

# Survival Analysis

Chao Cheng

December 12, 2022

## Contents

<b>1</b>	<b>Basic knowledge</b>	<b>1</b>
1.1	Survival and hazard . . . . .	1
1.2	Censor . . . . .	3
1.2.1	Right censor . . . . .	3
1.2.2	Left censor . . . . .	3
1.2.3	Interval censor . . . . .	3
<b>2</b>	<b>MLE</b>	<b>3</b>
2.1	Parametric MLE . . . . .	4
2.1.1	One-sample setting . . . . .	4
2.1.2	Two-sample setting . . . . .	6
2.2	Nonparametric MLE . . . . .	6
2.2.1	Discrete time points . . . . .	6
2.2.2	Continuous time points . . . . .	8
2.2.3	Some extensions . . . . .	9

## 1 Basic knowledge

### 1.1 Survival and hazard

Let  $T$  denote the time to an event that we are interested in. Then we know the c.d.f.

$$F_T(t) = P(T \leq t),$$

and the corresponding p.d.f.

$$f_T(t) = \frac{d}{dt} F_T(t).$$

Here to simplify the discussion, we assume  $T$  is a continuous random variable. In the context of survival analysis, the *event* often refers to death. Then  $T$  represents the lifespan of the subject. So  $F_T(t)$  represents the probability that the death occurs before  $t$ . In another word, we know the probability that the subject survives passes  $t$  is

$$S_T(t) = 1 - F_T(t) = P(T > t).$$

$S_T(t)$  is often called the **survival function?** and clearly

$$f_T(t) = -\frac{d}{dt} S_T(t).$$

The **hazard function**  $h(t)$  is defined as

$$h(t) = \lim_{\Delta \rightarrow 0} \frac{P(T \leq t + \Delta | T > t)}{\Delta} = \lim_{\Delta \rightarrow 0} \frac{F_T(t + \Delta) - F_T(t)}{\Delta \cdot S_T(t)} = \frac{f_T(t)}{S_T(t)}. \quad (1)$$

$h(t)$  represents the **instant hazard? unified probability?** that the subject will be dead instantly after  $t$  given the fact that it's alive at  $t$ . And the **cummulative hazard function** is

$$H(t) = \int_0^t h(x) dx = \int_0^t \frac{f_T(x)}{S_T(x)} dx = \int_0^t \frac{-dS_T(x)}{S_T(x)} = -\log(S_T(x))|_0^t = -\log(S_T(t)).$$

**Proposition 1.** *The random variable  $H(T)$  follows unit exponential distribution  $EXP(1)$ .*

*Proof.*

$$\begin{aligned} P(H(T) \leq t) &= P(-\log S(T) \leq t) \\ &= P(1 - F(T) \geq e^{-t}) \\ &= P(T \leq F^{-1}(1 - e^{-t})) \\ &= F(F^{-1}(1 - e^{-t})) \\ &= 1 - e^{-t}, \end{aligned}$$

which is the c.d.f of  $EXP(1)$ . Here to simplify the deduction we make some assumptions that

- $F(t)$  is continuous.
- $F^{-1}(t)$  is well defined.

Also to simplify the notation and avoid confusion, we use  $S(\cdot)$  and  $F(\cdot)$  instead of  $S_T(\cdot)$  and  $F_T(\cdot)$  like before.  $\square$

**1. Exponential distribution:** Denote  $T \sim EXP(\lambda)$ . Then

$$\begin{aligned} f(t) &= \lambda e^{-\lambda t} \\ F(t) &= 1 - e^{-\lambda t} \quad S(t) = e^{-\lambda t} \\ h(t) &= \lambda \quad \text{constant hazard} \\ H(t) &= \lambda t \\ E(T) &= 1/\lambda \quad \text{Var}(T) = 1/\lambda^2 \end{aligned}$$

**2. Weibull distribution:** Denote  $T \sim W(p, \lambda)$ . Then

$$\begin{aligned} f(t) &= p\lambda^p t^{p-1} \exp(-(\lambda t)^p) \\ F(t) &= 1 - \exp(-(\lambda t)^p) \quad S(t) = \exp(-(\lambda t)^p) \\ h(t) &= p\lambda^p t^{p-1} \\ H(t) &= (\lambda t)^p \\ E(T) &= \frac{1}{\lambda} \cdot \Gamma\left(1 + \frac{1}{p}\right) \quad \text{Var}(T) = \frac{1}{\lambda^2} \left( \Gamma\left(1 + \frac{2}{p}\right) - \Gamma\left(1 + \frac{1}{p}\right)^2 \right) \\ E(T^m) &= \frac{1}{\lambda^m} \Gamma\left(1 + \frac{m}{p}\right) \end{aligned}$$

## 1.2 Censor

### 1.2.1 Right censor

- Type I: an i.i.d sample  $T_1, \dots, T_n \sim F$  and a **fixed** constant  $c$ . And the observed data is  $(U_i, \delta_i)$  for  $i = 1, \dots, n$  where

$$U_i = \min(T_i, c)$$
$$\delta_i = 1_{T_i \leq c}.$$

So the observed data consists of a **random** number,  $r$ , of uncensored observations, all of which are less than  $c$ . And  $n - r$  censored observations, all are  $c$ .

