

**SUMMER SCHOOL IN STATISTICS 2015
HANOÏ VIETNAM**



Survival Analysis

Part I

Jean-François Dupuy
Jean-Francois.Dupuy@ina-rennes.fr
and
Agnès Lagnoux
lagnoux@univ-tlse2.fr

Contents

- 1 Introduction** **5**
- 1.1 Why specific methods? 5
- 1.2 Lifetime models 6
 - 1.2.1 Lifetime and related functions 6
 - 1.2.2 Usual lifetime distributions 9
- 1.3 Censored data 10
 - 1.3.1 Censoring and truncation 10
 - 1.3.2 Some examples of censoring 11
 - 1.3.3 Right censoring 12
- 1.4 Likelihood of a censored sample 14
 - 1.4.1 Expression of the likelihood 14
 - 1.4.2 Estimation in parametric models 14
- 1.5 The point process $N(t)$ 15
- 1.6 Bibliography 16

- 2 Non parametric estimations** **17**
- 2.1 Without censoring 17
- 2.2 Kaplan-Meier estimator of the survival function 18
 - 2.2.1 Construction of the estimator 18
- 2.3 Breslow estimator of the cumulative hazard rate function 22
- 2.4 Nelson-Aalen estimator of the cumulative hazard function 24
- 2.5 Harrington and Fleming estimator of the survival function 25
- 2.6 Comments 25

- 3 Comparison of the survival functions of two (or more) groups** **27**
- 3.1 The weighted logrank tests 27
 - 3.1.1 Comparison of two groups 27
 - 3.1.2 Generalization to the comparison of K groups 30
- 3.2 Comparison with adjustment: stratified logrank test 30

Chapter 1

Introduction

In this lecture, our goal is to modelize times elapsed before the occurrence of an event. Historically, the studied events corresponded to deaths (that is why the term “lifetime” is currently used) but any event can be studied in this framework.

Numerous and various domains are concerned:

- medicine and health: survival time of the patient, remission time of a disease...
- reliability: component lifetimes,
- economy: loss of employment,
- psychology: learning time of some specific skill,
- sociology: succession life events (wedding, birth, divorce,...).

The statistical community has been particularly active on this research domain during the 20th century:

- Kaplan and Meier (1958): non parametric estimator of the survival function,
- Peto and Peto (1972): logrank test to compare the survival of two groups,
- Cox (1972): semi-parametric model that now is the most popular in survival data analysis,
- Gill (1980): introduction to the theory of processes ad martingales.

1.1 Why specific methods?

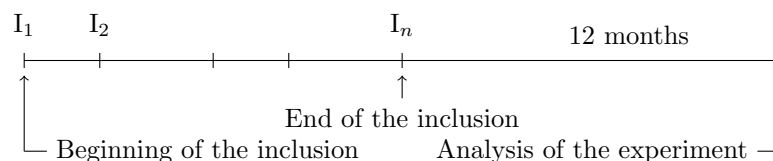
Example 1: We are interested in the survival of a group of mice after the injection of cancer cells; all the mice receive the injection at the same moment.

We can expect the death of all the mice and the exact lifetimes are observed and thus the standard methods apply.

Example 2: Therapeutic trial (clinical research)

We want to test the efficiency of a vaccine against the hepatitis B.

- 2 randomized groups: treatment/placebo.
- Criterion: arisen a hepatitis B during the year that follows the injection.
- Various injection dates for the different patients.



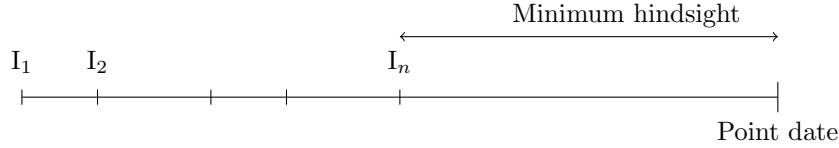
We observed 6 hepatitis on the 184 patients that received the vaccine but 7 patients has not been controlled during an entire year.

How can we estimate the rate of hepatitis after 12 months?

$$\frac{6}{184} ? \quad \frac{6}{184-7} ?$$

Example 3: Assume that we are interested in the survival after surgery of men operated for a bronchial cancer.

It is not realistic to wait for the death of all the subjects. Thus we fix a stopping time called **point date**.



Conclusion: Examples 2 and 3 deal with individuals for whom the event time is not known; the only thing that we know is that the event has occurred in a certain interval: this notion is called **censoring**. Only dealing with the non censored observations led to a loss of information which can be not insignificant.

↔ One has to introduce specific methods to take into account the censored data.

1.2 Lifetime models

1.2.1 Lifetime and related functions

Let X be a non negative random variable having a known continuous distribution function. Assume that X corresponds to the delay between a fixed original time and the occurrence of an event of interest.

Example : X is a survival variable and the event under concern corresponds to the death.

The usual functions such that the cumulative distribution function (c.d.f) F and the density probability function (p.d.f.) f are not used in this research domain. One prefers working with more adapted functions that are easier to interpret:

- The **survival (or reliability) function** is defined as the probability to survive until time t :

$$S(t) = 1 - F(t) = \mathbb{P}(X > t)$$

The survival function is sometimes denoted $\bar{F}(\cdot)$.

The mathematical expectation of X can be written in terms of the survival function:

$$\mathbb{E}[X] = \int_0^{+\infty} x dF(x) = \int_0^{+\infty} S(x) dx.$$

Indeed, by an integration by parts,

$$\begin{aligned} \mathbb{E}[X] &= \int_0^{+\infty} x dF(x) \\ &= [-S(x)x]_0^{+\infty} + \int_0^{+\infty} S(x) dx \\ &= \int_0^{+\infty} S(x) dx. \end{aligned}$$

- The **hazard rate (or risk) function** is given by

$$h(t) = \lim_{dt \rightarrow 0} \frac{\mathbb{P}(t \leq X < t + dt \mid X \geq t)}{dt}$$

for $t \geq 0$. $h(t)dt$ represents the probability to die during the interval t and $t + dt$ for an individual conditionally that he is still alive at time t ; in other words, the hazard rate at point t is the instantaneous probability of failure (or death) at time t given that failure (or death) has not occurred before. The hazard rate function may have different shapes: the most famous is called the **bathup curve**.

The function $h(\cdot)$ satisfies the following relations:

$$h(t) = \frac{f(t)}{S(t)} = -\frac{S'(t)}{S(t)} = -\frac{d}{dt} \ln S(t).$$

As a consequence, one has

$$S(t) = \exp \left\{ -\int_0^t h(s) ds \right\},$$

$$f(t) = h(t) \exp \left\{ -\int_0^t h(s) ds \right\}.$$

Examples of hazard rate functions:

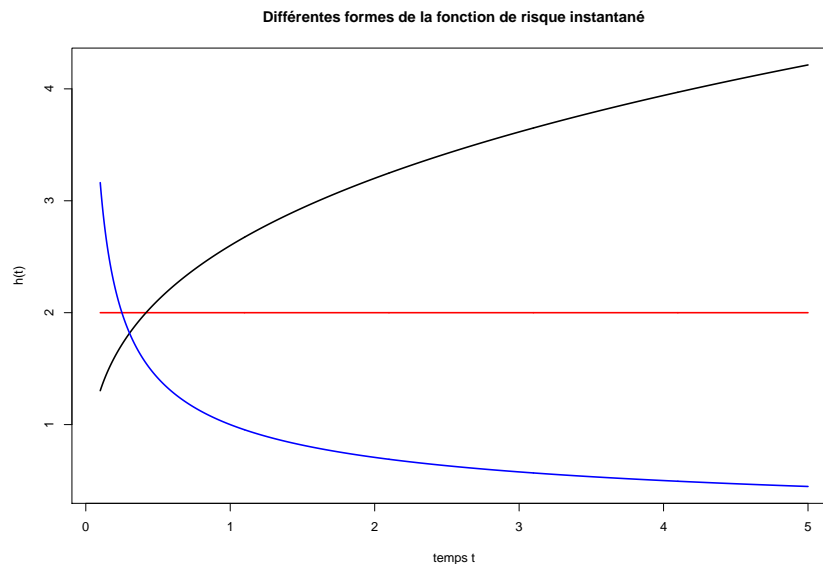
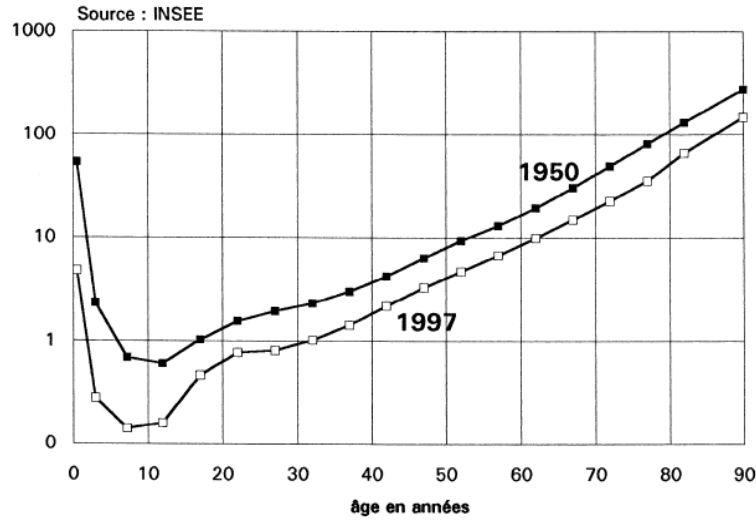
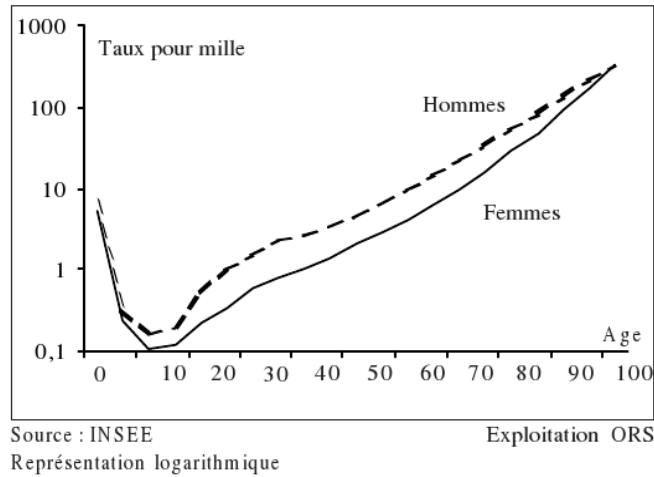


