

STK4080 SURVIVAL AND EVENT HISTORY ANALYSIS

Slides 4: Introduction to nonparametric estimation

Bo Lindqvist
Department of Mathematical Sciences
Norwegian University of Science and Technology
Trondheim

<https://www.ntnu.no/ansatte/bo.lindqvist>
bo.lindqvist@ntnu.no

University of Oslo, Autumn 2019

NONPARAMETRIC ESTIMATION OF $S(t)$

We are interested in estimating the distribution of the lifetime T of some equipment or the time to some given event in a medical context.

We have indicated how parametric models like exponential and Weibull can be fitted to data.

Now we shall instead see how in particular $S(t)$ can be estimated without making parametric assumptions.

Thus, instead of having to restrict to estimation of one or two parameters, we now have an infinite number of possible functions $S(t)$ to choose from. (Essentially, the only restriction is that it is decreasing, starts in 1 and converges to 0 as $t \rightarrow \infty$.)

In this case our observations are the exact failure times T_1, \dots, T_n , assumed to be i.i.d. observations of a lifetime T .

Hence we can estimate $S(t) = P(T > t)$ for a given $t > 0$ by the relative proportion of lifetimes that exceed t :

$$\hat{S}(t) = \frac{\text{number of } T_i > t}{n}$$

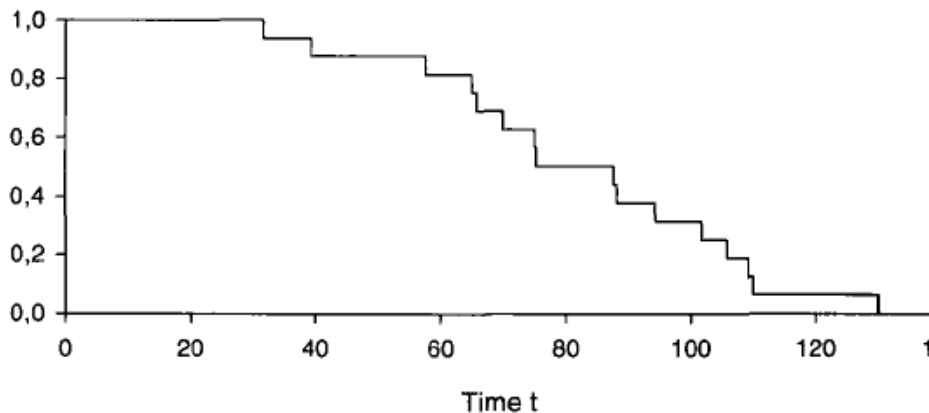
This is called the *empirical survival function*.

If we order the observations as $T_{(1)} < T_{(2)} < \dots < T_{(n)}$, then $\hat{S}(t)$ starts at 1 for $t = 0$ and makes a downward jump of $1/n$ at $T_{(1)}$, a new downward jump of $1/n$ at $T_{(2)}$, and so on until it jumps from $1/n$ to 0 at $T_{(n)}$.

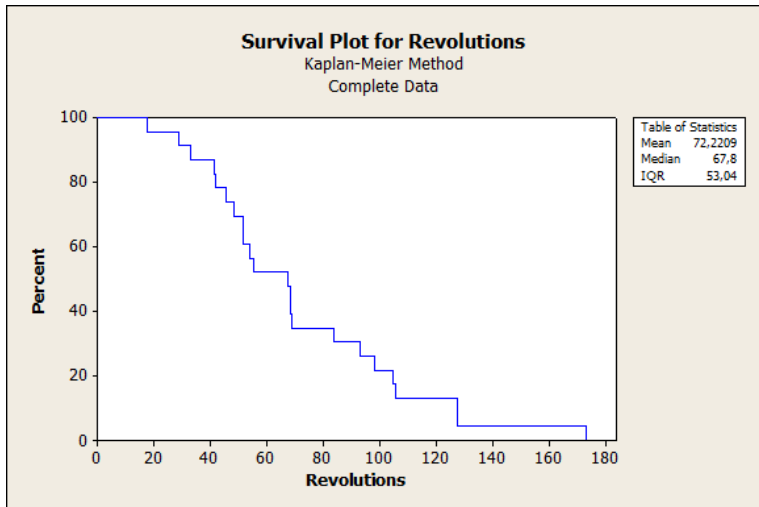
EXAMPLE OF EMPIRICAL SURVIVAL PLOT, $\hat{S}(t)$