- Type II: an i.i.d sample  $T_1, \dots, T_n \sim F$  and a **pre-defined** number of failure  $r$ . The observation is stopped when  $r$  failure occurs and the stopping time is  $c$ . The observed data is still the form  $(U_i, \delta_i)$  for  $i = 1, \dots, n$ , the same as that in Type I censor. But in actuality, we observe the first  $r$  **order statistics**

$$T_{(1,n)}, \dots, T_{(r,n)}.$$

Note that here  $(U_1, \delta_1), \dots, (U_n, \delta_n)$  are **dependent** whereas they are independent for Type I.

- Type III (Random censor): The underlying data is

$$c_1, \dots, c_n \text{ constant}$$
$$T_1, \dots, T_n \sim F.$$

And the observed data is  $(U_i, \delta_i)$  for  $i = 1, \dots, n$ , where

$$U_i = \min(T_i, c_i)$$
$$\delta_i = 1_{T_i \leq c_i}.$$

**Note:** for inference,  $c_i$  is often treated as constant. For study design or studying the asymptotic property, they are often treated as i.i.d random variables  $C_1, \dots, C_n$ .

### 1.2.2 Left censor

$T_i$  is censored when  $T_i \leq l_i$ .

### 1.2.3 Interval censor

$l_i \leq T_i \leq u_i$ , but only  $l_i$  and  $u_i$  are observed.

## 2 MLE

There is an i.i.d survival time sample  $T_1, \dots, T_n$  with common and unknown c.d.f.  $F(\cdot)$  and the observed data is  $(U_i, \delta_i)$  for  $i = 1, \dots, n$ , where

$$U_i = \min(T_i, C_i)$$
$$\delta_i = 1(T_i \leq C_i)$$

and  $C_i$  is the (fixed or random) censoring time. Let  $\perp$  denote “is independent of”. We assume  $T_i \perp C_i$  (Non-informative censoring, the key assumption) and  $(U_i, \delta_i)$  are also i.i.d. The observed data consists of two parts.  $U_i$  is continuous while  $\delta_i$  is binary.

$$\begin{aligned} (U_i, \delta_i) &= (u_i, 1) & T_i \text{ is uncensored at } u_i \\ (U_i, \delta_i) &= (u_i, 0) & T_i \text{ is censored at } u_i \end{aligned}$$

When  $C_i$ s are known constants, the likelihood for  $(U_i, \delta_i)$  is

$$\begin{aligned} L_i(F) &= \begin{cases} f(u_i) & \text{if } \delta_i = 1 \\ 1 - F(u_i) & \text{if } \delta_i = 0 \end{cases} \\ &= f(u_i)^{\delta_i} (1 - F(u_i))^{1-\delta_i} \end{aligned}$$

Therefore

$$L(F) = \prod_{i=1}^n L_i(F) = \prod_{i=1}^n \left( f(u_i)^{\delta_i} (1 - F(u_i))^{1-\delta_i} \right) = \prod_{i=1}^n \left( h(u_i)^{\delta_i} S(u_i) \right). \quad (2)$$

The last equality relies on the fact that  $f(t) = h(t) S(t)$ .

When  $C_i$ s are i.i.d.  $\sim G$ , where  $G$  is continuous with p.d.f  $g$ . Then we have

$$P(U_i \leq u, \delta_i = 1) = P(T_i \leq u, T_i \leq C_i) = \int_0^u \int_t^\infty f(t) g(c) dc dt = \int_0^u f(t) (1 - G(t)) dt$$

Therefore the likelihood for  $\delta_i = 1$  is

$$L_i(F, G) = f(u_i) (1 - G(u_i)) \quad \text{when } \delta_i = 1.$$

And similarly, for  $\delta_i = 0$ , the likelihood is

$$L_i(F, G) = g(u_i) (1 - F(u_i)) \quad \text{when } \delta_i = 0.$$

Hence the full likelihood is

$$\begin{aligned} L(F, G) &= \prod_{i=1}^n \left\{ (f(u_i) (1 - G(u_i)))^{\delta_i} ((1 - F(u_i)) g(u_i))^{1-\delta_i} \right\} \\ &= \prod_{i=1}^n \left\{ f(u_i)^{\delta_i} (1 - F(u_i))^{1-\delta_i} \right\} \cdot \prod_{i=1}^n \left\{ g(u_i)^{1-\delta_i} (1 - G(u_i))^{\delta_i} \right\} \end{aligned} \quad (3)$$

So the core to maximize  $L(F, G)$  with respect to  $F$  in (3) is the same as that in (2).

