

Lecture 5: Survival Analysis

Instructor: Yen-Chi Chen

Note: in this lecture, we will use the notations T_1, \dots, T_n as the response variable and all these random variables are positive. These random variables will be called event time or death time. They often refer to certain ‘time’ characteristics of each individual, e.g., the time that the individual is dead/gets a disease.

5.1 Survival Function

We assume that our data consists of IID random variables $T_1, \dots, T_n \sim F$. The **survival function** $S(t)$ of this population is defined as

$$S(t) = P(T_1 > t) = 1 - F(t).$$

Namely, it is just one minus the corresponding CDF. Although this definition is extremely simple and seems to be very trivial from the CDF, later we will see that it turns out to be an elegant tool of modeling and interpreting the data.

In medical research, the quantity T_i often refers to certain time characteristic of individual i . For instance, the variable T may refer to the age that the individual i passes away. Then the survival function $S(t)$ can be interpreted as *the chance that an individual is still alive after age t* . If $S(60) = 0.8$, it means that there are 80% of the individuals in the population who will still be alive at the age 60. Namely, $S(t)$ is the probability that an individual will survive past time t .

Here are some basic properties about $S(t)$:

- $S(0) = 1$ and $S(\infty) = 0$.
- $S(t)$ is a non-increasing function.

5.1.1 Continuous case

A quantity that is often used along with the survival function is the hazard function. The **hazard function** is

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T_1 \leq t + \Delta t | T_1 \geq t)}{\Delta t} = \frac{p(t)}{S(t)},$$

where $p(t) = \frac{d}{dt}F(t)$ is the PDF of random variable T_1 . Note that you can also write the hazard function as

$$h(t) = -\frac{\partial \log S(t)}{\partial t}.$$

Note that sometimes you may see the definition of hazard as $\lim_{\Delta t \rightarrow 0} \frac{P(t < T_1 \leq t + \Delta t | T_1 > t)}{\Delta t}$ that does not involve equality; both definitions are equivalent for continuous random variable.

How can we interpret the hazard function? The hazard function describes the ‘intensity of death’ at the time t given that the individual has already survived past time t .

There is another quantity that is also common in survival analysis, the cumulative hazard function. The **cumulative hazard function** is

$$H(t) = \int_0^t h(s) ds.$$

You can interpret $H(t)$ as the cumulative amount of hazard up to time t . The cumulative hazard function and survival function as linked as follows:

$$H(t) = -\log S(t), \quad S(t) = e^{-H(t)} = e^{-\int_0^t h(s) ds}. \quad (5.1)$$

Example 1. What is the survival function and hazard function of an exponential R.V.? Let $T_1 \sim \text{Exp}(\lambda)$. Then

$$p(t) = \lambda e^{-\lambda t}, \quad F(t) = 1 - e^{-\lambda t} \text{ for } t \geq 0$$

Thus,

$$S(t) = e^{-\lambda t}$$

and

$$h(t) = \lambda, \quad H(t) = \lambda t.$$

Namely, in an exponential distribution, the hazard function is a constant and the cumulative hazard is just a linear function of time.

Example 2 (Weibull distribution). The Weibull distribution is a distribution with two parameters, λ and k , and it is a distribution for positive random variable. Its PDF is

$$p(t) = \lambda k \cdot (\lambda t)^{k-1} \cdot e^{-(\lambda t)^k}, t \geq 0.$$

When $k = 1$, it reduces to the exponential distribution. Its CDF and survival function are

$$F(t) = 1 - e^{-(\lambda t)^k}, \quad S(t) = e^{-(\lambda t)^k}.$$

And the hazard function and cumulative hazard function are

$$h(t) = \lambda k \cdot (\lambda t)^{k-1}, \quad H(t) = (\lambda t)^k.$$

5.1.2 Discrete case

When the time-to-event variable T is discrete, the hazard function is not easy to work with since $h(t)$ is like a PDF.

