

Midterm exam, Survival Analysis I, 2017 Spring [+27 points]

Name: 吳柏庭

+23

- Not only answer but also derivations

+4

Q1 [+4] Consider a two-parameter Weibull distribution whose survival function is given by $S_X(x) = \exp(-\lambda x^\alpha)$, where $\lambda > 0, \alpha > 0$. Derive the likelihood function under the following cases.

- + | 1) [+1] Left-truncated data (y_{ti}, x_i) , subject to $y_{ti} \leq x_i, i=1, \dots, n$, where y_{ti} is a left-truncation time (likelihood must be simplified).

$$P(X = x_i | X \geq y_{ti}) = \frac{f(x_i)}{S(y_{ti})}$$

$$L = \prod_{i=1}^n \frac{f(x_i)}{S(y_{ti})} = \prod_{i=1}^n \frac{\lambda \alpha x_i^{\alpha-1} \exp(-\lambda x_i^\alpha)}{\exp(-\lambda y_{ti}^\alpha)} = \prod_{i=1}^n \lambda \alpha x_i^{\alpha-1} \exp(-\lambda (x_i^\alpha - y_{ti}^\alpha))$$

$$= \lambda^n \alpha^n \left(\prod_{i=1}^n x_i \right)^{\alpha-1} \exp(-\lambda \sum_{i=1}^n (x_i^\alpha - y_{ti}^\alpha)) \quad \checkmark$$

- + | 2) [+1] Interval-censored data $(L_i, R_i), i=1, \dots, n$.

$$P(X \in (L_i, R_i)) = S(L_i) - S(R_i)$$

$$L = \prod_{i=1}^n [S(L_i) - S(R_i)] = \prod_{i=1}^n [\exp(-\lambda L_i^\alpha) - \exp(-\lambda R_i^\alpha)] \quad \times$$

- + | 3) [+1] Doubly-truncated data $(y_{ti}, x_i, y_{ri}), i=1, \dots, n$, subject to $y_{ti} \leq x_i \leq y_{ri}, i=1, \dots, n$, (likelihood must be simplified).

$$P(X = x_i | y_{ti} \leq X \leq y_{ri}) = \frac{f(x_i)}{S(y_{ti}) - S(y_{ri})}$$

$$L = \prod_{i=1}^n \frac{f(x_i)}{S(y_{ti}) - S(y_{ri})} = \prod_{i=1}^n \frac{\lambda \alpha x_i^{\alpha-1} \exp(-\lambda x_i^\alpha)}{\exp(-\lambda y_{ti}^\alpha) - \exp(-\lambda y_{ri}^\alpha)}$$

$$= \lambda^n \alpha^n \left(\prod_{i=1}^n x_i \right)^{\alpha-1} \exp(-\lambda \sum_{i=1}^n x_i^\alpha) / \prod_{i=1}^n [\exp(-\lambda y_{ti}^\alpha) - \exp(-\lambda y_{ri}^\alpha)] \quad \times$$

+ |

- 4) [+1] $n=4$ patients whose age-at-death occurring in intervals

(90, 120], (110, 115], (80, 100], (70, 75], subject to the entry condition Age>=50.

$$P(X \in (L_i, R_i) | X \geq y_{ti}) = \frac{S(L_i) - S(R_i)}{S(y_{ti})}$$

$$L = \frac{S(90) - S(120)}{S(50)} \cdot \frac{S(110) - S(115)}{S(50)} \cdot \frac{S(80) - S(100)}{S(50)} \cdot \frac{S(70) - S(75)}{S(50)}$$

$$= [\exp(-\lambda 90^\alpha) - \exp(-\lambda 120^\alpha)] \cdot [\exp(-\lambda 110^\alpha) - \exp(-\lambda 115^\alpha)]$$

$$\cdot [\exp(-\lambda 80^\alpha) - \exp(-\lambda 100^\alpha)] \cdot [\exp(-\lambda 70^\alpha) - \exp(-\lambda 75^\alpha)] / [\exp(-\lambda 50^\alpha)]^4 \quad \times$$

+5

S π +

Q2 [+6]

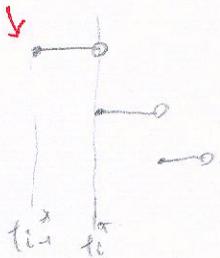
- + | 1) [+] Derive the product-limit form of a survival function $S(x) = \Pr(T > t)$.