Figure 10 - TAUX DE MORTALITE PAR AGE EN FRANCE METROPOLITAINE (pour 1000 personnes) EN 1950 ET EN 1997



Taux de mortalité en Ile-de-France par sexe et par âge en 1993-95



- The **cumulative hazard rate function** is given by

$$H(t) = \int_0^t h(s)ds, \quad \text{for all } t.$$

One has the following relationships:

$$S(t) = \exp \{-H(t)\},$$

$$f(t) = h(t)S(t) = h(t) \exp \left\{ - \int_0^t h(s)ds \right\}.$$

Particular case: If the random variable X is discrete (i.e. valued in a countable set $\{x_1, x_2, \dots, x_n, \dots\}$ with $x_1 < x_2 < \dots < x_n, \dots$), we have

$$F(t) = \sum_{i:x_i \leq t} p_i$$

where $p_i = \mathbb{P}(X = x_i)$. The hazard rate function is given by

$$h(x_i) = \mathbb{P}(X = x_i | X \geq x_i) = \frac{p_i}{S(x_{i-1})}$$

and the cumulative rate function by

$$H(t) = \sum_{i:x_i \leq t} h(x_i).$$

Finally,

$$S(t) = \prod_{i:x_i \leq t} [1 - h(x_i)].$$

- The **residual life at time t** , denoted by τ_t , is the random variable with distribution:

$$\mathbb{P}(\tau_t > s) = \mathbb{P}(X - t > s | X \geq t) = \frac{S(t+s)}{S(t)}.$$

- The **mean residual function** is defined for $t \geq 0$ by

$$m(t) = \mathbb{E}[\tau_t] = \mathbb{E}[X - t | X \geq t] = \frac{\int_t^{+\infty} S(s) ds}{S(t)}.$$

1.2.2 Usual lifetime distributions

Let X be a non negative random variable. We are interested in random variables whose support is \mathbb{R}^+ . As a consequence, the Gaussian distribution is no longer the reference distribution.

- The **exponential** $\mathcal{E}(\lambda)$ where $\lambda > 0$. One has, for any $t \geq 0$,

$$\begin{aligned} S(t) &= e^{-\lambda t} \\ f(t) &= \lambda e^{-\lambda t} \mathbb{1}_{t \geq 0} \\ h(t) &= \lambda. \end{aligned}$$

The risk is then a constant function with respect to the time: this distribution is currently used in models without aging. One has $\mathbb{E}[X] = \frac{1}{\lambda}$ and $\text{Var}(X) = \frac{1}{\lambda^2}$ and

$$\forall x, y > 0, \quad \mathbb{P}(X > x + y | X > y) = \mathbb{P}(X > x).$$

The exponential is said to have the memoryless property: the distribution of the survival times for times greater than y is not affected by the knowledge that the individual has survived until y .

- The **Weibull distribution** $\mathcal{W}(\alpha, \lambda)$, where $\lambda > 0$ and $\alpha > 0$, is a generalization of the exponential distribution (that is a particular case when $\alpha = 1$). It allows to have increasing or decreasing risks. The parameter λ is the scale parameter and α is the shape parameter.

Here one has, for any $t \geq 0$,

$$\begin{aligned} S(t) &= e^{-(\lambda t)^\alpha}, \\ f(t) &= \alpha \lambda^\alpha t^{\alpha-1} e^{-(\lambda t)^\alpha} \mathbb{1}_{t > 0}, \\ h(t) &= \alpha \lambda^\alpha t^{\alpha-1}. \end{aligned}$$

The risk is then a power of the time.

- when $0 < \alpha < 1$: the risk decreases from $+\infty$ to 0;
- when $\alpha = 1$: the risk is constant ($\mathcal{W}(1, \lambda) = \mathcal{E}(\lambda)$);
- when $\alpha \geq 1$: the risk increases from 0 to $+\infty$.

Property: If $X \sim \mathcal{W}(\alpha, \lambda)$ then $X^\alpha \sim \mathcal{E}(\lambda^\alpha)$.

Warning: in the software R, the scale parameter is $\frac{1}{\lambda}$.

- The **Gamma distribution** $\gamma(a, \lambda)$ where $\lambda > 0$ and $a > 0$:

$$f(t) = \frac{\lambda^a}{\Gamma(a)} t^{a-1} e^{-\lambda t} \mathbb{1}_{t>0}$$

where $\Gamma(a) = \int_0^{+\infty} x^{a-1} e^{-x} dx$.

The parameter λ is the scale parameter and α is the shape parameter.

Reminders on the Gamma function Γ :

$$\begin{aligned} \Gamma(1) &= 1, \\ \Gamma\left(\frac{1}{2}\right) &= \sqrt{\pi}, \\ \forall a > 0, \Gamma(a+1) &= a\Gamma(a), \\ \forall n \in \mathbb{N}^*, \Gamma(n) &= (n-1)!. \end{aligned}$$

In this context, there is no explicit expression for S nor h .

One has: $\mathbb{E}[X] = \frac{a}{\lambda}$ and $\text{Var}(X) = \frac{a}{\lambda^2}$.

- when $0 < a < 1$: the risk decreases from $+\infty$ to $\frac{1}{\lambda}$;
- when $a = 1$: the risk is constant ($\gamma(\lambda, 1) = \mathcal{E}(\lambda)$);
- when $a \geq 1$: the risk increases from 0 to λ .

There exists other distributions that lead to monotone risks. To get non monotone risk functions, one can use the following distributions:

- the log-normal distribution,
- the log-logistic distribution,
- the inverse Gaussian distribution,
- the generalized Weibull distribution.

1.3 Censored data

1.3.1 Censoring and truncation

1. The data are often censored, that means that the event under concern may not be observed and one only knows that the delay of interest lies in a given time interval. There exists 3 different types of **censoring**:

- right censoring (the most frequent): the event occurs after a given date C : $X > C$,

- left censoring: the event occurs before a given date C : $X < C$,
 - censoring by interval: the event occurs between two given dates C_1 and C_2 : $C_1 < X < C_2$.
2. The **truncation** is a condition that hide some individuals in such a way that the statistician is not aware of their existence: the individuals that not satisfy the condition are not included in the study.
- Left truncation: the condition $X > Y$ guarantees that the individual takes part to the study where Y corresponds to the delay of occurrence of another event.

Example: one is interested in the survival of individuals that have contracted a given disease and the sample is taken in an retirement home. The individuals prematurely dead (before being retired) are not taken into account. Here, the survival X (that is the age at death) is left truncated by Y that represents the age at the entry of the retirement home.

- Right truncation: only the individuals for whom the event under concern has occurred are included in the study.

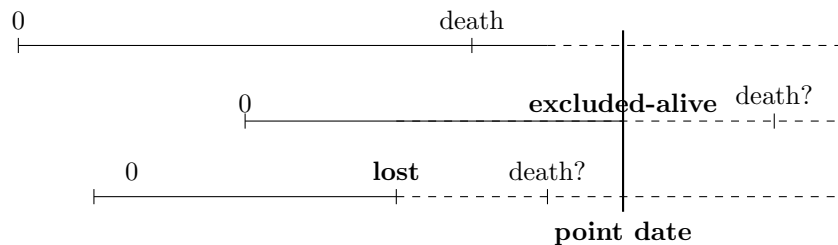
Example: Considering AIDS, one is interested in the distribution of the delay to develop the disease after the inclusion of the HIV. In this study, only the patients that have developed the HIV take part to the study: the healthy individuals that carry HIV are not known by the clinician.

Warning: Beware of not making the confusion between truncation and censoring!

In medicine, a right censoring and a left truncation appear in the most famous models.