Inspired by the PMF, for the case of discrete T , we define the **discrete hazard function** to be

$$\lambda(t) = P(T = t | T \geq t).$$

Suppose $T \in \{t_1, t_2, \dots, t_K\}$, the cumulative hazard for discrete T is

$$H(t) = \sum_{t_j \leq t} \lambda(t_j).$$

The recovery of survival function from discrete hazard is a bit different from continuous case. There is a famous formula for converting a discrete hazard function into a survival function:

$$S(t) = \prod_{t_j \leq t} (1 - \lambda(t_j)). \quad (5.2)$$

Here we show that equation (5.2) is indeed the survival function. Without loss of generality, assume that $T \in \{t_1, t_2, \dots, t_K\}$ such that $t_1 < t_2 < t_3 < \dots < t_K$. Since the survival function $S(t) = 1 - F(t)$, the survival function will be like the CDF that has multiple flat region and some drops (drops occur when the CDF has a jump, i.e., there is a probability mass). Therefore, we only need to focus on the function at each t_j .

For case $t = t_1$,

$$S(t_1) = P(T > t_1) = 1 - P(T = t_1) = 1 - \lambda(t_1)$$

so the formula works.

For case $t = t_2$, the product formula in equation (5.2) is

$$\begin{aligned} (1 - \lambda(t_1))(1 - \lambda(t_2)) &= [1 - P(T = t_1)] \times [1 - P(T = t_2 | T \geq t_2)] \\ &= [1 - P(T = t_1)] \times \left[1 - \frac{P(T = t_2)}{1 - P(T < t_2)}\right] \\ &= [1 - P(T = t_1)] \times \left[1 - \frac{P(T = t_2)}{1 - P(T = t_1)}\right] \\ &= 1 - P(T = t_1) - P(T = t_2) \\ &= P(T > t_2). \end{aligned}$$

Therefore, the formula works again.

You can then use induction to show that the formula works for any t_j . Therefore, equation (5.2) is indeed the formula for obtaining the survival function from discrete hazard.

5.2 Estimating the Survival Function

How do we estimate the survival function? There are four popular methods. The first method is a parametric approach. This method assumes a parametric model (e.g., exponential distribution) of the data and we estimate the parameter first then form the estimator of the survival function. A second approach is to compute the EDF first and then converted it to an estimator of the survival function. The third approach is a powerful nonparametric method called the Kaplan-Meier estimator and we will discuss it in the next section. Finally, there is another method call Nelson-Aalen estimator, which utilizes the cumulative hazard to estimate the survival function.

Parametric Approach. Assume that we model the distribution as an exponential distribution with unknown parameter λ . An estimator of λ is (you can check HW01 to see why this is an estimator)

$$\hat{\lambda} = \frac{1}{\bar{T}_n} = \frac{n}{\sum_{i=1}^n T_i}.$$

Then we estimate the survival function using

$$\hat{S}_1(t) = \hat{\lambda} e^{-\hat{\lambda}t} = \frac{e^{-\frac{t}{\bar{T}_n}}}{\bar{T}_n}, \quad t \geq 0.$$

EDF Approach. Recall that the EDF $\hat{F}(t)$ will be

$$\hat{F}(t) = \frac{1}{n} \sum_{i=1}^n I(T_i \leq t).$$

Then the survival function can be estimated by

$$\widehat{S}_2(t) = 1 - \widehat{F}(t) = \frac{1}{n} \sum_{i=1}^n I(T_i > t).$$

5.2.1 Kaplan-Meier estimator

Let $t_1 < t_2 < \dots < t_m$ be the time point where the observations T_1, \dots, T_n actually take values.

To see how the estimator is constructed, we do the following analysis. We partition the time axis into disjoint segments:

$$B_0 = [0, t_1), B_1 = [t_1, t_2), \dots, B_{m-1} = [t_{m-1}, t_m), B_m = [t_m, \infty).$$

Then we define

$$N_\ell = \text{number of individuals alive at (event happens after) the beginning of } B_\ell = \sum_{i=1}^n I(T_i \geq t_\ell)$$

and

$$D_\ell = \text{number of individuals die (event happens at) in } B_\ell = \sum_{i=1}^n I(T_i \in B_\ell).$$

Now we have converted T_1, \dots, T_n to $(N_0, D_0), \dots, (N_m, D_m)$. Formally, N_ℓ should be defined as the number of individuals *at risk* at the beginning of B_ℓ . Later we will explain what does the *at risk* means.