$$\begin{aligned} S(t_i^*) &= \frac{S(t_i^*)}{S(t_{i-1}^*)} \cdot \frac{S(t_{i-1}^*)}{S(t_{i-2}^*)} \cdots \frac{S(t_1^*)}{S(0)} \cdot S(0) \\ &= P(T > t_i^* | T > t_{i-1}^*) \cdot P(T > t_{i-1}^* | T > t_{i-2}^*) \cdots P(T > t_1^* | T > 0) \quad \text{①} \\ &= P(T > t_i^* | T > t_{i-1}^*) \cdot P(T > t_{i-1}^* | T > t_{i-2}^*) \cdots P(T > t_1^* | T > 0) \quad \text{②} \\ &= \left[1 - \frac{d(t_i^*)}{N(t_i^*)} \right] \cdot \left[1 - \frac{d(t_{i-1}^*)}{N(t_{i-1}^*)} \right] \cdots \left[1 - \frac{d(t_1^*)}{N(t_1^*)} \right] = \prod_{j=1}^i \left[1 - \frac{d(t_j^*)}{N(t_j^*)} \right]. \end{aligned}$$

- + | 2) [+] What assumption is made to derive the form? Explain why the assumption is necessary.

$S(\cdot)$ is a right continuous ^{step} function
Without this assumption $\textcircled{1} \not\Rightarrow \textcircled{2}$

Step function



- + | 3) [+] Define the Kaplan-Meier estimator for right-censored data (must define all notations clearly).

$$\hat{S}(t) = \prod_{i=1}^n \left[1 - \frac{d_i}{N_i} \right]$$

✓ d_i : numbers of deaths at time t_i^*

✓ N_i : numbers of survivors who were at risk at time t_i^*

- + | 4) [+] Define the product-limit estimator for left-truncated data (define all notations).

$$\hat{S}(t) = \prod_{i=1}^n \left[1 - \frac{d_i}{N_i} \right]$$

✓ d_i : numbers of deaths at time t_i^*

✓ N_i : numbers of survivors who have entered and were at risk at time t_i^*

- + | 5) [+] Explain why the size of the risk set under left-truncation is small initially.

In initial, a first patient enter into the study

Formula can better explain N_i .

But the other patients haven't entered into the study yet, so we cannot observe the events occurred from the other patients.

- + | 6) [+] Explain why one needs to estimate the conditional survival function under left-truncated data.

If event time $X < Y_i$, there is no information about X .

? Not clear enough

+5 Q3 [+5] The hazard function for the ovarian cancer patients follow the Cox model
 $h(t | z_1, z_2, z_3) = h_0(t) \exp(\beta_1 z_1 + \beta_2 z_2 + \beta_3 z_3)$, where the gene expression values are

$$z_1 = \begin{cases} 1 & \text{high value of NCOA3} \\ 0 & \text{low value of NCOA3} \end{cases}, \quad \beta_1 = 0.237$$

$$z_2 = \begin{cases} 1 & \text{high value of TEAD1} \\ 0 & \text{low value of TEAD1} \end{cases}, \quad \beta_2 = 0.223$$

$$z_3 = \begin{cases} 1 & \text{high value of YWHAB} \\ 0 & \text{low value of YWHAB} \end{cases}, \quad \beta_3 = 0.263$$

Compute the relative risk (RR).

+1 1) [+1] RR of (all three genes in high value) vs. (all three genes in low value).

$$RR = \frac{\exp(0.237 + 0.223 + 0.263)}{\exp(0)} = \exp(0.723) = 2.0606$$

+1 2) [+1] RR of (all three genes in high value) vs. (only NCOA3 in high value).

$$RR = \frac{\exp(0.237 + 0.223 + 0.263)}{\exp(0.237)} = \exp(0.486) = 1.6488$$

+1 3) [+1] RR of (only NCOA3 in high value) vs. (only YWHAB in high value).

$$RR = \frac{\exp(0.237)}{\exp(0.263)} = \exp(-0.026) = 0.974$$

+2 4) [+2] All RRs under different combinations risk factors (vs. the baseline risk). Make a table by sorting the RRs (from highest to lowest).