## 2.1 Parametric MLE

### 2.1.1 One-sample setting

Suppose  $T_1, \dots, T_n$  are i.i.d.  $Exp(\lambda)$ , and subject to noninformative right censoring. Then (2) becomes

$$L = L(\lambda) = \prod_{i=1}^n \left\{ (\lambda e^{-\lambda u_i})^{\delta_i} (e^{-\lambda u_i})^{1-\delta_i} \right\} = \lambda^{\sum_{i=1}^n \delta_i} e^{-\lambda \sum_{i=1}^n u_i} = \lambda^r e^{-\lambda W},$$

where  $r = \sum_{i=1}^n \delta_i$  is the number of observed events and  $W = \sum_{i=1}^n u_i$  is the total of observed time. Therefore  $\log L = r \log \lambda - \lambda W$  and the MLE for  $\lambda$  is

$$\hat{\lambda} = \frac{r}{W}.$$

Furthermore, we know that

$$\begin{cases} \frac{\partial \log L}{\partial \lambda} = \frac{r}{\lambda} - W \\ \frac{\partial^2 \log L}{\partial \lambda^2} = -\frac{r}{\lambda^2} \end{cases}.$$

Based on properties of fisher information ([See the notes about fisher information for more details.](#)), we know that at the **true underlying value**  $\lambda$ , it must satisfy

$$\begin{cases} E \frac{\partial \log L}{\partial \lambda} = \frac{Er}{\lambda} - EW = 0 \\ I(\lambda) = -E \frac{\partial^2 \log L}{\partial \lambda^2} = \frac{Er}{\lambda^2} \\ I^*(\lambda) = \frac{1}{n} I(\lambda) = \frac{Er}{n\lambda^2} \end{cases}. \quad (4)$$

Note that in (4),  $r$  and  $W$  are random variables. And the probability to observe an event is

$$p = P(\delta_i = 1) = P(U_i \leq \infty, \delta_i = 1) = \int_0^\infty f(t) (1 - G(t)) dt.$$

Therefore  $r \sim \text{binomial}(n, p)$ ,  $Er = np$ . And from property of MLE, we can write

$$\frac{\sqrt{n}(\hat{\lambda} - \lambda)}{\sqrt{I^*(\lambda)^{-1}}} = \frac{(\hat{\lambda} - \lambda)}{\sqrt{I(\lambda)^{-1}}} \xrightarrow{D} N(0, 1),$$

which means approximately

$$\hat{\lambda} \stackrel{\text{apx}}{\sim} N(\lambda, I(\lambda)^{-1}) = N\left(\lambda, \frac{\lambda^2}{np}\right).$$

Unfortunately, both  $\lambda$  and  $p$  (essentially  $G(\cdot)$ ) are unknown. We plug in the estimation  $\hat{\lambda} = r/W$  and  $\hat{p} = r/W$  and apply Slutsky's theorem. This means for the purpose of estimation, we use

$$\begin{cases} \hat{\lambda} = \frac{r}{W} \\ I(\hat{\lambda}) = \frac{r}{\hat{\lambda}^2}, \quad I^*(\hat{\lambda}) = \frac{r}{n\hat{\lambda}^2} \end{cases} \quad (5)$$

Not that unlike (4), here in (5),  $r$  and  $W$  are observations. And we have

$$\hat{\lambda} \stackrel{\text{apx}}{\sim} N\left(\lambda, \frac{r}{W^2}\right). \quad (6)$$

Note that it turns out that a better approximation is to assume  $\log \hat{\lambda}$  is normal. Using the delta method, this gives

$$\log \hat{\lambda} \stackrel{\text{apx}}{\sim} N\left(\log \lambda, \frac{1}{np}\right) \approx N\left(\log \lambda, \frac{1}{r}\right). \quad (7)$$

Now based on (6) or (7), we can construct CI on  $\lambda$ , which also means we can perform hypothesis testing about  $\lambda$ .

### 2.1.2 Two-sample setting

For two samples  $x_1, \dots, x_n$  and  $y_1, \dots, y_m$ , both follow exponential distribution with parameters  $\lambda_1$  and  $\lambda_2$ . Assume noninformative censoring in each group, using same tech in Section 2.1.1 we can get

$$Z = \frac{\log \hat{\lambda}_1 - \log \hat{\lambda}_2}{\sqrt{\frac{1}{r_1} + \frac{1}{r_2}}} \stackrel{\text{apx}}{\sim} N(0, 1).$$

## 2.2 Nonparametric MLE

The NPMLE of survivor function  $S(\cdot)$  based on i.i.d. survival time and non-informative right censoring is often known as Kaplan-Meier estimator or the Product-Limit Estimator. Here we provide some heuristic development, but formal proofs will be deferred to other notes. With the same notation as before, the observed data is

$$U_i = \min(T_i, C_i), \quad \delta_i = 1(T_i \leq C_i),$$

where  $T_i$ s are i.i.d survival times and  $C_i$ s are i.i.d **non-informative** censoring time. The full likelihood is already shown in (3).