1.3.2 Some examples of censoring

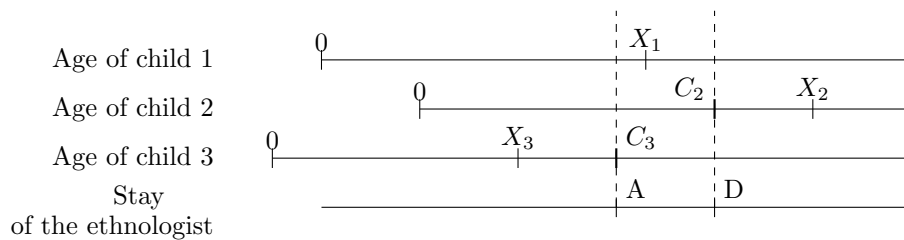
Example of right censoring:



We distinguish two different types of patients:

- the excluded-alive: patients still alive at the point date,
- the lost individuals: patients whose state is unknown at the point date.

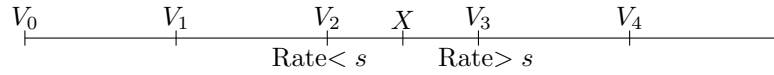
Example of left censoring: the ethnologist



- X_i is the age of child i when is able to do a given task,

- A and D are the arrival and departure times of the ethnologist,
- C_2 : age of child 2 when the ethnologist leaved,
- C_3 : age of child 3 at the arrival of the ethnologist.
- X_1 is not censored,
- X_2 is right censored by C_2 : $X_2 > C_2$,
- X_3 is left censored by C_3 : $X_3 < C_3$.

Examples of censoring by interval: we consider a group of patients that are controlled frequently during a clinical survey. Let X be the delay necessary for a biological variable to reach some level s . One only knows that X lives in an interval between two successive controls: here, $V_2 < X < V_3$.



The censoring by interval appears also in reliability if periodical inspections are lead to check the good functioning of the machines.

1.3.3 Right censoring

Since the right censoring is the most encountered in practice, we study it in details.

1. The censoring of type I: fixed

The delay X_i of the i -th individual is observed if and only if $X_i \leq c$ where c is a **fixed duration**. Otherwise, one only knows that $X_i > c$.

In other words, the censoring time c is known and fixed. One only observes the random variables $(T_i, \delta_i)_{i=1 \dots n}$ defined by

$$\begin{cases} T_i = \min(X_i, c) \\ \delta_i = \mathbb{1}_{\{X_i \leq c\}} \end{cases}$$

for $i = 1, \dots, n$.

2. The progressive censoring of type I

The censoring times c_i for $i = 1, \dots, n$ are known and fixed. One only observes the random variables $(T_i, \delta_i)_{i=1 \dots n}$ defined by

$$\begin{cases} T_i = \min(X_i, c_i) \\ \delta_i = \mathbb{1}_{\{X_i \leq c_i\}} \end{cases}$$

for $i = 1, \dots, n$.

Example: all the patients enter the study the same day. There are censored (excluded-alive) at the point date.

3. The censoring of type II : waiting (commonly used in reliability)

One considers the survival time X_i of n individuals until r events occur and then stops the study: the aim is to save time and money! The censoring occurs at $X_{(r)}$ the r -th order statistics. In other words, the censoring time is given by the time of the r -th failure observed in the sample. One only observes the random variables $(T_i, \delta_i)_{i=1 \dots n}$ defined by

$$\begin{cases} T_i = \min(X_i, X_{(r)}) \\ \delta_i = \mathbb{1}_{\{X_i \leq X_{(r)}\}} \end{cases}$$

for $i = 1, \dots, n$ and $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(r)}$ where are the r first order statistics.

4. The **censoring of type III**: random

To each patient, we associate a survival time X_i and a censoring delay C_i . There is a tradeoff between these two delays and one only observes the smaller denoted T_i :

- the observed delay is $T_i = \min(X_i, C_i)$;
- the event indicator is $\delta_i = \mathbb{1}_{X_i \leq C_i}$.

As a consequence, the sample consists in pairs (T_i, δ_i) , $i = 1, \dots, n$ and one only observes the random variables $(T_i, \delta_i)_{i=1 \dots n}$ defined by

$$\begin{cases} T_i = \min(X_i, C_i) \\ \delta_i = \mathbb{1}_{\{X_i \leq C_i\}} \end{cases}$$

for $i = 1, \dots, n$ and where C_1, \dots, C_n are random variables with c.d.f. G_1, \dots, G_n respectively. The random variables X_1, \dots, X_n and C_1, \dots, C_n are generally assumed to be independent. In this case, the c.d.f. of the random variable T_i is given by $S_{T_i} = S_X S_{C_i}$ for $i = 1, \dots, n$.

In the sequel, we consider a random censoring and we assume that X and C are two independent random variables. This assumption may be not satisfied in some practical examples: e.g. the treatment is stopped if some secondary effects appear or by tiredness of the failure in the studies on the infertility.

Vocabulary from medicine

We denote X the delay of occurrence of the event under concern. To define correctly this delay, one needs an origin and a date of end.

- The **origin date** (OD) corresponds to the original time of the individual. This is the principal source of the randomness of the censoring.

Example: date of birth, date of first symptoms of a disease, date of the diagnosis by the doctor, date of surgery, date of inclusion in a clinical experiment.

- The **date of last news** (DLN) corresponds to the date of death if the patient died. Otherwise it corresponds to the latest date of inspection, control, monitoring...
- The **monitoring duration** is the delay between the origin and the date of latest news.
- The **point date** (PD) corresponds to the date (common to all individuals) when we decide to check the state of the patients.

Remark: one does not take into account information posterior to the point date; otherwise a bias is introduced. Usually, one has recent news of non representative patients sub-groups (e.g. information on patients that relapse and no news of the ones that are in remission).

- The **participation time** corresponds to the “survival time” reported at the end of the analysis that is
 - if $DDN \leq PD$: the time of participation is the monitoring duration.
 - * if the patient is dead at the latest news, the time of participation is not censored.
 - * if the patient is still alive at the latest news, he is lost and his time of participation is censored.
 - if $DDN > PD$: the time of participation is censored and corresponds to the delay between the origin and the date point. The patient is then considered as an excluded-alive.
- The **hindsight** of an individual is the delay between the origin and the date point.

1.4 Likelihood of a censored sample

1.4.1 Expression of the likelihood

Let X be a lifetime random variable whose density and survival functions are given by $f(t)$ and $S(t)$. Let C be a censoring whose density and survival functions are given by $g(t)$ and $G(t)$. We assume that X and C are independent (we then deal with a model of censoring of type III).

As said before, the sample consists in n pairs $(T_i, \delta_i)_{1 \leq i \leq n}$ where $T_i = \min(X_i, C_i)$ corresponds to the observed delay and $\delta_i = \mathbb{1}_{X_i \leq C_i}$ is the event indicator for the patient i . The distribution of the pair (T_i, δ_i) is given by:

$$f(t, \delta) = [f(t)G(t)]^\delta [g(t)S(t)]^{1-\delta}.$$

The likelihood of the sample becomes

$$L_n((T_1, \delta_1), \dots, (T_n, \delta_n)) = \prod_{i=1}^n f(T_i, \delta_i).$$

If the distribution of the censoring does not depend on the parameters of the survival, the censoring is non informative and the “useful” part of the likelihood is simply:

$$\prod_{i=1}^n f(T_i)^{\delta_i} S(T_i)^{1-\delta_i}.$$

1.4.2 Estimation in parametric models

Let us consider a parametric model

$$\mathcal{F} = \{\mathbb{P}_\theta, \theta \in \Theta\} \quad \text{or} \quad \mathcal{F} = \{f_\theta, \theta \in \Theta\} \quad \text{or} \quad \mathcal{F} = \{F_\theta, \theta \in \Theta\}$$

for the lifetime X .

The question is how to estimate the (possibly multidimensional) parameter θ ? Many different methods are available...

- Without censoring, you already know how to do with a complete sample.
- With censoring, it is more difficult. Under the assumption of independence between lifetime and censoring, the likelihood is given by

$$L_n((T_1, \delta_1), \dots, (T_n, \delta_n); \theta) = \prod_{i=1}^n [f_\theta(T_i)(1 - G(T_i))]^{\delta_i} [S_\theta(T_i)g(T_i)]^{1-\delta_i},$$

where g (respectively G) still represents the p.d.f. (resp. c.d.f.) of the censoring C .

If the distribution of the censoring time does not depend on the parameter of interest θ , one can consider the “useful” likelihood given by

$$L_n((T_1, \delta_1), \dots, (T_n, \delta_n); \theta) = \prod_{i=1}^n [f_\theta(T_i)]^{\delta_i} [S_\theta(T_i)]^{1-\delta_i}.$$

The maximum likelihood estimator (MLE) is the value of θ that maximizes the likelihood.