Order	RR	NCOA3	TEAD1	YWHAB
1	2.0606	High	High	High
2	1.6488	High	Low	High
3	1.6258	Low	High	High

4 1.5841 High High Low

Write --> 5 1.301 Low Low High
✓ 6 1.267 High Low Low

7	1.2498	Low	High	Low
8	1	Low	Low	Low

+2 Q4 [+4] There are four stages of cancer (Stages I, II, III and IV). Define three different Cox models according to three different definitions of covariates.

1) [+1] Model 1

$$h(t|z) = h_0(t) \exp(\beta_2 z_2 + (\beta_3 z_3 + \beta_4 z_4))$$



\uparrow
don't start from β_2
(start from β_1)

$$z_1 z_2 = \begin{cases} 1 & , \text{Stage II} \\ 0 & , \text{o.w.} \end{cases}$$

$$z_2 z_3 = \begin{cases} 1 & , \text{Stage III} \\ 0 & , \text{o.w.} \end{cases}$$

$$z_3 z_4 = \begin{cases} 1 & , \text{Stage IV} \\ 0 & , \text{o.w.} \end{cases}$$

+1 RR (Stage IV vs. Stage I) = $\frac{\exp(\beta_4)}{\exp(0)} = \exp(\beta_4)$

2) [+1] Model 2

$$h(t|z) = h_0(t) \exp(\beta_1 z_1 + \beta_3 z_3 + \beta_4 z_4) \quad z_1 = \begin{cases} 1 & , \text{stage I} \\ 0 & , \text{o.w.} \end{cases}$$

All the same model as Model 1.

$$z_3 = \begin{cases} 1 & , \text{stage III} \\ 0 & , \text{o.w.} \end{cases}$$

RR (Stage IV vs. Stage I) = $\frac{\exp(\beta_4)}{\exp(\beta_1)} = \exp(\beta_4 - \beta_1) \quad z_4 = \begin{cases} 1 & , \text{stage IV} \\ 0 & , \text{o.w.} \end{cases}$

3) [+1] Model 3

$$h(t|z) = h_0(t) \exp(\beta_1 z_1 + \beta_2 z_2 + \beta_4 z_4) \quad z_1 = \begin{cases} 1 & , \text{stage I} \\ 0 & , \text{o.w.} \end{cases}$$

$$\frac{h(t|1,0,0) \cdot h(t|0,1,1)}{h(t|1,0,0) \cdot h(t|0,1,2)} =$$

$$z_2 = \begin{cases} 1 & , \text{stage II} \\ 0 & , \text{o.w.} \end{cases}$$

RR (Stage III-IV vs. Stage I-II) = $\frac{\exp(0) \exp(\beta_4)}{\exp(\beta_1 + \beta_2)} - \exp(\beta_4 - \beta_1 - \beta_2) \quad z_4 = \begin{cases} 1 & , \text{stage IV} \\ 0 & , \text{o.w.} \end{cases}$

+1 4) [+1] Which model do you prefer? State your opinion.

Model 1 $h(t|z) = h_0(t) \exp(\beta_2 z_2 + \beta_3 z_3 + \beta_4 z_4)$

Choose stage I as the baseline hazard.

It is natural to choose the lowest stage as the baseline hazard.

$$\begin{aligned} z_2 &= \begin{cases} 1 & , \text{stage II} \\ 0 & , \text{o.w.} \end{cases} \\ z_3 &= \begin{cases} 1 & , \text{stage III} \\ 0 & , \text{o.w.} \end{cases} \\ z_4 &= \begin{cases} 1 & , \text{stage IV} \\ 0 & , \text{o.w.} \end{cases} \end{aligned}$$

+17 Q5 [+8] Consider right-censored data (t_i, δ_i, z_i) , $i = 1, \dots, n$. Assume the Cox model

$h(t | z_i) = h_0(t) e^{\beta z_i}$, where β is one-dimensional (β is a scalar). Define ordered

times at deaths, $t_{(1)}^* \leq t_{(2)}^* \leq \dots \leq t_{(D)}^*$, where $D = \sum_{i=1}^n \delta_i$ (assume no ties).