### 2.2.1 Discrete time points

To begin with, let's assume  $F(\cdot)$  takes discrete values with mass points at  $\{v_i\}$ s:  $0 \leq v_1 < v_2 < \dots < \dots$ , and define the discrete hazard functions as

$$\begin{aligned} h_1 &= P(T = v_1) \\ h_j &= P(T = v_j | T > v_{j-1}) \quad j > 1. \end{aligned} \tag{8}$$

Note that (8) can be seen as discrete version of (1). And for  $t \in [v_j, v_{j+1})$ ,

$$\begin{aligned} S(t) &\stackrel{\text{def}}{=} P(T > t) = P(T > v_j) \\ &= P(T > v_j | T > v_{j-1}) P(T > v_{j-1}) \\ &= P(T > v_j | T > v_{j-1}) P(T > v_{j-1} | T > v_{j-2}) P(T > v_{j-2}) \\ &= \dots \\ &= P(T > v_1) \prod_{i=1}^{j-1} P(T > v_{i+1} | T > v_i) \\ &= \prod_{i=1}^j (1 - h_i) \quad j > 1. \end{aligned}$$

For discrete case, the p.m.f  $f(\cdot)$  is

$$\begin{aligned} f(v_1) &= P(T = v_1) = h_1 \\ f(v_j) &= P(T = v_j) = P(T = v_j | T > v_{j-1}) P(T > v_{j-1}) = h_j \prod_{i=1}^{j-1} (1 - h_i). \end{aligned}$$

Then if we want to estimate  $F(\cdot)$  from likelihood, either (2) or (3), we are just trying to maximizing

$$\begin{aligned} L(F) &= \prod_{i=1}^n \left\{ f(u_i)^{\delta_i} (1 - F(u_i))^{1-\delta_i} \right\} \\ &= \prod_{\{u_i|\delta_i=1\}} f(u_i) \prod_{\{u_i|\delta_i=0\}} S(u_i). \end{aligned}$$

Let  $I(\cdot)$  be an index mapping function that returns the index in  $v_i$ s that matches  $u_i$ , i.e.  $I(u_i) = j$  if and only if  $u_i \in [v_j, v_{j+1})$ . Then we know that  $u_i = v_{I(u_i)}$  and we can write

$$\begin{aligned} L(F) &= \prod_{\{u_i|\delta_i=1\}} f(v_{I(u_i)}) \prod_{\{u_i|\delta_i=0\}} S(v_{I(u_i)}) \\ &= \left[ \prod_{\{u_i|\delta_i=1, u_i=v_1\}} f(v_1) \right] \left[ \prod_{\{u_i|\delta_i=1, u_i \neq v_1\}} f(v_{I(u_i)}) \right] \left[ \prod_{\{u_i|\delta_i=0\}} \prod_{k=1}^{I(u_i)} (1 - h_k) \right] \\ &= \left[ \prod_{\{u_i|\delta_i=1, u_i=v_1\}} h_1 \right] \left[ \prod_{\{u_i|\delta_i=1, u_i \neq v_1\}} \left( h_{I(u_i)} \prod_{k=1}^{I(u_i)-1} (1 - h_k) \right) \right] \left[ \prod_{\{u_i|\delta_i=0\}} \prod_{k=1}^{I(u_i)} (1 - h_k) \right] \\ &= \left[ \prod_{\{u_i|\delta_i=1\}} h_{I(u_i)} \right] \left[ \prod_{\{u_i|\delta_i=1, u_i \neq v_1\}} \prod_{k=1}^{I(u_i)-1} (1 - h_k) \right] \left[ \prod_{\{u_i|\delta_i=0\}} \prod_{k=1}^{I(u_i)} (1 - h_k) \right]. \end{aligned} \tag{9}$$

Note that in (9), the first part is

$$\prod_{\{u_i|\delta_i=1\}} h_{I(u_i)} = \prod_{j=1}^{\infty} h_j^{d_j}, \tag{10}$$

where  $d_j = \sum_{i=1}^n \delta_i \cdot 1(u_i = v_j)$  is the number of event at  $v_j$ . The second and third part in (9) is