The **Kaplan-Meier (KM) estimator** estimates $S(t)$ using

$$\widehat{S}_{KM}(t) = \prod_{\ell: t_\ell \leq t} \left(1 - \frac{D_\ell}{N_\ell}\right).$$

What is the intuition of the KM estimator? We now consider t in different time segments and see if we can gain some intuitions. Recall that the survival function

$$S(t) = P(T > t) = \text{Probability of surviving past time } t.$$

For $t \in B_0 = [0, t_1)$, there is no event happens within this interval so $\widehat{S}_{KM}(t) = 1$.

For $t \in B_1 = [t_1, t_2)$, the survival function

$$S(t) = P(T > t) = P(\text{survives past time } t) = P(\text{survives in } [0, t_1) \text{ and in } [t_1, t)) = P(\text{survives in } B_0 \text{ and in } B_1).$$

Now recall that for two events A and B , $P(A \text{ and } B) = P(A)P(B|A)$. Thus,

$$S(t) = P(\text{survives in } B_0 \text{ and in } B_1) = P(\text{survives in } B_0)P(\text{survives in } B_1 | \text{survives in } B_0).$$

The probability $P(\text{survives in } B_1 | \text{survives in } B_0)$ can be estimated using

$$\widehat{P}(\text{survives in } B_1 | \text{survives in } B_0) = \frac{N_1 - D_1}{N_1} = 1 - \frac{D_1}{N_1}$$

and because no event occurs in B_0 , $P(\text{survives in } B_0) = 1$. Thus,

$$\widehat{S}_{KM}(t) = 1 \times \left(1 - \frac{D_1}{N_1}\right).$$

Now for the next time segment B_2 , we apply the same intuition. Namely, for $t \in B_2$,

$$S(t) = P(\text{survives in } B_0)P(\text{survives in } B_1|\text{survives in } B_0)P(\text{survives in } B_2|\text{survives in } B_1),$$

where we can estimate $P(\text{survives in } B_2|\text{survives in } B_1)$ via

$$\widehat{P}(\text{survives in } B_2|\text{survives in } B_1) = 1 - \frac{D_2}{N_2},$$

which leads to

$$\widehat{S}_{KM}(t) = 1 \times \left(1 - \frac{D_1}{N_1}\right) \times \left(1 - \frac{D_2}{N_2}\right).$$

For the other segments, we can apply the same procedure to obtain the estimator. This gives you the intuition of how the KM estimator is constructed.

Why does the Kaplan-Meier estimator work? A simple explanation is the formula in equation (5.2). The Kaplan-Meier estimator is essentially using an estimated discrete hazard function to construct a survival function estimator. Recall that a discrete hazard is

$$\lambda(t) = P(T = t|T \geq t) = \frac{P(T = t)}{P(T \geq t)}.$$

Therefore, a simple estimator of $\lambda(t)$ is

$$\widehat{\lambda}(t) = \frac{\widehat{P}(T = t)}{\widehat{P}(T \geq t)} = \frac{\text{Number of } \{T_i = t\}}{\text{Number of } \{T_i \geq t\}} = \frac{D(t)}{N(t)},$$

where $D(t)$ is the number of events at time $T = t$ and $N(t)$ is the number of individual at risk at time $T = t$. Namely, $D_j = D(t_j)$ and $N_j = N(t_j)$.

Note that when we observe every individual's event time (namely, there is no censoring – a mechanism we will discuss later), the KM estimator and the EDF approach are the same.

Example. Suppose we have a data such that the time observed is $T_i = 1, 3, 3, 5, 7, 8, 11, 11, 14, 15$. The following table summarizes the computation of the KM estimator.

| | | | | | | | | |
|-----------------------------------|------|------|------|------|------|------|------|------|
| Time (t_j) | 1 | 3 | 5 | 7 | 8 | 11 | 14 | 15 |
| Events (D_j) | 1 | 2 | 1 | 1 | 1 | 2 | 1 | 1 |
| At Risk (N_j) | 10 | 9 | 7 | 6 | 5 | 4 | 2 | 1 |
| Dis. Haz. (D_j/N_j) | 0.10 | 0.22 | 0.14 | 0.17 | 0.20 | 0.50 | 0.50 | 1.00 |
| Cum. Haz. $\widehat{H}(t)$ | 0.10 | 0.32 | 0.46 | 0.63 | 0.83 | 1.33 | 1.83 | 2.83 |
| KM Est. $\widehat{S}(t)$ | 0.90 | 0.70 | 0.60 | 0.50 | 0.40 | 0.20 | 0.10 | 0.00 |

5.2.2 Nelson-Aalen estimator

Nelson-Aalen (NA) estimator is another powerful estimator of the survival function. It not only estimates the survival function but also provides an estimate of the cumulative hazard. Actually, NA estimator first estimate the cumulative hazard function and then convert it into an estimate of the survival function using the relation $S(t) = e^{-H(t)}$. Here is an intuition about how this estimator is constructed.

