{kh},$$
(2.2)

Let W be a random variable taking values $\sum_{i=1}^{n_1} V_{k_i}$, where $\{k_1, ..., k_{n_1}\}$ is a selection of n_1 distinct integers from $\{1, ..., n_1 + n_2\}$.

$$V = \sum_{k=1}^{n_1} V_k$$
 (2.3)

is a value of W. Under this sample space (permutation sample space, each permutation has equal probability),

$$E(W) = 0 \text{ and } \sigma_W^2 = n_1 n_2 \sum_{k=1}^{n_1+n_2} \frac{V_k^2}{(n_1+n_2)(n_1+n_2-1)}.$$

$$\left(E(W^2) = \sum_{k_1,\dots,k_{n_1}} \frac{1}{\binom{n_1+n_2}{n_i}} (\sum_{i=1}^{n_1} V_{k_i})^2\right).$$
s that

Mantel suggests that

$$\frac{W}{\sigma_W} \xrightarrow{D} N(0,1) \text{ as } n \to \infty.$$

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.

$$V = \sum_{k=1}^{n_1} \sum_{h=1}^{n_1+n_2} V_{kh}$$

= $\sum_{k=1}^{n_1} \sum_{h=1}^{n_1} V_{kh} + \sum_{k=1}^{n_1} \sum_{h>n_1}^{n_1+n_2} V_{kh}$
= $\sum_{k=1}^{n_1} \sum_{h=1}^{n_1} V_{kh} + \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} U_{ij}$
= $0 + U$ (as $V_{kh} = -V_{hk}$).

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

$$\psi = \mathbf{1}_{\left(\frac{U}{\sigma_W} > z_\alpha\right)}.\tag{2.5}$$

For IC data we are sure that observation-k > observation-h iff

either
$$L_k > R_h$$
 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, ..., n. $\hat{S}_X(t) = \prod_{i: M_{(i)} \leq t} (1 - \frac{\delta_{(i)}}{n-i+1}) \rightarrow S_X(t)$ a.s.,

$$\hat{S}_M(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(M_i > t) \to S_M(t) \text{ a.s.}.$$

 $\hat{S}_M(t) = \hat{S}_X(t)\hat{S}_Y(t) \text{ and } P(M > t) = P(X > t)P(Y > t).$

 $\hat{S}_M(t)/\hat{S}_X(t) \to S_Y(t) \text{ a.s.},$ $\hat{S}_M(t)/\hat{S}_X(t) = \hat{S}_Y(t) ?$ $\hat{S}_Y(t) = \prod_{i: \ M_{(i)} \le t} (1 - \frac{1 - \delta_{(i)}}{n - i + 1}) \to S_Y(t) \text{ a.s.},$ Thus 1 $\delta = 1(Y \le Y) \neq 1(Y \le Y)$ does no

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:

 $c_1, ..., c_{m+n}$ are pooled RC observations (sample 1, sample 2).

 $\begin{array}{l}
\text{Among the pooled-sample} \begin{cases}
t_1 < \cdots < t_k & \text{are distinct exact times,} \\
d_{.,i} & \text{is the $\#$ of deaths at t_i, $i = 1, ..., k,} \\
\mathcal{R}(t_i) & \text{is the set of individuals in risk at time t_i.} \\
d_{j,i} = $\#$ of deaths from group j at time t_i.} \\
d_{j,.} = $\#$ of deaths in group j.} \\
r_{j,i} = $\#$ of elements of group j in $\mathcal{R}(t_i)$.} \\
r_{.,i} = r_{1,i} + r_{2,i}$.} \\
U = \sum_{j=1}^k [d_{1,j} - d_{.,j} \frac{r_{1,j}}{r_{.,j}}]. \\
\frac{U}{\hat{\sigma}_U} \xrightarrow{D} N(0,1), \text{ where } \hat{\sigma}_U^2 = \sum_{i=1}^k d_{.,i} \frac{r_{2,i}r_{1,i}(r_{.,i} - d_{.,i})}{r_{.,i}^2(r_{.,i} - 1)}
\end{array}$ (2.6)

The level- α logrank test is $\phi = \mathbf{1}_{(\frac{U}{\tilde{\sigma}_{II}} < -z_{\alpha})}$.

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. $\mathcal{D} = \{i : \text{ observation } i \text{ in pooled-sample died}\}.$

$$S_1 = S(t|0)$$
 and $S_2 = S(t|1)$, where $S(t|z) = (S_o(t))^{e^{\beta z_i(t)}}$, (2.7)

$$\begin{split} H_0: S_1 &= S_2 \text{ is equivalent to } \beta = 0.\\ lik &= \prod_{i \in \mathcal{D}} \frac{e^{\beta z_i(c_i)}}{\sum_{h \in \mathcal{R}(c_i)} e^{\beta z_h(c_i)}} \end{split}$$

$$U(\beta) = \frac{d \ln lik}{d\beta} = \sum_{i \in \mathcal{D}} [z_i(c_i) - \frac{\sum_{j \in \mathcal{R}(c_i)} e^{\beta z_j(c_i)} z_j(c_i)}{\sum_{h \in \mathcal{R}(c_i)} e^{\beta z_h(c_i)}}], --- \text{ the score function.}$$
$$U(0) = \sum_{i \in \mathcal{D}} [z_i(c_i) - \frac{\sum_{j \in \mathcal{R}(c_i)} z_j(c_i)}{|\mathcal{R}(c_i)|}] = \sum_{i \in \mathcal{D}} [z_i(c_i) - \sum_{i \in \mathcal{D}} \frac{\sum_{j \in \mathcal{R}(c_i)} z_j(c_i)}{|\mathcal{R}(c_i)|}]$$
$$U(0) = \sum_{j=1}^k w(t_j) [d_{1,j} - d_{\cdot,j} \frac{r_{1,j}}{r_{\cdot,j}}], \text{ as } z_i = 0 \text{ for } i > n_1.$$

If $w(t_i) = 1$, $U(0) = \sum_{i=1}^{n_1} [d_{1,j} - \frac{d_{\cdot,j}r_{1,j}}{r_{\cdot,j}}]$, corresponding to the logrank test.