+1 1) [+1] Define Cox's partial likelihood for β .

$$L(\beta) = \prod_{i=1}^n \left(\frac{\exp(\beta z_i)}{\sum_{j \in R_i} \exp(\beta z_j)} \right)^{\delta_i} \quad R_i = ?$$

+1 2) [+1] Derive the score function

$$\ell(\beta) = \ell_j L(\beta) = \sum_{i=1}^n \delta_i (\beta z_i - \log(\sum_{j \in R_i} \exp(\beta z_j)))$$

$$U(\beta) = \frac{d}{d\beta} \ell(\beta) = \sum_{i=1}^n \delta_i \left(z_i - \frac{\sum_{j \in R_i} z_j \exp(\beta z_j)}{\sum_{j \in R_i} \exp(\beta z_j)} \right)$$

+1 3) [+1] Derive the information matrix

$$I(\beta) = -\frac{\partial^2}{\partial \beta^2} U(\beta) = \sum_{i=1}^n \delta_i \left[\frac{(\sum_{j \in R_i} z_j^2 \exp(\beta z_j))(\sum_{j \in R_i} \exp(\beta z_j)) - (\sum_{j \in R_i} z_j \exp(\beta z_j))^2}{(\sum_{j \in R_i} \exp(\beta z_j))^2} \right]$$

$$= \sum_{i=1}^n \delta_i \left[\frac{\sum_{j \in R_i} z_j^2 \left(\frac{\exp(\beta z_j)}{\sum_{k \neq i} \exp(\beta z_k)} \right) - \left(\sum_{j \in R_i} z_j \left(\frac{\exp(\beta z_j)}{\sum_{k \neq i} \exp(\beta z_k)} \right) \right)^2}{\left(\sum_{j \in R_i} \exp(\beta z_j) \right)^2} \right]$$

+1 4) [+1] Prove that the information matrix is nonnegative.

Define indicator I_i : $P(W_i = z_j) = \frac{\exp(\beta z_j)}{\sum_{k \neq i} \exp(\beta z_k)} = P_j$.

$$\therefore \sum_{j \in R_i} z_j^2 P_j - \left(\sum_{j \in R_i} z_j P_j \right)^2 = \text{Var}(W_i) \geq 0 \quad \therefore I(\beta) = \sum_{i=1}^n \delta_i \text{Var}(W_i) \geq 0$$

+2 4) [+2] Show that the log-rank statistics is equivalent to the score function at $\beta = 0$ under $z_i = 1$ (group 1) and $z_i = 0$ (group 2).

$\log\text{-rank statistic}$ $= \sum_{i=1}^n [d_{i1} - \frac{Y_{i1}}{T_i}]$ $\quad Y_{i1} \frac{d_{i1}}{T_i} (1 - \frac{d_{i1}}{T_i}) \cdot \frac{T_i - Y_{i1}}{T_i - 1}$	$L(0) = \sum_{i=1}^n \delta_i \left(z_{i1} - \frac{\sum_{j \in R_i} z_j}{\sum_{j \in R_i} 1} \right)$ $= \sum_{i=1}^n \delta_i (d_{i1} - \frac{Y_{i1}}{T_i})$ $= \sum_{i=1}^n \left[d_{i1} - \frac{Y_{i1}}{T_i} \right] \quad (d_{i1} = 1)$ $= \log\text{-rank statistic of } z_1$
---	---

+1 5) [+2] In the above setting, show that the variance of the log-rank statistics is equivalent to the information matrix.

$\boxed{\sum_{i=1}^n \frac{d_{i1}}{T_i} (1 - \frac{d_{i1}}{T_i}) \cdot \frac{T_i - Y_{i1}}{T_i - 1}}$ $= ?$	$I(0) = \sum_{i=1}^n \delta_i \left[\frac{\sum_{j \in R_i} z_j^2}{\sum_{j \in R_i} 1} \cdot \frac{1}{\sum_{j \in R_i} 1} - \left(\sum_{j \in R_i} z_j \left(\frac{1}{\sum_{j \in R_i} 1} \right) \right)^2 \right]$ $= \sum_{i=1}^n \delta_i \left[\frac{Y_{i1}}{T_i} - \left(\frac{Y_{i1}}{T_i} \right)^2 \right]$ $= \sum_{i=1}^n \delta_i \left[\frac{T_i Y_{i2}}{T_i T_i} \right]$ $= \sum_{i=1}^n \frac{Y_{i1} Y_{i2}}{T_i^2} \quad (d_{i1} = 1)$ $= \log\text{-rank statistic of } \text{Var}(z_1) \quad \leftarrow \text{write clearly}$
--	---