$$\begin{aligned} & \left[ \prod_{\{u_i|\delta_i=1, u_i \neq v_1\}} \prod_{k=1}^{I(u_i)-1} (1 - h_k) \right] \left[ \prod_{\{u_i|\delta_i=0\}} \prod_{k=1}^{I(u_i)} (1 - h_k) \right] \\ &= \left[ \prod_{k=1}^{\infty} \prod_{\{i|\delta_i=1, I(u_i)-1 \geq k\}} (1 - h_k) \right] \left[ \prod_{k=1}^{\infty} \prod_{\{i|\delta_i=0, I(u_i) \geq k\}} (1 - h_k) \right] \\ &= \left[ \prod_{k=1}^{\infty} (1 - h_k)^{\sum_{i=1}^n \delta_i \cdot 1(I(u_i)-1 \geq k)} \right] \left[ \prod_{k=1}^{\infty} (1 - h_k)^{\sum_{i=1}^n (1-\delta_i) \cdot 1(I(u_i) \geq k)} \right] \\ &= \left[ \prod_{k=1}^{\infty} (1 - h_k)^{\sum_{i=1}^n \delta_i \cdot [1(I(u_i) \geq k) - 1(I(u_i) = k)]} \right] \left[ \prod_{k=1}^{\infty} (1 - h_k)^{\sum_{i=1}^n (1-\delta_i) \cdot 1(I(u_i) \geq k)} \right] \\ &= \prod_{k=1}^{\infty} (1 - h_k)^{Y(v_k) - d_k}, \end{aligned} \tag{11}$$

where

$$Y(v_k) = \sum_{i=1}^n 1(I(u_i) \geq k) = \sum_{i=1}^n 1(u_i \geq v_k)$$

is the number of subjects that are “at risk” at time  $v_k$ . **Note:** by the word “at risk”, we also count the subjects that died just at  $v_k$ , which means  $Y(v_j) \geq d_j$ . But we do NOT count the subjects that are censored before  $v_j$ .

Then from (10) and (11) we know that (9) can be written as

$$L(F) = \prod_{j=1}^{\infty} h_j^{d_j} (1 - h_j)^{Y(v_j) - d_j}. \quad (12)$$

And the NPMLE is just

$$\hat{h}_j = \frac{d_j}{Y(v_j)} \quad (13)$$

for  $j = 1, \dots, \infty$  and  $Y(v_j) > 0$ . (13) implies some properties of this discrete NPMLE:

1. This estimation makes sense: the probability of dying at  $v_j$  given the fact you live past  $v_{j-1}$  can be estimated by the proportion of subjects die at  $v_j$  over the number of “at risk” at  $v_j$ .
2.  $\hat{h}_j$  is only defined at time points where  $Y(v_j) > 0$ . Therefore, for large enough  $v_j$ , there will be no observation, no matter event or censoring, resulting inability to make estimation about hazard at those time points.
3. For time points where  $Y(v_j) > 0$  but no event occurs, the hazard is estimated to be 0.

This means

$$\hat{S}(t) = \begin{cases} 1 & t < v_1 \\ \prod_{j=1}^k (1 - \hat{h}_j) & v_k \leq t < v_{k+1} \end{cases} \quad (14)$$

**Note:**  $S(\cdot)$  is defined to be **right-continuous**.

Let  $v_g$  denotes the largest time point with observation, which means  $Y(v_g) > 0$  and  $Y(v_{g+1}) = 0$ . Then either  $d_g = Y(v_g)$  or  $d_g < Y(v_g)$ . If  $d_g = Y(v_g)$ , then  $\hat{h}_g = 1$  and  $\hat{S}(t) = 0$  for  $t \geq v_g$ . But if  $d_g < Y(v_g)$ , then  $\hat{S}(t) > 0$  for  $v_g \leq t < v_{g+1}$  and  $\hat{S}(t)$  is undefined on  $t \in [v_{g+1}, \infty)$ .

Here one might say that the KM estimator is undefined on  $t \in [v_{g+1}, \infty)$ . Or another explanation is that NPMLE is not unique and any survival function that is identical to  $\hat{S}$  at previous time is the NPMLE.

### 2.2.2 Continuous time points

Now, if we don’t know in advance the times at which  $F$  had mass, or even did not want to assume  $F$  was discrete distribution? The core of likelihood still takes the form of (2), but now we have to maximize it over all c.d.f., including discrete, continuous and mixed distributions.

Kaplan and Meier argue that the solution must be a discrete distribution with mass on the observed times  $u_i$  only. That is, the KM (product-limit) estimator of  $F(\cdot)$  is

$$\hat{S}(t) = \begin{cases} 1 & t < v_1 \\ \prod_{j=1}^k \left(1 - \frac{d_j}{Y(v_j)}\right) & v_k \leq t < v_{k+1} \end{cases} \quad (15)$$



Note that (15) only puts weight at the observed (**uncensored**) failure time. Another (equivalent) representation of  $\hat{S}(t)$  is given by

$$\hat{S}(t) = \prod_{j:v_j \leq t} \left( \frac{Y(v_j) - d_j}{Y(v_j)} \right) \quad \text{for } t \leq \max(v_i), \quad (16)$$

where  $v_1 < v_2 < \dots$  are the distinct observed failure time.

### 2.2.3 Some extensions

**Other perspective for this NPMLE** Besides the KM-estimator, there is also Efron's "Redistribution of Mass" algorithm that gives the same results.