Recall that the KM estimator uses

$$\widehat{S}_{KM}(t) = \prod_{\ell: t_\ell \leq t} \left(1 - \frac{D_\ell}{N_\ell}\right).$$

as an estimate of $S(t)$. When D_ℓ is much smaller than N_ℓ , we have

$$e^{-\frac{D_\ell}{N_\ell}} \approx 1 - \frac{D_\ell}{N_\ell}.$$

Therefore,

$$\begin{aligned} \widehat{H}_{KM}(t) &= -\log \widehat{S}_{KM}(t) \\ &= -\log \prod_{\ell:t_\ell \leq t} \left(1 - \frac{D_\ell}{N_\ell}\right) \\ &= -\sum_{\ell:t_\ell \leq t} \log \left(1 - \frac{D_\ell}{N_\ell}\right) \\ &\approx -\sum_{\ell:t_\ell \leq t} \log e^{-\frac{D_\ell}{N_\ell}} \\ &= \sum_{\ell:t_\ell \leq t} \frac{D_\ell}{N_\ell}. \end{aligned}$$

Using the above derivation, the NA estimator estimates the cumulative hazard function by

$$\widehat{H}_{NA}(t) = \sum_{\ell:t_\ell \leq t} \frac{D_\ell}{N_\ell}$$

and then estimate the survival function as

$$\widehat{S}_{NA}(t) = e^{-\widehat{H}_{NA}(t)} = e^{-\sum_{\ell:t_\ell \leq t} \frac{D_\ell}{N_\ell}} = \exp\left(-\sum_{\ell:t_\ell \leq t} \frac{D_\ell}{N_\ell}\right).$$

You can view the NA estimator as a combination of discrete and continuous hazard. We use sample discrete hazard function $\lambda(t)$ to estimate the cumulative hazard and then use the exponential formula from continuous hazard in equation (5.1) to obtain the survival function.

5.3 Censoring

However, in reality, our data may not be so nice. We may not be able to observe the actual event time T_i because of many complications. For instance, in a medical research, individuals may leave the study (called dropout) so we only observe their leaving time instead of the actual death time. The phenomena that we sometimes cannot observe the actual time but a ‘censoring time’ is called **censoring** in Statistics.

To model this process, we often need to introduce two other variables: Y and C . The T is the actual event time of interest and C is the censoring time that is competing with T and Y is the actual observing time.

In most cases, we will consider the **right-censoring** problem where the three variables are related by

$$Y = \min\{T, C\}.$$

We will assume that T and C are independent. Note that if what we observe is $Y = \max\{T, C\}$, this problem is called a left-censoring problem. Moreover, we not only observe Y , we also know if this Y comes from the event time or censoring time. Namely, we have one extra variable Δ such that $\Delta = I(T < C)$.

When we only observe $(Y_1, \Delta_1), \dots, (Y_n, \Delta_n)$ instead of T_1, \dots, T_n , how can we infer the survival function T_1 ? This is the central question to many biostatistical research.

Because we have several R.V.s now, we will add subscript to denote the functions associated to each random variable. Namely, F_T, S_T, h_T, H_T are the CDF, survival function, hazard function, and cumulative hazard function of random variable T and F_C, S_C, h_C, H_C are those of random variable C and F_Y, S_Y, h_Y, H_Y are those of random variable Y .

Here are some relations among these functions.

- $S_Y(t) = P(Y > t) = P(\min\{T, C\} > t) = P(T > t)P(C > t) = S_T(t)S_C(t)$.
Namely, the survival function of Y is the product of the other two survival functions.
- $F_Y(t) = 1 - (1 - F_T(t))(1 - F_C(t)) = F_T(t) + F_C(t) - F_T(t)F_C(t)$.
- $p_Y(t) = p_T(t) + p_C(t) - p_T(t)F_C(t) - p_C(t)F_T(t) = p_T(t)S_C(t) + p_C(t)S_T(t)$.
The PDF of Y is the sum of the weighted PDF of the other two and the weight is the survival function.
- $h_Y(t) = h_T(t) + h_C(t)$.
Namely, the hazard function of Y is the summation of the other two.
- $H_Y(t) = H_T(t) + H_C(t)$.
Similarly, the cumulative hazard is also the sum of the other two.