1.5 The point process $N(t)$

Instead of studying the lifetime X , we may define the point process $N(t)$ that is equal to 0 while the event has yet occurred and is equal to 1 after:

$$N(t) = \mathbb{1}_{\{X \leq t\}}.$$

The process $N(t)$ jumps by 1 in $t = x$ when $X = x$. We denote $dN(t)$ the process variation on the interval $[t, t + dt[$:

$$\begin{aligned} dN(t) &= N((t + dt)^-) - N(t^-) \\ \lim_{dt \rightarrow 0} \frac{\mathbb{P}(dN(t) = 1/N(t^-) = 0)}{dt} &= h(t) \\ \mathbb{P}(dN(t) = 1/N(t^-) = 1) &= 0. \end{aligned}$$

Then one has

$$\lim_{dt \rightarrow 0} \frac{P(dN(t) = 1/N(t^-))}{dt} = h(t)\mathbb{1}_{X \geq t}.$$

We denote $\lambda(t) = h(t)\mathbb{1}_{X \geq t}$: λ is the intensity of the counting process N . The cumulative intensity is the function Λ defined by

$$\Lambda(t) = \int_0^t \lambda(u)du = \int_0^t \lambda(u)\mathbb{1}_{u \leq X} du = H(t \wedge X).$$

Example: For a lifetime X exponentially distributed with parameter $\theta > 0$,

$$\begin{aligned} f_X(t) &= \theta e^{-\theta t} \mathbb{1}_{t > 0}, \quad S_X(t) = e^{-\theta t}, \quad t > 0 \\ h(t) &= \theta, \quad H(t) = \theta t, \\ \lambda(t) &= \theta \mathbb{1}_{X \geq t}, \quad \Lambda(t) = \theta(t \wedge X). \end{aligned}$$

The difference between $N(t)$ and $\Lambda(t)$ is a martingale: $M(t) = N(t) - \Lambda(t)$.

$$\mathbb{E}[dM(t)/\mathcal{F}_{t-}] = \mathbb{E}[dN(t) - \lambda(t)dt/\mathcal{F}_{t-}] = \mathbb{E}[dN(t)/\mathcal{F}_{t-}] - \lambda(t)dt = 0$$

where $\mathcal{F}_{t-} = \sigma(N(u), u < t)$.

Point processes for censored data:

In the previous example, we add a right censoring and one observes $T = \min(X, C)$ instead of X . Let $\delta = \mathbb{1}_{\{X \leq C\}}$. The indicator of "presence at risk of a subject" becomes:

$$Y(t) = \mathbb{1}_{\{t \leq X \wedge C\}}.$$

If there is n patients, we denote for $i = 1, \dots, n$,

$$\begin{aligned} Y_i(t) &= \mathbb{1}_{T_i \geq t} \\ N_i(t) &= \mathbb{1}_{T_i \leq t, \delta_i = 1}. \end{aligned}$$

We denote $\bar{Y} = \sum_{i=1}^n Y_i(t)$ and $\bar{N}(t) = \sum_{i=1}^n N_i(t)$.

Moreover, if there is a left truncation (the subject is seen in the study only if $X > U$), we get

$$Y(t) = \mathbb{1}_{U \leq t \leq X \wedge C}.$$

1.6 Bibliography

Andersen, Per Kragh, Borgan, Ørnulf, Gill, Richard D. and Keiding, Niels. *Statistical models based on counting processes*. Springer Series in Statistics. Springer-Verlag, New York, 1993.

Cox, D. R. and Oakes, D. *Analysis of survival data*. Monographs on Statistics and Applied Probability. Chapman & Hall, London, 1984.

Fleming, Thomas R. and Harrington, David P. *Counting processes and survival analysis*. Wiley Series in Probability and Mathematical Statistics: Applied Probability and Statistics. John Wiley & Sons, Inc., New York, 1991.

Kalbfleisch, John D. and Prentice, Ross L. *The Statistical Analysis of Failure Time Data*. Second Edition Wiley Series in Probability and Statistics, 2011.

Klein, John P., Moeschberger, Melvin L. *Survival Analysis. Techniques for Censored and Truncated Data*. Springer Science & Business Media, 2003.

Therneau, Terry M. and Grambsch, Patricia M. *Modeling survival data: extending the Cox model*. Statistics for Biology and Health. Springer-Verlag, New York, 2000.

Chapter 2

Non parametric estimations

2.1 Without censoring

Assume that we observe a sample X_1, X_2, \dots, X_n of the lifetime X . The empirical cumulative distribution function

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

is an estimator of the c.d.f. $F(\cdot)$. As a consequence, a natural estimator of $S(\cdot) = 1 - F(\cdot)$ is given through the empirical cumulative distribution function \hat{F}_n :

$$\hat{S}_n(x) = 1 - \hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i > x}.$$

It is equivalent to estimate the hazard rate $\lambda(\cdot)$ by

$$\begin{aligned} \lambda_n(X_{(i)}) &= \frac{1}{n - i + 1} \quad \text{for } i = 1, \dots, n \\ \lambda_n(x) &= 0 \quad \text{for all } x \text{ that is not an observation} \end{aligned}$$

where the sample has been ordered: $X_{(1)} \leq \dots \leq X_{(n)}$.

The cumulative hazard rate function can be estimated by

$$\hat{\Lambda}_n(x) = \frac{1}{n} \sum_{i:x_{(i)} \leq x} \frac{1}{n - i + 1}.$$

Theorem 2.1.1 (Glivenko-Cantelli Theorem) *The empirical c.d.f. is consistent: as n goes to infinity*

$$\sup_{x \geq 0} \left| \hat{F}_n(x) - F(x) \right| \xrightarrow{a.s.} 0$$

Theorem 2.1.2 (Donsker Theorem) *The empirical cdf is weakly convergent:*

$$\sqrt{n} \left(\hat{F}_n(\cdot) - F(\cdot) \right) \xrightarrow{\mathcal{L}} B(\cdot)$$

where $B(\cdot)$ is a centered Gaussian process with covariance function

$$\text{Cov}(B(s), B(t)) = F(s \wedge t) - F(t)F(s).$$

Remark 2.1.3 $B(\cdot)$ is called the **Brownian bridge** and is also defined by $B(t) = W(t) - tW(1)$ for any $0 \leq t \leq 1$ and $W(\cdot)$ is a standard Brownian motion.

Now one can deduce a confidence band or a pointwise confidence interval for $F(\cdot)$. For example,

$$\left[\hat{F}_n(x) - z_{1-\alpha/2} \sqrt{\frac{\hat{F}_n(x) \hat{S}_n(x)}{n}}, \hat{F}_n(x) + z_{1-\alpha/2} \sqrt{\frac{\hat{F}_n(x) \hat{S}_n(x)}{n}} \right]$$

is an asymptotic $(1 - \alpha)$ confidence interval for $F(x)$ where z_α is the α -quantile of the standard Gaussian distribution.

2.2 Kaplan-Meier estimator of the survival function

With right censored observations, the estimation of S only based on the uncensored observations does not lead to an estimator that converges to S . More precisely, assume that we observe a sample of possibly right censored data $(T_i, \delta_i)_{i=1, \dots, n}$. Let $n_1 = \sum_{i=1}^n \delta_i$ be the number of uncensored data in the sample. One may want to use

$$\hat{S}_n^{(1)}(t) = \frac{1}{n_1} \sum_{i=1}^n \mathbb{1}_{T_i > t, \delta_i = 1}$$

as an estimator of $S(\cdot)$. But this is not a good idea since one can show that

$$\hat{S}_n^{(1)}(t) \longrightarrow \int_t^{+\infty} (1 - G(t)) dF(t) \neq S(t)$$

except if $G(t) = 0$ for $x > t$ which means that there is no censoring.

2.2.1 Construction of the estimator

The estimator of Kaplan-Meier (1958) is based on the following statement: being alive after t amounts to being alive just before t and not dying at time t . Hence one has,

$$\begin{aligned} S(t) &= \mathbb{P}(X \geq t) \\ &= \mathbb{P}(X \geq t \mid X \geq t-1) \mathbb{P}(X \geq t-1) \\ &= \dots \\ &= \mathbb{P}(X \geq t \mid X \geq t-1) \mathbb{P}(X \geq t-1 \mid X \geq t-2) \dots \mathbb{P}(X \geq 1 \mid X \geq 0) \mathbb{P}(X \geq 0) \\ &= Q_t \times Q_{t-1} \dots \times Q_1 \times 1 \end{aligned} \tag{2.1}$$

where $Q_j = \mathbb{P}(X \geq j \mid X \geq j-1)$ represents the survival probability at j conditioned on being alive just before j .

Then one estimates $S(t)$ by the product of the estimations \hat{Q}_j of Q_j where \hat{Q}_j is the observed proportion of patients that are still alive after the j -th day among those alive just before j .

- If m_j deaths have occurred at j , then

$$\hat{Q}_j = \frac{n_j - m_j}{n_j}.$$

- If at day j , no death occurred, then $\hat{Q}_j = 1$.

As a consequence, only the death times appear in the estimation of $S(\cdot)$ and the estimator of $S(\cdot)$ is constant between two death times.