$n = 16$ observed lifetimes:

31.7, 39.2, 57.5, 65.0, 65.8, 70.0, 75.0, 75.2, 87.7, 88.3, 94.2, 101.7,
105.8, 109.2, 110.0, 130.0



EMPIRICAL SURVIVAL PLOT FOR BALL BEARING DATA



Consider n individuals, where the i th individual has potential lifetime T_i and potential censoring time C_i . We *observe* the pair (\tilde{T}_i, D_i) , where

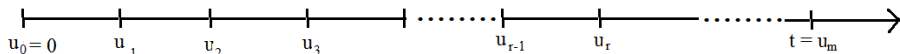
$$\begin{aligned}\tilde{T}_i &= \min(T_i, C_i) \\ D_i &= \begin{cases} 1 & \text{if } \tilde{T}_i = T_i \\ 0 & \text{if } \tilde{T}_i = C_i \end{cases}\end{aligned}$$

Assume:

- T_1, T_2, \dots, T_n are *independent and identically distributed* with common reliability function $S(t)$.
- The censoring mechanism satisfies the property of *independent censoring*.

The estimator is constructed in the following.

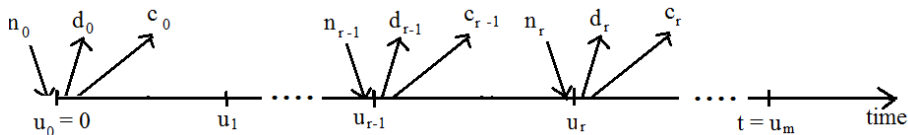
MAIN IDEA OF CONSTRUCTION



Assume first that time is measured on a discrete scale with values $u_0 = 0 \leq u_1 \leq u_2 \leq \dots$, so that all T_i , C_i and hence \tilde{T}_i are among these. Let $t = u_m$. Then

$$\begin{aligned} S(t) &= P(T > t) = P(T > u_m) \\ &= P(T > u_m \cap T > u_{m-1} \cap \dots \cap T > u_2 \cap T > u_1 \cap T > u_0) \\ &= P(T > u_0) \cdot P(T > u_1 \mid T > u_0) \cdot P(T > u_2 \mid T > u_1 \cap T > u_0) \\ &\quad \dots P(T > u_r \mid T > u_{r-1} \cap T > u_{r-2} \dots \cap T > u_0) \dots \\ &\quad \dots P(T > u_m \mid T > u_{m-1} \cap \dots \cap T > u_0) \\ &= P(T > u_0) \cdot P(T > u_1 \mid T > u_0) \cdot P(T > u_2 \mid T > u_1) \\ &\quad \dots P(T > u_r \mid T > u_{r-1}) \dots P(T > u_m \mid T > u_{m-1}) \end{aligned}$$

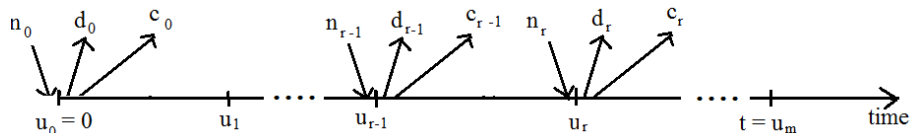
Idea: Estimate each factor $P(T > u_r \mid T > u_{r-1})$, from data (Y_i, δ_i) ; $i = 1, \dots, n$.



Define:

- n_r = number at risk at time u_r ; i.e. number that can fail at u_r ; counted immediately before u_r .
- d_r = number failing at u_r (those with $Y = u_r$, $\delta = 1$)
- c_r = number censored at u_r (those with $Y = u_r$, $\delta = 0$); assumed to be censored right after u_r , and by convention after all failures at u_r (in practice in the interval following u_r)

CONSTRUCTION OF ESTIMATOR (CONT.)