Another point of view is that the KM estimator can be seen as the self-consistency estimator. If there's no censoring, the survival function can be estimated as

$$\hat{S}(t) = n^{-1} \sum_{i=1}^n 1(T_i > t).$$

In the precense of censoring, and the observed data  $\{(U_i, \delta_i), i = 1, \dots, n\}$ , the survival function can be estimated as

$$\hat{S}(t) = n^{-1} \sum_{i=1}^n E\{1(T_i > t) | U_i, \delta_i\}$$

where

$$\begin{aligned} E\{1(T > t) | U_i, \delta_i = 1\} &= 1(U_i > t) \\ E\{1(T > t) | U_i, \delta_i = 0\} &= E(1(T > t, U_i \leq t) + 1(T > t, U_i > t) | U_i, \delta_i = 0) \\ &= E(1(T > t) | U_i = u_i \leq t, \delta_i = 0) 1(u_i \leq t) + 1(U_i > t) \\ &= E(1(T > t) | T > u_i) 1(u_i \leq t) + 1(U_i > t) \\ &= P(T > t | T > u_i) 1(u_i \leq t) + 1(U_i > t) \\ &= S(t) / S(u_i) 1(t \geq u_i) + 1(t < u_i). \end{aligned}$$

Unfortunately,  $S(\cdot)$  is unknown, and we calculate  $\hat{S}(\cdot)$  iteratively via

$$\begin{aligned} \hat{S}_{new}(t) &= n^{-1} \sum_{i=1}^n \left\{ \delta_i \cdot 1(U_i > t) + (1 - \delta_i) \cdot 1(U_i \leq t) \cdot \frac{\hat{S}_{old}(t)}{\hat{S}_{old}(U_i)} + (1 - \delta_i) \cdot 1(U_i > t) \right\} \\ &= n^{-1} \sum_{i=1}^n \left\{ 1(U_i > t) + (1 - \delta_i) \cdot 1(U_i \leq t) \cdot \frac{\hat{S}_{old}(t)}{\hat{S}_{old}(U_i)} \right\} \end{aligned}$$

And the limit  $\hat{S}(t)$  solves

$$\hat{S}(t) = n^{-1} \sum_{i=1}^n \left\{ 1(U_i > t) + (1 - \delta_i) \cdot 1(U_i \leq t) \cdot \frac{\hat{S}(t)}{\hat{S}(U_i)} \right\},$$

which gives the same results as KM estimator.

**NP estimator for  $H(t)$ :** Since  $H(t) = -\log(S(t))$ , it follows that a nonparametric estimator for  $H(t)$  is

$$\tilde{H}(t) = -\log(\hat{S}(t)) = -\sum_{i=1}^k \log(1 - \hat{h}_i) \quad \text{for } v_k \leq t < v_{k+1}.$$

Note that  $\log(1 - x) \approx \log(1) + \frac{-1}{1-x}|_{x=0} \cdot x = -x$ , therefore an alternative estimator (for  $k \geq 1$ ) is

$$\hat{H}(t) = \sum_{i=1}^k \hat{h}_i \quad \text{for } v_k \leq t < v_{k+1}.$$

And this is called the **Nelson-Aalen** estimator of  $H(\cdot)$ .

**Inference on  $h_j$  and  $S(t)$ :** Now what if we want to approximate the distribution of  $\hat{S}(t)$ ? One can use the large-sample property of MLE (but the ordinary regularity conditions do not hold here since it is not a finite-dimensional parameter space). Nevertheless, let's proceed as if this is not a problem. Then from (12) we know that

$$\begin{aligned} -\frac{\partial^2 \log L}{\partial h_i \partial h_j} &= 0 \quad \text{for } i \neq j \\ -\frac{\partial^2 \log L}{\partial h_j^2} &= -\frac{\partial^2}{\partial h_j^2} (d_j \log h_j + (Y(v_j) - d_j) \log(1 - h_j)) \\ &= -\frac{\partial}{\partial h_j} \left( \frac{d_j}{h_j} - \frac{Y(v_j) - d_j}{1 - h_j} \right) \\ &= -(-d_j h_j^{-2} - (Y(v_j) - d_j)(1 - h_j)^{-2}) \end{aligned}$$

and

$$-\frac{\partial^2 \log L}{\partial h_i \partial h_j} \Big|_{h_j = \hat{h}_j} = \frac{Y(v_j)}{\hat{h}_j (1 - \hat{h}_j)}$$

since  $\hat{h}_j = \frac{d_j}{Y(v_j)}$ . So the hessian matrix of  $\log L$  is diagonal, which means  $\hat{h}_1, \hat{h}_2, \dots$  are approximately uncorrelated, with the approximated means  $h_1, h_2, \dots$  and

$$\text{Var}(\hat{h}_j) \approx \frac{\hat{h}_j (1 - \hat{h}_j)}{Y(v_j)} = \frac{d_j (Y(v_j) - d_j)}{Y(v_j)^3}.$$