Examples

1. We observe 5 patients dead at days 3, 4, 6, 6 and 7. We want to estimate the survival probability after the 6-th day. We get

$$\begin{aligned}\hat{Q}_1 &= \hat{Q}_2 = 5/5 = 1 \\ \hat{Q}_3 &= 4/5 = 0.80 & m_3 = 1 & n_3 = 5 \\ \hat{Q}_4 &= 3/4 = 0.75 & m_4 = 1 & n_4 = 4 \\ \hat{Q}_5 &= 3/3 = 1 \\ \hat{Q}_6 &= 1/3 = 0.33 & m_6 = 2 & n_6 = 3\end{aligned}$$

Then

$$\hat{S}_5(6) = \hat{Q}_6 \hat{Q}_5 \hat{Q}_4 \hat{Q}_3 \hat{Q}_2 \hat{Q}_1 = \hat{Q}_6 \hat{Q}_4 \hat{Q}_3 = 4/5 \times 3/4 \times 1/3 = 1/5.$$

In this example, the probability to survive after day 6 is 20%. Since the data are uncensored, this probability can be computed making the ratio of the persons that are still alive after the day 6 and the total number of people.

2. Assume now that a right censoring has occurred at time 4. The survival times are then 3, 4⁺, 6, 6 and 7. One has

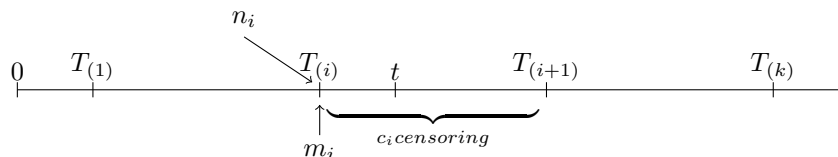
$$\begin{aligned}\hat{Q}_1 &= \hat{Q}_2 = 5/5 = 1 \\ \hat{Q}_3 &= 4/5 = 0.80 & m_3 = 1 & n_3 = 5 \\ \hat{Q}_4 &= 4/4 = 1 \\ \hat{Q}_5 &= 3/3 = 1 \\ \hat{Q}_6 &= 1/3 = 0.33 & m_6 = 2 & n_6 = 3\end{aligned}$$

Then

$$\hat{S}_5(6) = \hat{Q}_6 \hat{Q}_3 = 4/5 \times 1/3 = 0.27.$$

Then the probability to survive after day 6 is 27%.

More generally, to estimate the survival function with a n -sample, one has to order the observations increasingly according to their participation time. Let $T_{(1)} \leq \dots \leq T_{(n)}$ be the n ordered observed times and $\delta_{(1)}, \dots, \delta_{(n)}$ the corresponding indicators. Denote n_i the number of subjects at risk at $T_{(i)}$ (not dead nor censored) and m_i the number of patients dead in t_i .



At the origin, $T_{(0)} = 0$, $m_0 = 0$ and c_0 is the number of censored observations between 0 and $T_{(1)}$. As a consequence, one has

$$n_i = n_{i-1} - m_{i-1} - c_{i-1}$$

and thus

$$n_i = n_0 - \sum_{j=1}^{i-1} m_j - \sum_{j=1}^{i-1} c_j$$

that means that the number of subjects at risk at $T_{(i)}$ is the number of subjects at the origin minus the number of subjects dead before $T_{(i)}$ minus the number of subjects censored before $T_{(i)}$. A natural estimator of the survival function is given by

$$\hat{S}_{n,KM}(t) = \prod_{i:T_{(i)} \leq t} \hat{Q}_{T_{(i)}} = \prod_{i:T_{(i)} \leq t} \frac{n_i - m_i}{n_i} = \prod_{i:T_{(i)} \leq t} \left(1 - \frac{m_i}{n_i}\right).$$

using the decomposition (2.1) in product of S with $i = T_{(i)}$.

For $t < T_{(1)}$, one has by convention $\hat{S}_{n,KM}(t) = 1$.

If there is no ex-æquo in the sample

One has $n_i = n - i + 1$ and $m_i = 1$ at each death time and one can estimate the hazard rate function $h(\cdot)$ by

$$\begin{aligned}\hat{h}_n(T_{(i)}) &= \frac{\delta_{(i)}}{n - i + 1} \quad \text{for } i = 1, \dots, n \\ \hat{h}_n(t) &= 0 \quad \text{otherwise.}\end{aligned}$$

The **Nelson-Aalen estimator** of the cumulative hazard rate function $H(\cdot)$ is then naturally given by

$$\hat{H}_{n,NA}(t) = \sum_{i:T_{(i)} \leq t} \frac{\delta_{(i)}}{n - i + 1}.$$

We present in Section 2.4 the expression of the estimator in the general case.

The **Kaplan-Meier estimator** of the survival function $S(\cdot)$ is

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

If there are some ex-æquo in the sample

Let $T'_1 \leq \dots \leq T'_k$ be the k distinct and ordered observed times in the sample. The **Kaplan-Meier estimator** of the survival function $S(\cdot)$ is written in its general form as

$$\hat{S}_{n,KM}(t) = \prod_{i:T'_{(i)} \leq t} \left(1 - \frac{m_i}{n_i}\right).$$

If there is no censoring in the sample

If there is no censoring before $T_{(i)}$,

$$\begin{aligned}n_1 &= n_0 - m_0 = n_0, \\ n_j &= n_0 - \sum_{k=1}^{j-1} m_k \quad \text{for any } 1 < j \leq i.\end{aligned}$$

Thus, for any $T_{(j)} < T_{(i)}$, one has

$$\begin{aligned}\hat{S}_{n,KM} &= \frac{n_0 - m_1}{n_0} \frac{n_0 - m_1 - m_2}{n_0 - m_1} \times \dots \times \frac{n_0 - m_1 - m_2 - \dots - m_j}{n_0 - m_1 - m_2 - \dots - m_{j-1}} \\ &= \frac{n_0 - m_1 - m_2 - \dots - m_j}{n_0}.\end{aligned}$$

As a consequence, without censoring, $\hat{S}_{n,KM}$ is the observed proportion of the subjects still alive at t that is simply the empirical survival function $\hat{S}_n = 1 - \hat{F}_n$.

Remark 2.2.1 1. *This estimator is also called limit-product estimator since it can be obtained as the limit of a product.*

2. The estimator $\hat{S}_{n,KM}$ is a stepwise function with jumps at the observed death times. It reaches 0 only if the last observed delay $T_{(n)}$ of the sample corresponds to a real event (death or failure..) instead of a censoring i.e. if $\delta_{(n)} = 1$. When 0 is not reached, one observes a stage that traduces in general a few number of subjects that are controlled during a long period. That corresponds to a lack of information instead to a disappearance of the long-term risk.

3. Expression with the point processes:

$$\hat{S}_{n,KM}(t) = \prod_{i:T_{(i)} \leq t} \left(1 - \frac{\Delta \bar{N}(T_{(i)})}{\bar{Y}(T_{(i)})} \right)$$

where

- $\bar{Y}(t) = \sum_{i=1}^n Y_i(t)$ is the number of subjects at risk before t ,
- $\bar{N}(t) = \sum_{i=1}^n N_i(t)$ is the number of observed deaths before t ,
- $\Delta \bar{N}(t)$ corresponds to the number of observed deaths in t .

Remark 2.2.2 One can find another expression of the Kaplan-Meier estimator in the literature:

$$\hat{S}_{n,KM}(t) = \prod_{i:T_{(i)} \leq t} \left(1 - \frac{1}{n - i + 1} \right)^{\delta_{(i)}} \mathbb{1}_{\{t \leq T_{(n)}\}}.$$

Example of Freireich data: Freireich, in 1963, realized a therapeutic study in order to compare the remission durations (in weeks) of patients that suffer from leukemia. The patients are divided into two subgroups: some of them received a medicine (6 M-P) and the others a placebo. The results are the following:

6-MP	6, 6, 6, 6 ⁺ , 7, 9 ⁺ , 10, 10 ⁺ , 11 ⁺ , 13, 16, 17 ⁺ , 19 ⁺ , 20 ⁺ , 22, 23, 25 ⁺ , 32 ⁺ , 32 ⁺ , 34 ⁺ , 35 ⁺ .
Placebo	1, 1, 2, 2, 3, 4, 4, 5, 5, 8, 8, 8, 8, 11, 11, 12, 12, 15, 17, 22, 23.

The patients with a + sign correspond to lost subjects at the considered time of observation: they are censored, "excluded-alive" of the study and one only knows that their remission duration is greater than the observed delay.

Properties of the estimator

Proposition 2.2.3 (Optimality) $\hat{S}_{n,KM}$ maximizes the non parametric likelihood.