Note: The d_i , c_i are found directly from the data, while the n_i are found recursively as:

$$n_0 = n$$

$$n_1 = n_0 - d_0 - c_0$$

...

$$n_r = n_{r-1} - d_{r-1} - c_{r-1}$$

Then estimate,

$$P(T > u_r \mid T > u_{r-1}) = 1 - P(T = u_r \mid T > u_{r-1}) \approx 1 - \frac{d_r}{n_r} = \frac{n_r - d_r}{n_r}$$

$$\& \quad P(T > u_0) = 1 - P(T = u_0) \approx 1 - \frac{d_0}{n_0} = \frac{n_0 - d_0}{n_0}$$

It follows that $S(t) = P(T > t)$ can be estimated by

$$\hat{S}(t) = \frac{n_0 - d_0}{n_0} \cdot \frac{n_1 - d_1}{n_1} \cdots \frac{n_r - d_r}{n_r} \cdots \frac{n_m - d_m}{n_m}$$

Note that these factors are 1, whenever $d_r = 0$. Thus

$$\hat{S}(t) = \prod_{\substack{\text{all } u_r \leq t \\ \text{with } d_r \geq 1}} \frac{n_r - d_r}{n_r}$$

In practice we have continuous time. But this case can be approximated by making the grid $u_1 < u_2 < \cdots$ finer and finer.

Thus in general the KM-estimator is given by:

If $T_{(1)} < T_{(2)} < \cdots$, are the times with at least one failure, and n_i, d_i are, respectively, the number at risk and the number of failures at $T_{(i)}$, then

$$\hat{S}(t) = \prod_{i: T_{(i)} \leq t} \frac{n_i - d_i}{n_i}$$

GREENWOOD'S FORMULA FOR VARIANCE OF THE KM-ESTIMATOR

$$\widehat{\text{Var}}(\hat{S}(t)) = (\hat{S}(t))^2 \cdot \sum_{T_{(i)} \leq t} \frac{d_i}{n_i(n_i - d_i)}$$

It can be shown that for large n , $\hat{S}(t)$ is approximately normally distributed,

$$\hat{S}(t) \approx N(S(t), \widehat{SD}(\hat{S}(t)))$$

Thus an approximate 95% confidence interval can be obtained for each t by

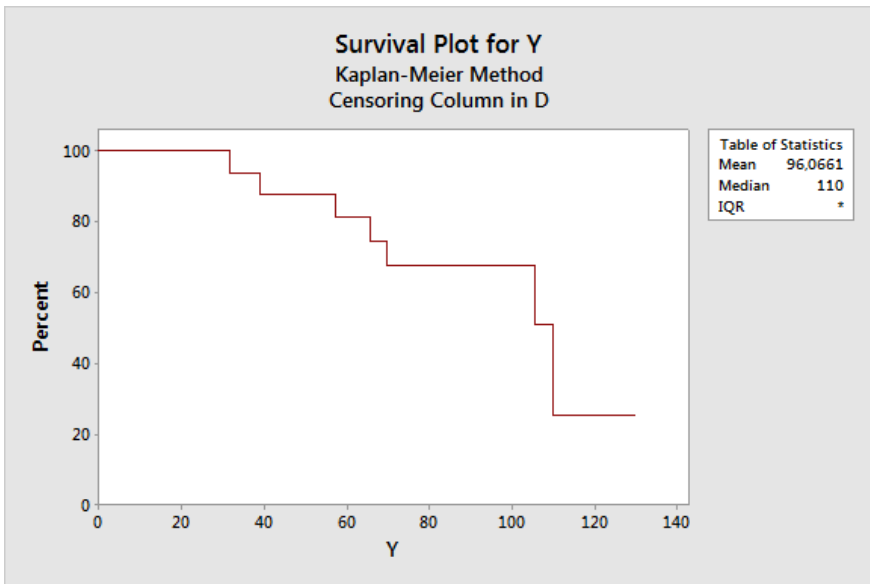
$$P(\hat{S}(t) - 1.96 \cdot \widehat{SD}(\hat{S}(t)) \leq S(t) \leq \hat{S}(t) + 1.96 \cdot \widehat{SD}(\hat{S}(t)))$$