The asymptotic variance of U(0) can be obtained by

$$\sigma_{U}^{2} \approx -\frac{dU(\beta)}{d\beta}\Big|_{\beta=0} = \sum_{i\in\mathcal{D}} \left[\frac{\sum_{j\in\mathcal{R}(c_{i})} e^{\beta z_{j}} z_{j}^{2}}{\sum_{h\in\mathcal{R}(c_{i})} e^{\beta z_{h}}} - \frac{\sum_{j\in\mathcal{R}(c_{i})} e^{\beta z_{j}} z_{j} \sum_{l\in\mathcal{R}(c_{i})} e^{\beta z_{l}} z_{l}}{(\sum_{h\in\mathcal{R}(c_{i})} e^{\beta z_{h}})^{2}}\right]\Big|_{\beta=0}$$
$$= \sum_{i=1}^{k} (w(t_{i}))^{2} d_{\cdot,i} \left[\frac{r_{1,i}}{r_{\cdot,i}} - \frac{r_{1,i}^{2}}{r_{\cdot,i}^{2}}\right]$$
$$= \sum_{i=1}^{k} [(w(t_{i}))^{2} d_{\cdot,i} \frac{r_{1,i}}{r_{\cdot,i}} \frac{r_{2,i}}{r_{\cdot,i}}]$$
(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.

$$\hat{\sigma}_{U}^{2} = -\frac{dU(\beta)}{d\beta}\Big|_{\beta=0} = \sum_{i=1}^{k} [d_{\cdot,i} \frac{r_{1,i}}{r_{\cdot,i}} \frac{r_{2,i}}{r_{\cdot,i}}] \text{ if } w(t_{i}) = 1,$$

*

**

which equals (2.6) as $d_{.i} = 1$ by the continuity assumption.

A two-sided test $\psi = \mathbf{1}_{(|\frac{U(0)}{\hat{\sigma}_U}| > z_{\alpha/2})}$ is called a linear rank test, where $w(t_i) = r_{\cdot,i}$ — 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.

$x = coxph(Surv(time) \sim ag + log(wbc), data = leuk)$ summary(x)

coefexp(coef)se(coef)Pr(>|z|) \boldsymbol{z} agpresent -1.06910.34330.4293-2.4900.01276 log(wbc)0.3677 1.4444 0.13602.7030.00687 exp(-coef)lower.95 exp(coef)upper.95agpresent 0.3433 2.91260.1480.7964log(wbc)1.4444 1.1060.69231.8857Concordance = 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.q. 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)). The logrank test needs the assumption that the data are from a PH model.

A test statistics for checking the PH model in the two-sample set-up 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 the pooled sample (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, $\hat{S}_1(t) \frac{1}{\sum_{i=1}^n \mathbf{1}(Z_i=1)} \sum_{j=1}^n \mathbf{1}(M_i > t, Z_i = 1).$

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. $S_{X,Z}$ is a joint cdf, which does not need to be from any PH model. S_X is the marginal distribution.

If $S_{X|Z}$ is from a PH model, then

 $S_X = E(S_{X|Z}(\cdot|Z)) = E((S_1(t))^{\beta Z}),$

which can be estimated by $\hat{S}_{X^*}(t) = \frac{1}{n} \sum_{i=1}^n (\hat{S}_1(t))^{\beta Z_i}$. Otherwise, $S_{X^*} = E((S_1(t))^{\beta Z})$ satisfies the PH model.

For IC data

Self and Grossman (1986) consider the linear regression problem with IC data for a given distribution with location-scale parameters and propose a marginal likelihood approach, the marginal distribution of the ranks of the underlying survival times. Note that their problem is essentially parametric even though they used a nonparametric approach. Rabinowitz et al. (1995) proposed a class of score statistics to estimate parameters of the accelerated failure time model with IC data. Finkelstein (1986, Biometrics) extended the logrank test to the IC data. Satten (1996) consider rank-based inference in the proportional hazards model with IC data. Both extended the logrank test for RC data.

Reference.

- * Self, S.G. and Grossman, E.A. (1986). Linear rank tests for interval-censored data with application to PCB levels in adipose tissue of transformer repair workers. *Biometrics*, 42, 521-530.
- * Gehan, E.A. (1965). A generalized Wilcoxon test for comparing arbitrarily singlycensored samples. Biomerika 52 203-223.
- * Geskus, R. B. and Groeneboom, P. (1999). Asymptotically optimal estimation of smooth functionals for interval censoring, case 2. Ann. Statist. 27, 627-675.
- * Mantel, N. (1966). Evaluation of survival data and two new rank order statistics arising in its consideration. Cancer Chemother. Rep., 50, 163-170.
* Satten, G.A. (1996). Rank-based inference in the proportional hazards model for interval censored data. *Biometrika*, 83, 355-370.