Therefore approximately

$$\hat{h}_j \stackrel{\text{apx}}{\sim} N \left( h_j, \frac{d_j (Y(v_j) - d_j)}{Y(v_j)^3} \right).$$

Furthermore, from (14) we know that in the discrete time setting  $\hat{S}(t)$  is approximately

unbiased. Now for the variance of  $\hat{S}(t)$ , we can write, for  $v_j \leq t < v_{j+1}$  that

$$\begin{aligned}
\text{Var} \left( \log \hat{S}(t) \right) &= \text{Var} \left( \sum_{j=1}^k \log (1 - \hat{h}_j) \right) \\
&\approx \sum_{j=1}^k \text{Var} \left( \log (1 - \hat{h}_j) \right) \quad \hat{h}_j\text{s are approximatly uncorrelated} \\
&\approx \sum_{j=1}^k \text{Var} \left( \hat{h}_j \right) \cdot \frac{1}{(1 - \hat{h}_j)^2} \quad \delta - method \\
&= \sum_{j=1}^k \frac{d_j}{Y(v_j)(Y(v_j) - d_j)}.
\end{aligned}$$

Also from  $\delta$ -method we know that

$$\text{Var} \left( \hat{S}(t) \right) \approx \text{Var} \left( \log \hat{S}(t) \right) \cdot \left( e^{\log \hat{S}(t)} \right)^2 \approx \hat{S}(t)^2 \sum_{j=1}^k \frac{d_j}{Y(v_j)(Y(v_j) - d_j)} \quad \text{for } v_j \leq t < v_{j+1}. \quad (17)$$

And (17) is often called the **Greenwood's formula**. And from this we can construct a confidence interval for  $S(t)$ . One choice for an 95% CI would be

$$\hat{S}(t) \pm 1.96 \sqrt{\text{Var} \left( \hat{S}(t) \right)}.$$

However, it's all possible for this type of CI to be out the range of  $[0, 1]$ . One alternative is to make CI on  $\log(-\log S(t))$ , which takes value from  $-\infty$  to  $\infty$ , then transform back to the scale of  $S(t)$ . Again by  $\delta$ -method

$$\text{Var} \left( \log \left( -\log \hat{S}(t) \right) \right) \approx \frac{\sum_{j=1}^k \frac{d_j}{Y(v_j)(Y(v_j) - d_j)}}{\left( \log \hat{S}(t) \right)^2}$$

Then a 95% CI for  $S(t)$  can be expressed as

$$\begin{aligned}
& \log \left( -\log \left( \hat{S}(t) \right) \right) \pm 1.96 \sqrt{\frac{\sum_{j=1}^k \frac{d_j}{Y(v_j)(Y(v_j)-d_j)}}{(\log \hat{S}(t))^2}} \\
\Rightarrow & \left( -\log \hat{S}(t) \right) \cdot \exp \left( \pm 1.96 \sqrt{\frac{\sum_{j=1}^k \frac{d_j}{Y(v_j)(Y(v_j)-d_j)}}{(\log \hat{S}(t))^2}} \right) \\
\Rightarrow & \hat{S}(t) \exp \left( \pm 1.96 \sqrt{\frac{\sum_{j=1}^k \frac{d_j}{Y(v_j)(Y(v_j)-d_j)}}{(\log \hat{S}(t))^2}} \right) \\
\Rightarrow & \left[ \hat{S}(t) \exp \left( 1.96 \sqrt{\frac{\sum_{j=1}^k \frac{d_j}{Y(v_j)(Y(v_j)-d_j)}}{(\log \hat{S}(t))^2}} \right), \hat{S}(t) \exp \left( -1.96 \sqrt{\frac{\sum_{j=1}^k \frac{d_j}{Y(v_j)(Y(v_j)-d_j)}}{(\log \hat{S}(t))^2}} \right) \right].
\end{aligned}$$

Breslow and Crowley (1974)(need reference here) show that as  $n \rightarrow \infty$ ,

$$\sqrt{n} \left( \hat{S}(\cdot) - S(\cdot) \right) \xrightarrow{w} \text{zero mean Gaussian process},$$

which can be easily proved by if we express the KM estimator as a martingale process.

I think this previous work is summarised in Brookmeyer and Crowley(1982)(need reference here), that we can consider a transform function  $g$ , such that  $g(S(t))$  follows normal distribution. Again from  $\delta$ -method, we know that

$$\text{Var} \left( g \left( \hat{S}(t) \right) \right) \approx \left( g' \left( \hat{S}(t) \right) \right)^2 \text{Var} \left( \hat{S}(t) \right).$$

Some popular transformations are

linear:	$g(x) = x$
log:	$g(x) = \log x$
log-log:	$g(x) = \log(-\log x)$
hall-werner, loghall, epband, logep:	This is simultaneous CI, unlike before are point-wise.