Note that Δ is just a Bernoulli random variable with probability being 1 as $P(T < C)$.

5.3.1 Estimating the Survival Function in Censoring

When there is censoring, the EDF approach no longer works. However, the KM and NA estimators are still valid. Essentially, the estimator is the same but we need to modify a little bit about N_ℓ and D_ℓ . As we have mentioned, formally, N_ℓ should be defined as

$$N_\ell = \text{number of individuals at risk at the beginning of } B_\ell.$$

What does the phrase *at risk* means? It refers to as being alive *and* not censored so it can be modified by replacing T_i with Y_i . Thus,

$$N_\ell = \sum_{i=1}^n I(Y_i \geq t_\ell).$$

For the quantity D_ℓ , it is still the number of events in the interval B_ℓ but we need to modify it by the number of *observed* events in the interval. Therefore,

$$D_\ell = \sum_{i=1}^n I(Y_i \in B_\ell, \Delta_i = 1).$$

Using these two modifications, the KM estimator and NA estimator are

$$\hat{S}_{KM}(t) = \prod_{\ell: t_\ell \leq t} \left(1 - \frac{D_\ell}{N_\ell}\right)$$

$$\hat{S}_{NA}(t) = \exp\left(-\sum_{\ell: t_\ell \leq t} \frac{D_\ell}{N_\ell}\right).$$

Note that parametric models may still be applicable during the censoring case and the estimator is often done using a maximum likelihood approach, which is beyond the scope of this course so we will not cover it here.

Example. Now consider the previous example but with censoring. The observed time is

$$Y_i = 1, 3, 3*, 5, 7*, 8, 11, 11, 14*, 15.$$

The star sign indicates the censoring event. Namely, $\Delta_3 = \Delta_5 = \Delta_9 = 0$ and other $\Delta_i = 1$. The following table summarizes the computation of the KM estimator in this case:

| | | | | | | | | |
|-----------------------------------|------|------|------|------|------|------|------|------|
| Time (t_j) | 1 | 3 | 5 | 7 | 8 | 11 | 14 | 15 |
| Events (D_j) | 1 | 1 | 1 | 0 | 1 | 2 | 0 | 1 |
| At Risk (N_j) | 10 | 9 | 7 | 6 | 5 | 4 | 2 | 1 |
| Dis. Haz. (D_j/N_j) | 0.10 | 0.11 | 0.14 | 0.00 | 0.20 | 0.50 | 0.00 | 1.00 |
| Cum. Haz. $\widehat{H}(t)$ | 0.10 | 0.21 | 0.35 | 0.35 | 0.55 | 1.05 | 1.05 | 2.05 |
| KM Est. $\widehat{S}(t)$ | 0.90 | 0.80 | 0.69 | 0.69 | 0.55 | 0.27 | 0.27 | 0.00 |

A feature is that when the unique time is censored (such as the case of $T = 7, 14$), the estimator is not updated since there is no observed event. However, the number of individual at risk is updated.

5.3.2 Why hazard can be estimated under censoring?

A key reason that Kaplan-Meier and Nelson-Aalen estimators work in censoring case is that *the hazard function can still be estimated even if we have censoring*.

To see how it works, we consider a discrete random time, i.e., T, Y, C are all discrete random variables. The new hazard estimator is using the ratio $\frac{D_\ell}{N_\ell}$, which can be clearly viewed as an estimator of $\frac{P(Y=t_\ell, \Delta=1)}{P(Y \geq t_\ell)}$. Namely,

$$\frac{D_\ell}{N_\ell} \approx \frac{P(Y = t_\ell, \Delta = 1)}{P(Y \geq t_\ell)}.$$

Using the fact that $Y = \min\{T, C\}$ and $\Delta = I(Y = T)$, the denominator

$$P(Y \geq t_\ell) = P(\min\{T, C\} \geq t_\ell) = P(T \geq t_\ell)P(C \geq t_\ell)$$