KM-ESTIMATOR FOR CENSORED DATA (MINITAB)

Row	Y	D
1	31,7	1
2	39,2	1
3	57,5	1
4	65,0	0
5	65,8	1
6	70,0	1
7	75,0	0
8	75,2	0
9	87,5	0
10	88,3	0
11	94,2	0
12	101,7	0
13	105,8	1
14	109,2	0
15	110,0	1
16	130,0	0

Time	Number at Risk	Number Failed	Survival Probability	Standard Error	95,0% Normal CI Lower	Upper
31,7000	16	1	0,9375	0,0605	0,8189	1,0000
39,2000	15	1	0,8750	0,0827	0,7130	1,0000
57,5000	14	1	0,8125	0,0976	0,6213	1,0000
65,8000	12	1	0,7448	0,1105	0,5283	0,9613
70,0000	11	1	0,6771	0,1194	0,4431	0,9111
105,8000	4	1	0,5078	0,1718	0,1711	0,8445
110,0000	2	1	0,2539	0,1990	0,0000	0,6440

KM-PLOT FOR CENSORED DATA



KM-PLOT WITH CONFIDENCE LIMITS

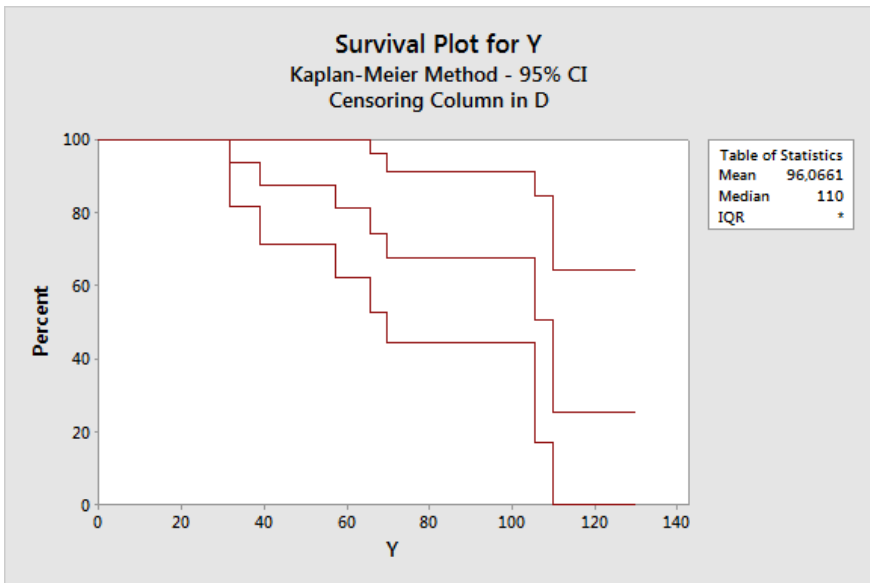


Table 1.2 *Survival times of women with tumours that were negatively or positively stained with HPA.*

Negative staining	Positive staining	
23	5	68
47	8	71
69	10	76*
70*	13	105*
71*	18	107*
100*	24	109*
101*	26	113
148	26	116*
181	31	118
198*	35	143
208*	40	154*
212*	41	162*
224*	48	188*
	50	212*
	59	217*
	61	225*

BREAST CANCER DATA: KM PLOTS (Collett)