Theorem 2.2.4 (Consistency) Under some assumptions (fulfilled in particular as soon as X is absolutely continuous) and if the censoring times are i.i.d. with common c.d.f. G , we have

$$\sup_{t \geq 0} \left| \hat{F}_n(t) - F(t) \right| \xrightarrow{a.s.} 0,$$

Theorem 2.2.5 (Asymptotic normality) Under some assumptions (fulfilled in particular as soon as X is absolutely continuous) and if the censoring times are i.i.d. with common c.d.f. G , we have

$$\sqrt{n} \left(\hat{F}_n(\cdot) - F(\cdot) \right) \xrightarrow{\mathcal{L}} S(\cdot)B(W(\cdot)),$$

where $B(\cdot)$ is a Brownian motion on \mathbb{R}^+ and $W(\cdot)$ is the function defined by

$$W(t) = \int_0^t \frac{dF(u)}{S^2(u)(1 - G(u^-))}, \quad \text{for all } t > 0.$$

One can see that the process $Z_2(\cdot) := S(\cdot)B(W(\cdot))$ is a centered Gaussian process with covariance function

$$\text{Cov}(Z_2(s), Z_2(t)) = S(s)S(t) \int_0^{s \wedge t} \frac{dF(u)}{S^2(u)(1 - G(u^-))}.$$

Let us denote the **Greenwood estimator** $\hat{\sigma}_n^2(t)$ by

$$\hat{\sigma}_{n,KM}^2(t) = \hat{S}_{n,KM}(t)^2 \sum_{i:T'_{(i)} \leq t} \frac{m_i}{n_i(n_i - m_i)}.$$

In particular, if there is no ex-æquo, this expression becomes

$$\hat{\sigma}_n^2(t) = \hat{S}_{n,KM}(t)^2 \sum_{i:T'_{(i)} \leq t} \frac{\delta_i}{(n-i)(n-i+1)}.$$

Variance estimation

One can show that the Greenwood estimator $\hat{\sigma}_n^2(t)$ is an estimation of the variance of the Kaplan-Meier estimator and is a uniformly consistent estimator of the asymptotic variance function

$$S(t)^2 \int_0^t \frac{dF(u)}{S^2(u)(1 - G(u^-))}.$$

Without censoring, one recover the estimation of the variance of a proportion:

$$\frac{\hat{S}_{n,KM}(t)(1 - \hat{S}_{n,KM}(t))}{n_0}.$$

In Section 2.3, we explain the origin of this estimator of the variance of $\hat{S}_{n,KM}(t)$.

Confidence interval for $S(t)$

Assuming the gaussianity of the asymptotic behavior of $\hat{S}_{n,KM}(t)$ an asymptotic $(1 - \alpha)$ -confidence interval for $S(t)$ is given by

$$\left[\hat{S}_{n,KM}(t) - z_{1-\alpha/2} \hat{\sigma}_{n,KM}(t), \hat{S}_{n,KM}(t) + z_{1-\alpha/2} \hat{\sigma}_{n,KM}(t) \right]$$

where $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of the standard Gaussian distribution.

This interval is the one given by the option `conf.type = plain` of the `survfit` function that is available in R.

As an estimator of the expectation $\mathbb{E}[X]$ one can use

$$\hat{\mu}_n = \int_0^{+\infty} t d\hat{F}_n(t) = \sum_{i=1}^k T'_{(i)} \Delta \hat{F}_n(T'_{(i)}) = \sum_{i=1}^k T'_{(i)} \frac{m_i}{n_i} \prod_{j=1}^{i-1} \left(1 - \frac{m_j}{n_j} \right).$$

We have

$$\hat{\mu}_n \xrightarrow{a.s.} \int x dF(x) = \mathbb{E}[X], \quad \text{as } n \rightarrow +\infty.$$

2.3 Breslow estimator of the cumulative hazard rate function

One can deduce from $\hat{S}_{n,KM}$ an estimator of the cumulative hazard rate function using the relation $H(t) = -\ln S(t)$:

$$\hat{H}_{n,BR}(t) = -\ln \hat{S}_{n,KM}(t) = -\sum_{i:T'_{(i)} \leq t} \ln \hat{Q}_{T'_{(i)}} = -\sum_{i:T'_{(i)} \leq t} \ln \frac{n_i - m_i}{n_i}.$$

Theorem 2.3.1 (Consistency) *Under some assumptions (fulfilled in particular as soon as X is absolutely continuous) and if the censoring times are i.i.d. with common c.d.f. G , we have*

$$\sup_{t \geq 0} \left| \hat{H}_n(t) - H(t) \right| \xrightarrow{a.s.} 0.$$

Theorem 2.3.2 (Asymptotic normality) *Under some assumptions (fulfilled in particular as soon as X is absolutely continuous) and if the censoring times are i.i.d. with common c.d.f. G , we have*

$$\sqrt{n} \left(\hat{H}_n(\cdot) - H(\cdot) \right) \xrightarrow{\mathcal{L}} B(W(\cdot)),$$

where $B(\cdot)$ is a Brownian motion on \mathbb{R}^+ and $W(\cdot)$ is the function defined by

$$W(t) = \int_0^t \frac{dF(u)}{S^2(u)(1-G(u^-))}, \quad \text{for all } t > 0.$$

Variance estimation

One has

$$\text{Var} \left(\hat{H}_{n,BR}(t) \right) = \text{Var} \left(\sum_{i: T_{(i)} \leq t} \ln \hat{Q}_{T_{(i)}} \right).$$

The random variables $\hat{Q}_{T_{(i)}}$ are not independent, but assuming it is the case, one gets the following variance approximation:

$$\text{Var} \left(\hat{H}_{n,BR}(t) \right) \approx \sum_{i: T_{(i)} < t} \text{Var} \left(\ln \hat{Q}_{T_{(i)}} \right).$$

Now we apply the delta method [1] to approximate $\text{Var} \left(\ln \hat{Q}_{T_{(i)}} \right)$: for any regular function f ,

$$\text{Var}(f(X)) \approx f'(\mathbb{E}[X])^2 \text{Var}(X)$$

and we use the fact that the random variables $n_i \hat{Q}_{T_{(i)}}$ are binomial with parameters n_i and $Q_{T_{(i)}}$. We get

$$\text{Var} \left(\hat{H}_{n,BR}(t) \right) \approx \sum_{i: T_{(i)} < t} \frac{1 - \hat{Q}_{T_{(i)}}}{n_i \hat{Q}_{T_{(i)}}} \approx \sum_{i: T_{(i)} < t} \frac{m_i}{n_i(n_i - m_i)}.$$

Thus an approximation of the variance of $\hat{H}_{n,BR}(t)$ is given by

$$\hat{\sigma}_{n,BR}^2(t) := \text{Var}(\widehat{\hat{H}_{n,BR}(t)}) = \sum_{i: T_{(i)} \leq t} \frac{m_i}{n_i(n_i - m_i)}.$$

Confidence intervals for $H(t)$ and $S(t)$

Assuming the gaussianity of the asymptotic behavior of $\hat{H}_{n,BR}(t)$, we deduce an asymptotic $(1 - \alpha)$ -confidence interval for $H(t)$:

$$\left[\hat{H}_{n,BR}(t) - z_{1-\alpha/2} \hat{\sigma}_{n,BR}(t), \hat{H}_{n,BR}(t) + z_{1-\alpha/2} \hat{\sigma}_{n,BR}(t) \right],$$

where $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of the standard Gaussian distribution.

Using the relation that links S to H : $S(t) = \exp\{-H(t)\}$, one deduces an asymptotic $(1 - \alpha)$ -confidence interval for $S(t)$:

$$\left[\exp \left\{ -(\hat{H}_{n,BR}(t) + z_{1-\alpha/2} \hat{\sigma}_{n,BR}(t)) \right\}, \exp \left\{ -(\hat{H}_{n,BR}(t) - z_{1-\alpha/2} \hat{\sigma}_{n,BR}(t)) \right\} \right]. \quad (2.2)$$

This interval is the one given by the option `conf.type = log` of the `survfit` function that is available in R.

Remark 2.3.3 (Estimation of the variance of $\hat{S}_{n,KM}$) The Greenwood estimator $\hat{\sigma}_{n,KM}^2(t)$ of the variance of $\hat{S}_{n,KM}$ can be deduced from the previous variance. It has been obtained by the following considerations. We use once again the delta method [1]: for any regular function f ,

$$\text{Var}(f(X)) \approx f'(\mathbb{E}[X])^2 \text{Var}(X)$$

to get the following approximation

$$\begin{aligned} \text{Var}\left(\hat{S}_{n,KM}(t)\right) &\approx \hat{S}_{n,KM}(t)^2 \text{Var}\left(\ln \hat{S}_{n,KM}(t)\right) \\ &= \hat{S}_{n,KM}(t)^2 \text{Var}\left(\hat{H}_{n,BR}(t)\right) \\ &\approx \hat{S}_{n,KM}(t)^2 \sum_{i:T(i) < t} \frac{m_i}{n_i(n_i - m_i)} \\ &\approx \hat{S}_{n,KM}(t)^2 \hat{\sigma}_{n,BR}^2(t). \end{aligned}$$

One deduces an asymptotic $(1 - \alpha)$ -confidence interval for $S(t)$:

$$\left[\hat{S}_{n,KM}(t) (1 - z_{1-\alpha/2} \hat{\sigma}_{n,BR}(t)), \hat{S}_{n,KM}(t) (1 + z_{1-\alpha/2} \hat{\sigma}_{n,BR}(t)) \right]$$

and since, $\hat{S}_{n,KM}(t) = \exp(-\hat{H}_{n,BR}(t))$, this interval corresponds also to:

$$\left[\exp\{-\hat{H}_{n,BR}(t)\} (1 - z_{1-\alpha/2} \hat{\sigma}_{n,BR}(t)), \exp\{-\hat{H}_{n,BR}(t)\} (1 + z_{1-\alpha/2} \hat{\sigma}_{n,BR}(t)) \right]. \quad (2.3)$$

If $\hat{\sigma}_{n,BR}(t)$ is close to 0, one gets an equivalent expression for both intervals in (2.2) and (2.3).