So we can construct the CI based on  $g(S(t))$  then transform back to original scale:

$$\begin{aligned}
& g^{-1} \left\{ g(S(t)) \pm z_{\alpha/2} \times \sqrt{\text{Var}(g(S(t)))} \right\} \\
& \approx g^{-1} \left\{ g(\hat{S}(t)) \pm z_{\alpha/2} \times g'(\hat{S}(t)) \sqrt{\text{Var}(\hat{S}(t))} \right\}.
\end{aligned} \tag{18}$$

**CI for survival time at given survival rate:** The KM estimator can be used to estimated survival time and construct corresponding CI at given survival probability. For example, the median survival time  $t_m$  can be estimated as  $\hat{t}_m$ , the solution such as

$$\hat{S}(\hat{t}_m) = 0.5.$$

Also from previous discussion, we can get  $\hat{S}_L(\cdot)$  and  $\hat{S}_U(\cdot)$  as the lower and upper bound of the 95% CI for  $S(\cdot)$ , i.e.

$$P(\hat{S}_L(t) \leq S(t)) = P(\hat{S}_U(t) \geq S(t)) = 0.975 \quad \forall t.$$

For  $\hat{S}_L(t)$  and  $\hat{S}_U(t)$  we can also find  $\hat{t}_{ml}$  and  $\hat{t}_{mu}$  such that

$$\hat{S}_L(\hat{t}_{ml}) = \hat{S}_U(\hat{t}_{mu}) = 0.5.$$

Then we can deduct that

$$\begin{aligned} P(\hat{t}_{ml} \leq t_m) &= P(\hat{S}_L(\hat{t}_{ml}) \geq \hat{S}_L(t_m)) = P(0.5 \geq \hat{S}_L(t_m)) = P(S(t_m) \geq \hat{S}_L(t_m)) = 0.975 \\ P(\hat{t}_{mu} \geq t_m) &= P(\hat{S}_U(\hat{t}_{mu}) \leq \hat{S}_U(t_m)) = P(0.5 \leq \hat{S}_U(t_m)) = P(S(t_m) \leq \hat{S}_U(t_m)) = 0.975, \end{aligned}$$

which means a 95% CI for  $t_m$  is just  $[\hat{t}_{ml}, \hat{t}_{mu}]$ .

**CI for survival rate difference at given time point:** From (18) and we choose the linear transformation  $g(x) = x$ , then basically we are saying  $\hat{S}(t)$  approximately follows normal distribution and the variance is given by Greenwood's formula (17). Then the CI for survival rate at given time point can be seen as the CI for difference of two normal distributed variables.

- It's harder to construct CI for survival rate difference when choosing other  $g(x)$ .
- One can consider using bootstrap to construct this CI. Some reference R scripts are

```
set.seed(1234)

library(survival)
library(boot)

km <- survfit(Surv(time, status) ~ sex, data = lung)
year1_surv_est <- summary(km, t = 365)

## Method 1
year1_point_male <- year1_surv_est$surv[1]
year1_std_male <- year1_surv_est$std.err[1]

year1_point_female <- year1_surv_est$surv[2]
year1_std_female <- year1_surv_est$std.err[2]

diff_est <- year1_point_male - year1_point_female
## > -.19 (male survival is 19% lower than females at 1 year)
```

```

diff_std <- sqrt(year1_std_male^2 + year1_std_female^2)
## > .073 (standard error around this estimate of .073)

diff_est + c(-1, 1) * 1.96 * diff_std
## > (-.335, -.046) (95% CI)

## Bootstrapping
boot_func <- function(data, index) {
  data <- data[index, ]

  km <- survfit(Surv(time, status) ~ sex, data = data)

  year1_surv_est <- summary(km, t = 365)
  year1_point_male <- year1_surv_est$surv[1]
  year1_point_female <- year1_surv_est$surv[2]

  year1_point_male - year1_point_female
}

boot_obj <- boot(lung, boot_func, R = 10000)
boot.ci(boot_obj, type = "perc")
## < (-0.334, -0.043) (95% CI via percentile)

```

This discussion can be found at <https://discourse.datamethods.org/t/kaplan-meier-se-for-absolute-difference-in-time-point-survival/5035/14>

**Restricted mean survival time:** The area under survival time curve for any given  $\tau$  is defined as

$$\mu = \int_0^{\tau} S(t) dt.$$

We can show that

$$\begin{aligned}
 \mu &= t(1 - F(t))|_0^{\tau} - \int_0^{\tau} t(-f(t)) dt \\
 &= \tau S(\tau) + \int_0^{\tau} tf(t) dt \\
 &= \int_{\tau}^{\infty} \tau f(t) dt + \int_0^{\tau} tf(t) dt \\
 &= E(\min(\tau, T)),
 \end{aligned}$$

which can be interpreted as restricted mean survival time. And  $\mu$  can be estimated by

$$\hat{\mu} = \int_0^{\tau} \hat{S}(t) dt.$$

## References