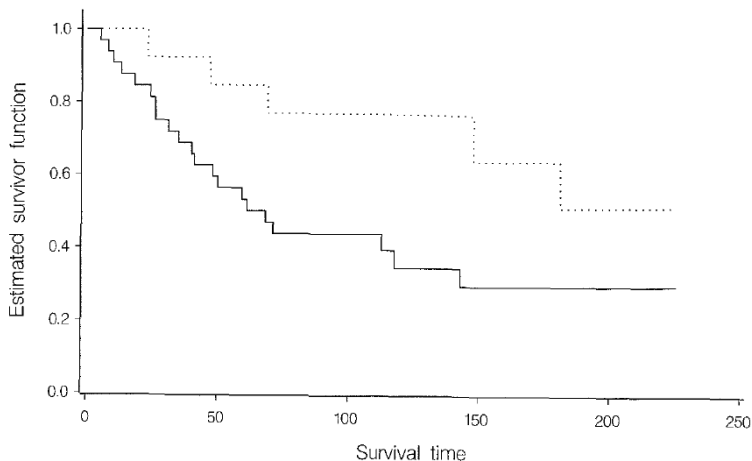


Figure 2.9 *Kaplan-Meier estimate of the survivor functions for women with tumours that were positively stained (—) and negatively stained (···).*

NONPARAMETRIC ESTIMATION OF THE CUMULATIVE HAZARD $A(t)$

Why is an estimate of $A(t)$ useful?

Note first that $A'(t) = \alpha(t)$. Thus,

- T is IFR $\Leftrightarrow \alpha(t)$ is *increasing* $\Leftrightarrow A(t)$ is *convex*
- T is DFR $\Leftrightarrow \alpha(t)$ is *decreasing* $\Leftrightarrow A(t)$ is *concave*

Thus a plot of an estimate $\hat{A}(t)$ can give us information on whether the distribution of T is IFR (*increasing failure rate*) or DFR (*decreasing failure rate*).

More generally, a plot of $\hat{A}(t)$ can give us information on the shape of the hazard rate.

ESTIMATING $A(t)$ BY THE KM-ESTIMATOR

Recall that $S(t) = e^{-A(t)}$, so

$$A(t) = -\ln S(t)$$

Thus, if $\hat{S}_{KM}(t)$ is the KM-estimator for $S(t)$, then we can define,

$$\begin{aligned}\hat{A}_{KM}(t) &= -\ln \hat{S}_{KM}(t) \\ &= -\ln \prod_{T_{(i)} \leq t} \frac{n_i - d_i}{n_i} \\ &= -\sum_{T_{(i)} \leq t} \ln \left(1 - \frac{d_i}{n_i}\right) \\ &\approx \sum_{T_{(i)} \leq t} \frac{d_i}{n_i}\end{aligned}$$

where we used that for small x is

$$-\ln(1 - x) \approx x$$

The Nelson-Aalen estimator (NA-estimator) is simply defined by

$$\hat{A}_{NA}(t) = \sum_{T_{(i)} \leq t} \frac{d_i}{n_i}$$

It can then be shown that its variance can be estimated by

$$\widehat{\text{Var}}(\hat{A}_{NA}(t)) = \sum_{T_{(i)} \leq t} \frac{d_i}{n_i^2}$$

Note: *A more thorough treatment of the Nelson-Aalen and Kaplan-Meier estimators are at the main core of the course – using theory for counting processes and martingales.*

EXAMPLE: NELSON-AALEN ESTIMATOR

Row	C1	C2
1	31,7	1
2	39,2	1
3	57,5	1
4	65,0	0
5	65,8	1
6	70,0	1
7	75,0	0
8	75,2	0
9	87,5	0
10	88,3	0
11	94,2	0
12	101,7	0
13	105,8	1
14	109,2	0
15	110,0	1
16	130,0	0

Row	Time	Numb at risk	1/Numb at risk	Cum Haz Nelson	Survival Nelson
1	31,7	16	0,062500	0,06250	0,939413
2	39,2	15	0,066667	0,12917	0,878827
3	57,5	14	0,071429	0,20060	0,818244
4	65,8	12	0,083333	0,28393	0,752820
5	70,0	11	0,090909	0,37484	0,687401
6	105,8	4	0,250000	0,62484	0,535348
7	110,0	2	0,500000	1,12484	0,324705

Nelson Plot