For the numerator, we have

$$P(Y = t_\ell, \Delta = 1) = P(T = t_\ell, \Delta = 1) = P(T = t_\ell, C \geq t_\ell) = P(T = t_\ell)P(C \geq t_\ell)$$

Therefore, the ratio

$$\frac{D_\ell}{N_\ell} \approx \frac{P(Y = t_\ell, \Delta = 1)}{P(Y \geq t_\ell)} = \frac{P(T = t_\ell)P(C \geq t_\ell)}{P(T \geq t_\ell)P(C \geq t_\ell)} = \frac{P(T = t_\ell)}{P(T \geq t_\ell)} = P(T = t_\ell | T \geq t_\ell) = \lambda_T(t_\ell),$$

which is the discrete hazard of T at t_ℓ , the time-to-event variable of interest, not the variable Y !

Since we know that Kaplan-Meier and Nelson-Aalen estimators are based on the hazard estimator, the ratio $\frac{D_\ell}{N_\ell}$ approximates the true hazard, so the two estimators are still applicable in the censoring case.

5.4 Greenwood's formula

For Kaplan-Meier estimator $\widehat{S}_{KM}(t)$, a common approach to assess its uncertainty is via the *Greenwood's formula*. Simply put, the variance of $\widehat{S}_{KM}(t)$ can be estimated via

$$\text{Var}(\widehat{S}_{KM}(t)) \approx \widehat{S}_{KM}^2(t) \cdot \sum_{t_j \leq t} \frac{D_j}{N_j(N_j - D_j)}. \quad (5.3)$$

This formula is obtained via a number of approximation and the so-called *delta method*. Here is a simple derivation of it.

First, recall that the Kaplan-Meier estimator can be written as

$$\widehat{S}_{KM}(t) = \prod_{t_j \leq t} \left(1 - \frac{D_j}{N_j}\right) = \prod_{t_j \leq t} (1 - \widehat{H}_j),$$

where $H_j = P(T = t_j | T \geq t_j)$ is the discrete hazard (consider no censoring case for simplicity). Therefore, the log of the KM estimator is

$$\log \widehat{S}_{KM}(t) = \sum_{t_j \leq t} \log(1 - \widehat{H}_j).$$

We first consider the variance of this logarithm and approximate it with

$$\text{Var}(\log \widehat{S}_{KM}(t)) \approx \sum_{t_j \leq t} \text{Var}(\log(1 - \widehat{H}_j)).$$

The reason why the above is an approximation is due to the fact that terms in the summation are dependent but the dependency is rather weak so we may approximate them as if they are independent. The above approximation, we now need to compute the variance $\text{Var}(\log(1 - \widehat{H}_j))$.

The estimated hazard \widehat{H}_j is a sample proportion (observed event divided by the total individual at risk) estimator of the population proportion/probability H_j , so

$$\sqrt{N_j}(\widehat{H}_j - H_j) \approx N(0, H_j(1 - H_j))$$

or equivalently,

$$\widehat{H}_j \approx N\left(H_j, \frac{H_j(1 - H_j)}{N_j}\right).$$

Then we utilize a method called Delta method. For a random variable $Z \approx N(c, \sigma^2)$ and a smooth function f , we have

$$\text{Var}(f(Z)) \approx \text{Var}[f(c) + (Z - c)f'(c)] = \text{Var}(Z)|f'(c)|^2.$$

With these result, we can approximate $\text{Var}(\log(1 - \widehat{H}_j))$ via

$$\text{Var}(\log(1 - \widehat{H}_j)) \approx \left(\frac{1}{1 - H_j}\right)^2 \frac{H_j(1 - H_j)}{N_j} = \frac{H_j}{N_j} \frac{1}{1 - H_j} = \frac{D_j}{N_j(N_j - D_j)}.$$

Thus,

$$\text{Var}(\log \widehat{S}_{KM}(t)) \approx \sum_{t_j \leq t} \text{Var}(\log(1 - \widehat{H}_j)) \approx \sum_{t_j \leq t} \frac{D_j}{N_j(N_j - D_j)}.$$

Finally, we use the delta method again on

$$W = \log \widehat{S}_{KM}(t) \approx N \left(\log S(t), \text{Var}(\log \widehat{S}_{KM}(t)) \right)$$

so that

$$\begin{aligned} \text{Var} \left(\widehat{S}_{KM}(t) \right) &= \text{Var}(e^W) \\ &= S^2(t) \cdot \text{Var}(W) \\ &= S^2(t) \cdot \sum_{t_j \leq t} \frac{D_j}{N_j(N_j - D_j)} \\ &\approx \widehat{S}_{KM}^2(t) \cdot \sum_{t_j \leq t} \frac{D_j}{N_j(N_j - D_j)}, \end{aligned}$$

where we replace the true survival function $S(t)$ with the Kaplan-Meier estimator in the last approximation. The final result is exactly the Greenwood formula.