2.4 Nelson-Aalen estimator of the cumulative hazard function

For ds small, one as

$$H(s + ds) - H(s) \sim h(s)ds \sim P(s < X < s + ds \mid X > s)$$

that is naturally estimated by $(\bar{N}(s + ds) - \bar{N}(s)) / \bar{Y}(s)$.

Summing these quantities over the sub intervals of $(0, t]$ and making their length tends to 0 in such a way that they contains only at most one element, we derive the **Nelson-Aalen estimator** of the cumulative hazard function:

$$\hat{H}_{n,NA}(t) = \int_0^t \frac{d\bar{N}(s)}{\bar{Y}(s)}$$

that can be written in a discrete form:

$$\hat{H}_{n,NA}(t) = \sum_{i:T(i) \leq t} \frac{d\bar{N}(T(i))}{\bar{Y}(T(i))}$$

By convention, $\hat{H}_{n,NA}(t) = 0$ for $t < T_{(1)}$.

An equivalent expression without the point processes is the following:

$$\hat{H}_{n,NA}(t) = \sum_{i:T(i) \leq t} \frac{m_i}{n_i}.$$

Property and interpretation:

1. $\hat{H}_{n,NA}$ is a stepwise function whose jumps are equal to $\frac{m_i}{n_i}$ at each observed death $T_{(i)}$.
2. $H(t)$ represents the mean number of deaths (or failures) over $(0, t]$ for an individual perpetually at risk over the interval $[0, t]$: an electric component is working during t hours; at each failure, we replace the current component by another component that has functioned the same duration than the one we are replacing (in reliability, this method is called protocol of minimal reparation).

3. The derivative of H represents the hazard rate function h . Since the estimator $\hat{H}_{n,NA}$ is a stepwise function, it is not possible to compute its derivative. As a consequence, to estimate h , as for the estimation of any density probability function, one needs to proceed to the smoothing of $\hat{H}_{n,NA}$.

Variance estimation

A consistent estimator of the variance of $\hat{H}_{n,NA}(t)$ is given by

$$\sigma_{n,NA}^2(t) := \text{Var}(\widehat{\hat{H}_{n,NA}}(t)) = \sum_{i:T(i) \leq t} \frac{d\bar{N}(t_i)}{\bar{Y}(t_i)^2} = \sum_{i:T(i) \leq t} \frac{m_i}{n_i^2}.$$

Confidence interval for $H(t)$

Assuming the gaussianity of $\hat{H}_{n,NA}(t)$, one deduces an asymptotic $(1 - \alpha)$ -confidence interval for $H(t)$:

$$\left[\hat{H}_{n,NA}(t) - z_{1-\alpha/2} \sigma_{n,NA}(t), \hat{H}_{n,NA}(t) + z_{1-\alpha/2} \sigma_{n,NA}(t) \right].$$

2.5 Harrington and Fleming estimator of the survival function

Since $S(t) = \exp\{-H(t)\}$, we deduce from the Nelson-Aalen estimator of the cumulative hazard rate function the **Harrington and Fleming estimator** of the survival function:

$$\hat{S}_{n,HF}(t) = \exp\left\{-\hat{H}_{n,NA}(t)\right\}.$$

The relation $\hat{H}_{n,NA}(t) = \sum_{i:T(i) \leq t} \frac{m_i}{n_i}$ yields

$$\hat{S}_{n,HF}(t) = \exp\left\{-\sum_{i:T(i) \leq t} \frac{m_i}{n_i}\right\} = \prod_{i:T(i) \leq t} \exp\left\{-\frac{m_i}{n_i}\right\}$$

that we should compare to the Kaplan-Meier estimation $\hat{S}_{n,KM}(t) = \prod_{i:T(i) \leq t} \frac{n_i - m_i}{n_i}$ of the survival function.

Variance estimation

By the delta method [1], $\text{Var}(\hat{S}_{n,HF}(t)) = \hat{S}_{n,HF}(t)^2 \text{Var}(\hat{H}_{n,NA}(t))$ that leads to an estimator of the variance of $\hat{S}_{n,HF}$:

$$\sigma_{n,HF}^2(t) := \text{Var}(\widehat{\hat{S}_{n,HF}}(t)) = \exp\left\{-2 \sum_{i:T(i) \leq t} \frac{m_i}{n_i}\right\} \sum_{i:T(i) \leq t} \frac{m_i}{n_i^2}.$$

Confidence interval for $S(t)$

Assuming the gaussianity of $\hat{S}_{n,HF}(t)$, one deduces an asymptotic $(1 - \alpha)$ -confidence interval for $S(t)$:

$$\left[\hat{S}_{n,HF}(t) - z_{1-\alpha/2} \sigma_{n,HF}(t), \hat{S}_{n,HF}(t) + z_{1-\alpha/2} \sigma_{n,HF}(t) \right].$$

2.6 Comments

We introduced two different estimators of the cumulative hazard rate function H :

1. the Breslow estimator deduced from the Kaplan-Meier estimation of the survival function,
2. the Nelson-Aalen estimator.

One may find in the literature a third estimation of H that comes from the actuarial method. These three different estimators of H are asymptotically equivalent and in practice one may prefer use the first one because of its relation with the Kaplan-Meier estimator.

Chapter 3

Comparison of the survival functions of two (or more) groups

3.1 The weighted logrank tests

These tests allows us to test the equality of K survival distributions through a sample of censored data. Here,

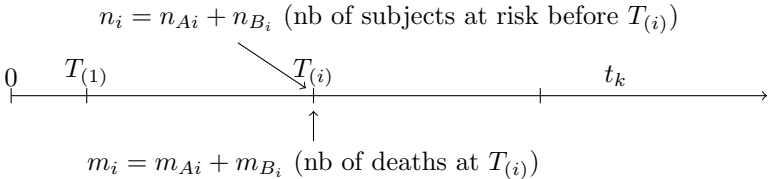
$$\begin{aligned}
 H_0 : S_1(t) = S_2(t) = \dots = S_K(t) \quad \forall t \quad \text{versus} \\
 H_1 : \text{ at least two groups do not share the same survival function.}
 \end{aligned}$$

3.1.1 Comparison of two groups

The goal is to test

$$\begin{aligned}
 H_0 : S_A(t) = S_B(t) \quad \forall t \quad \text{versus} \\
 H_1 : \exists t \quad \text{such that} \quad S_A(t) \neq S_B(t).
 \end{aligned}$$

Let $T_{(1)} < T_{(2)} < \dots < T_{(k)}$ be the k ordered times of death of the sample that gathers the data of the two groups A and B .



where n_i represents the number of subjects at risk in $T_{(i)}$ and m_i the number of subjects that die in $T_{(i)}$.

For any time of death $T_{(i)}$, one can resume the information in the following 2×2 tabular:

	Dead	Alive after $T_{(i)}$	Total
Group A	m_{Ai}	$n_{Ai} - m_{Ai}$	n_{Ai}
Group B	m_{Bi}	$n_{Bi} - m_{Bi}$	n_{Bi}
Total	m_i	$n_i - m_i$	n_i

Under the null hypothesis H_0 , in $T_{(i)}$, the proportion of dead among the subjects at risk is identical in both groups. Considering that the marginals of the tabular are fixed, under the null hypothesis H_0 , m_{Ai} is distributed as an hypergeometric random variable whose parameters are (m_i, n_{Ai}, n_{Bi}) . Then,

$$\begin{aligned}\mathbb{E}(m_{Ai}) &= e_{Ai} = m_i \frac{n_{Ai}}{n_i}, \\ \text{Var}(m_{Ai}) &= v_{Ai} = m_i \frac{n_i - m_i}{n_i - 1} \frac{n_{Ai} n_{Bi}}{n_i^2}.\end{aligned}$$

The random variable $m_{Ai} - e_{Ai}$ is centered under H_0 and one can prove that m_{Ai} and m_{Aj} are uncorrelated as soon as $i \neq j$. Then we naturally define the statistics of the weighted logrank test by

$$U_A = \sum_{i=1}^k w_i \left(m_{Ai} - m_i \frac{n_{Ai}}{n_i} \right)$$

where the weight w_i depends on $T_{(i)}$. Obviously,

$$LR = \frac{U_A}{\sqrt{\text{Var}(U_A)}} = \frac{\sum_{i=1}^k w_i \left(m_{Ai} - m_i \frac{n_{Ai}}{n_i} \right)}{\sqrt{\sum_{i=1}^k w_i^2 \left(m_i \frac{n_i - m_i}{n_i - 1} \frac{n_{Ai} n_{Bi}}{n_i^2} \right)}}$$

converges under H_0 to a standard Gaussian random variable and equivalently

$$\frac{U_A^2}{\text{Var}(U_A)} \xrightarrow[H_0]{\mathcal{L}} \chi_1^2.$$

Remark 3.1.1 1. *The weighted logrank statistics are ordered statistics: they only depend on the ranks of the observations instead of their exact value.*

2. *Breslow proposed a simplified form of the logrank statistics that comes from another way to take into account the ex-æquo:*

$$LR_{BR} = \frac{U_A}{\sqrt{\text{Var}U_A}} = \frac{\sum_{i=1}^k w_i \left(m_{Ai} - m_i \frac{n_{Ai}}{n_i} \right)}{\sqrt{\sum_{i=1}^k w_i^2 m_i \frac{n_{Ai} n_{Bi}}{n_i^2}}}.$$

Examples of weighting:

1. The **logrank test** (also called **Mantel-Haenszel test**): $w_i = 1 \forall i$.

All the deaths have the same weights. It is the simplest weighting. The test amounts to compare the observed number O_A of deaths in the group A to the expected number E_A of deaths. Under H_0 , we have

$$O_A = \sum_{i=1}^k m_{Ai} \quad \text{and} \quad E_A = \sum_{i=1}^k e_{Ai}.$$

The test statistics is then

$$\frac{(O_A - E_A)^2}{\sum_{i=1}^k v_i}$$

2. The **Gehan test**: $w_i = n_i \forall i$ (available in the software SAS under the name Wilcoxon since it is a generalization of the Wilcoxon-Mann-Whitney test in the context of censoring).

The first deaths have greater weights.

3. The **Peto and Prentice test** (close to the Kaplan-Meier estimator): $w_i = S_i \quad \forall i$ with $S_i^* = \prod_{j=1}^i \frac{n_j}{n_j + m_j}$ (available in SAS).
The first deaths have greater weights.
4. The **Tarone and Ware test**: $w_i = \sqrt{n_i} \quad \forall i$.
5. The **Harrington and Fleming test**: $w_i = \hat{S}_{n,KM}(T_{(i)})^\rho \quad \forall i$ with $0 \leq \rho \leq 1$.
The logrank test is a particular case of the Harrington and Fleming test when $\rho = 0$ and for $\rho = 1$, we get a test that is close to the one of Peto and Prentice.

These tests are available in R (function `survdif`).

Approximated logrank test

The observed total number of deaths in both groups is equal, under H_0 , to the expected total number of deaths:

$$O_A + O_B = E_A + E_B$$

and

$$\text{Var}(O_A - E_A) = \text{Var}(O_B - E_B).$$

The logrank statistics can be rewritten equivalently

$$\frac{(O_A - E_A)^2}{\text{Var}(O_A - E_A)} \quad \text{or} \quad \frac{(O_B - E_B)^2}{\text{Var}(O_B - E_B)}.$$

One can show that the statistics

$$(LRA)^2 = \frac{(O_A - E_A)^2}{E_A} + \frac{(O_B - E_B)^2}{E_B}$$

is always less than or equal to the logrank one.

1. This statistics is easier to compute than the previous one and is called the **approximated logrank statistics**.
2. It reminds the χ_2 statistics used to compare two observed proportions.
3. It is a conservative statistics: anytime it rejects H_0 , the classical logrank statistics would have given the same conclusion. However the approximated logrank test is less powerful than the classical one.

Comparison criterion

In order to quantify the mortality difference between the groups and in analogy with the notion of relative risk, the quantity

$$RR_1 = \frac{O_B/E_B}{O_A/E_A}$$

has been proposed as an estimation of the ratio of the risk functions in each group. One might be cautious using this estimation since it is biased and the bias increases with the real ratio of the risk functions and the sample size.

One may compute this relative risk using the Gehan or Peto and Prentice weighting introducing

$$O_g = \sum_{i=1}^k w_i m_{gi} \quad \text{and} \quad E_g = \sum_{i=1}^k w_i e_{gi}.$$

the letter g stands for the group: A or B .

Other estimators have been introduced: for example,

$$RR_2 = \frac{\sum_{i=1}^k m_{Bi}(n_{Ai} - m_{Ai})/n_i}{\sum_{i=1}^k m_{Ai}(n_{Bi} - m_{Bi})/n_i} \quad \text{or} \quad RR_3 = \frac{\sum_{i=1}^k m_{Bi}n_{Ai}/n_i}{\sum_{i=1}^k m_{Ai}n_{Bi}/n_i}.$$

Weighting choice

As usual, we are looking for the most powerful test. As a consequence the choice of the weighting depends on the alternative hypothesis H_1 . One can prove that, for a given alternative, the optimal asymptotic weights are proportional to $\ln\left(\frac{h_A(t)}{h_B(t)}\right)$.

- Logrank test ($w_i = 1$): optimal for the proportional alternatives

$$H_0 : h_A(t) = h_B(t) \quad \forall t \quad \text{versus} \quad H_1 : h_B(t) = rh_A(t) \quad \forall t.$$

These alternatives correspond to hypothesis of the Cox model (we see in the chapter dedicated to the Cox model that the logrank test corresponds to the scoring test of the Cox model with two groups).

An equivalent expression is given by: $S_B(t) = S_A(t)^r$.

- Peto and Prentice test: optimal for the alternatives such that $\ln\left(\frac{h_A(t)}{h_B(t)}\right) = S_A(t)$.

The Gehan and Peto and Prentice tests are pretty well adapted when there are numerous premature deaths.

3.1.2 Generalization to the comparison of K groups

The goal is to test

$$H_0 : S_1(t) = S_2(t) = \dots = S_K(t) \quad \forall t$$

versus

$$H_1 : \quad \text{at least two groups do not share the same survival distribution.}$$

We compute the statistics U_1, U_2, \dots, U_{K-1} , sommes pondérées en chaque temps de décès observé de la différence entre le nombre de décès observés dans le groupe et le nombre de décès attendus sous H_0 .

The test statistic is given by

$$X = (U_1, U_2, \dots, U_{K-1})\Sigma^{-1}(U_1, U_2, \dots, U_{K-1})'$$

where Σ is the variance-covariance matrix of the vector $(U_1, U_2, \dots, U_{K-1})$.

Under the null hypothesis H_0 , X converges to a χ^2 distribution with $K - 1$ degrees of freedom.

3.2 Comparison with adjustment: stratified logrank test

Assume that the survival is strongly linked to the age of the patient (which is quite natural), then it is more relevant to compare groups at a fixed given age.

In that view, we introduce a qualitative factor with M strata where each strata s corresponds to a sub sample.



One defines

$$U_s = \sum_{i=1}^{k_s} \left(m_{Asi} - m_{si} \frac{n_{Asi}}{n_{si}} \right);$$

then

$$\text{Var}(U_s) = \sum_{i=1}^{k_s} \left(m_{si} \frac{n_{si} - m_i}{n_{si} - 1} \frac{n_{Asi}n_{Bsi}}{n_{si}^2} \right)$$

The stratified logrank test statistics is then defined by

$$U_{\text{stratified}} = \frac{\sum_{s=1}^M U_s}{\sqrt{\sum_{s=1}^M \text{Var}(U_s)}} \xrightarrow[H_0]{\mathcal{L}} \mathcal{N}(0, 1).$$

Expression through the point processes

The corresponding notation are:

$$\begin{aligned} \bar{Y}_A(T_{(i)}) &\equiv n_{Ai}, \\ \bar{Y}_B(T_{(i)}) &\equiv n_{Bi}, \\ \bar{Y}(T_{(i)}) &\equiv n_i, \\ \Delta \bar{N}_A(T_{(i)}) &\equiv m_{Ai}, \\ \Delta \bar{N}_B(T_{(i)}) &\equiv m_{Bi}, \\ \Delta \bar{N}(T_{(i)}) &\equiv m_i. \end{aligned}$$

We rewrite the logrank statistics as the sum of the differences at each death time between the observed hazard rate in group A and the one expected under the null hypothesis H_0 . More precisely, the observed hazard rate group A in $T_{(i)}$ is

$$\frac{\Delta \bar{N}_A(T_{(i)})}{\bar{Y}_A(T_{(i)})}$$

while the expected hazard rate under H_0 is

$$\frac{\Delta \bar{N}(T_{(i)})}{\bar{Y}(T_{(i)})}.$$

Then the weighted logrank test statistics is

$$U_A = \sum_{i=1}^k W_i(T_{(i)}) \left(\frac{\Delta \bar{N}_A(T_{(i)})}{\bar{Y}_A(T_{(i)})} - \frac{\Delta \bar{N}(T_{(i)})}{\bar{Y}(T_{(i)})} \right)$$

where $W_i(T_{(i)})$ is a non negative weighting function that has usually the following form $W_i(T_{(i)}) = \bar{Y}_A(T_{(i)})W(T_{(i)})$ with W a common weight to both groups.

One directly gets the following equivalent expression:

$$U_A = \sum_{i=1}^k W(T_{(i)}) \left(\Delta \bar{N}_A(T_{(i)}) - \bar{Y}_A(T_{(i)}) \frac{\Delta \bar{N}(T_{(i)})}{\bar{Y}(T_{(i)})} \right).$$

Bibliography

- [1] A. W. van der Vaart. *Asymptotic statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 1998.