5.5 Cox Model (optional)

In reality, we often not only observe the event time for an individual but also have access to other covariates of this individual. We often are interested in understanding how these covariates affect the survival function of the event.

For instance, in a cancer study, we may have each individual's age when they got cancer (the event time T) and this individual's gender, BMI, smoking habit, and education level. The other variables are the covariates in this study. Health scientists are often interested in how these covariates change the survival function. Let X denotes the covariates. A parameter of interest will be the survival function of T given X . Namely, it is the conditional survival function

$$S(t|x) = P(T > t | X = x).$$

For instance, we may be interested in

$$S(\text{Age} = t | (\text{gender, BMI, smokinghabit, educationlevel}) = (\text{male, 20, neversmoke, college})).$$

We can then define the conditional hazard function and conditional cumulative hazard function as

$$h(t|x) = -\frac{\partial \log S(t|x)}{\partial t}, \quad H(t|x) = -\log S(t|x).$$

The **Cox (proportional hazard) model** is one of the most popular model combining the covariates and the survival function. It starts with modeling the hazard function $h(t|X = x)$:

$$h(t|X = x) = h_0(t) \exp(x^T \beta),$$

where β is the vector of coefficients of each covariate. The function $h_0(t)$ is called the baseline hazard function. Namely, the Cox model assumes that the covariates have a linear multiplication effect on the hazard function and the effect stays the same across time.

This implies the conditional hazard function being

$$H(t|x) = \exp(x^T \beta) \int_0^t h_0(s) ds = \exp(x^T \beta) H_0(t),$$

where $H_0(t)$ is the baseline cumulative hazard function. This further yields the conditional survival function

$$S(t|x) = \exp(-H(t|x)) = \exp(-\exp(x^T\beta)H_0(t)) = \exp(-H_0(t))^{\exp(x^T\beta)} = S_0(t)^{\exp(x^T\beta)},$$

where $S_0(t)$ is called the baseline survival function.

Why it is called a *proportional* hazard model? Here is an intuition about it. Consider two individuals with different covariates that one has $X = x_1$ and the other has $X = x_2$. The ratio of their hazard function

$$\frac{h(t|x_1)}{h(t|x_2)} = \frac{h_0(t)\exp(x_1^T\beta)}{h_0(t)\exp(x_2^T\beta)} = \frac{\exp(x_1^T\beta)}{\exp(x_2^T\beta)} = \exp((x_1 - x_2)^T\beta)$$

is a constant over time. Namely,

$$h(t|x_1) = \exp((x_1 - x_2)^T\beta) \times h(t|x_2) \propto h(t|x_2) \quad \forall t \geq 0.$$

Thus, their hazard is always proportional to each other regardless of the value of time t .

Estimation of the parameter β is often done by maximizing the *partial likelihood function*:

$$\hat{L}_n(\beta) = \prod_{i=1}^n L_i(\beta),$$

where

$$L_i(\beta) = \frac{h(T_i|X_i)}{\sum_{j:T_j \geq T_i} h(T_j|X_j)} = \frac{\exp(X_i^T\beta)}{\sum_{j:T_j \geq T_i} \exp(X_j^T\beta)}.$$

Namely, our estimator

$$\hat{\beta}_n = \operatorname{argmax}_{\beta} \hat{L}_n(\beta).$$

This estimator turns out to be an unbiased estimator and has variance shrinking at rate $O(n^{-1})$ and has asymptotic normality under suitable condition. An interesting fact is that *we do not need to know the baseline hazard function $h_0(t)$ to estimate β !* (estimating $h_0(t)$ is not easy and the convergence rate is often slow; we will discuss a similar pattern in density estimation) The property that we can estimate parameter of interest without estimating the entire model is related to the topic *semi-parametric model*¹.

¹https://en.wikipedia.org/wiki/Semiparametric_model