

## Survival functions, cumulative hazards, and product integrals: the general case

Uncensored survival time  $T$

Survival function:  $S(t) = P(T > t)$

For the **absolute continuous case**, the hazard is given by:

$$\alpha(t) = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} P(T < t + \Delta t | T \geq t)$$

Cumulative hazard:  $A(t) = \int_0^t \alpha(u) du$

We have the relations:

$$\alpha(t) = A'(t) = -\frac{S'(t)}{S(t)} \quad S(t) = \exp\{-A(t)\}$$

1

For a **general distribution** the hazard rate is not defined, but we may define the cumulative hazard rate as (generalizing the first relation above):

$$A(t) = -\int_0^t \frac{dS(u)}{S(u-)}$$

For the **discrete case**  $A(t)$  is a step function with increments

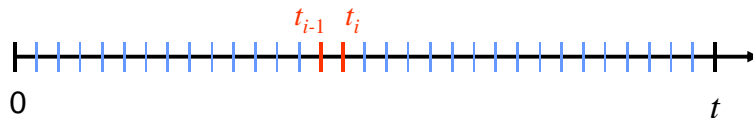
$$\begin{aligned} \Delta A(u) &= -\frac{\Delta S(u)}{S(u-)} \\ &= P(T = u | T \geq u) \end{aligned}$$

How can the second relation above be generalized?

2

Need product-integrals to achieve this generalization

Partition  $[0, t]$  into small time intervals:



$$S(t) = \lim_{\max |t_i - t_{i-1}| \rightarrow 0} \prod (1 - \{A(t_i) - A(t_{i-1})\})$$

$$\stackrel{\text{def}}{=} \mathcal{P}_{0 \leq u \leq t} (1 - dA(u))$$

The limit is a **product-integral**

3

For the **continuous case** we have:

$$\mathcal{P}_{0 \leq u \leq t} (1 - dA(u)) = \exp\{-A(t)\}$$

For the **discrete case** we have:

$$\mathcal{P}_{0 \leq u \leq t} (1 - dA(u)) = \prod_{u \leq t} (1 - \Delta A(u))$$

where  $\Delta A(u) = P(T = u | T \geq u)$  is the increment of the cumulative hazard at time  $u$

For the general case we have a mixture of the two

4

## The Kaplan-Meier estimator

For right censored survival data we observe:

$$\tilde{T}_i = \min\{\text{survival time } T_i, \text{ censoring time } C_i\}$$

$$D_i = I\{\tilde{T}_i = T_i\}$$

**Model:** the uncensored survival times  $T_i$  are *iid* with hazard rate  $\alpha(t)$

Counting and intensity processes:

$$N_i(t) = I\{\tilde{T}_i \leq t, D_i = 1\}$$

$$\lambda_i(t) = I\{\tilde{T}_i \geq t\} \alpha(t) = Y_i(t) \alpha(t)$$

5

Aggregated counting process:

$$N(t) = \sum_{i=1}^n N_i(t)$$

Intensity process:

$$\lambda(t) = \sum_{i=1}^n \lambda_i(t) = Y(t) \alpha(t)$$

with

$$Y(t) = \sum_{i=1}^n I\{\tilde{T}_i \geq t\}$$

the number at risk just before time  $t$

6

Nelson-Aalen estimator:

$$\hat{A}(t) = \int_0^t \frac{dN(u)}{Y(u)} = \sum_{u \leq t} \frac{\Delta N(u)}{Y(u)} = \sum_{T_j \leq t} \frac{1}{Y(T_j)}$$

Plug this into the product-integral expression for the survival function (Nelson-Aalen is a step-function):

$$\begin{aligned} \hat{S}(t) &= \prod_{0 \leq u \leq t} (1 - d\hat{A}(u)) = \prod_{u \leq t} (1 - \Delta \hat{A}(u)) \\ &= \prod_{u \leq t} \left(1 - \frac{\Delta N(u)}{Y(u)}\right) = \prod_{T_j \leq t} \left(1 - \frac{1}{Y(T_j)}\right) \end{aligned}$$

This is the **Kaplan-Meier estimator**

7

## Example 3.8: Second births

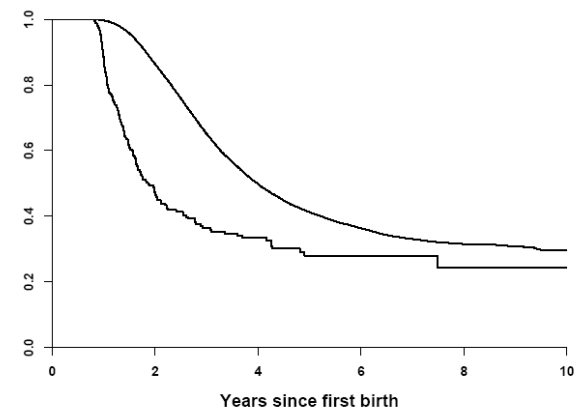


Fig. 3.11 Kaplan-Meier estimates for the time between first and second birth. Upper curve: first child survived one year; lower curve: first child died within one year.

8

### Example 3.9: Third births

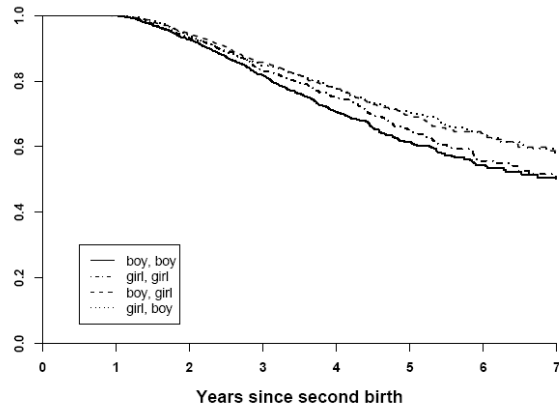


Fig. 3.12 Kaplan-Meier estimates for the time between the second and third births depending on the gender of the two older children.

9

An alternative estimator of the survival function is

$$\begin{aligned} \tilde{S}(t) &= \exp\{-\hat{A}(t)\} \\ &= \exp\left\{-\sum_{T_j \leq t} \frac{1}{Y(T_j)}\right\} \\ &= \prod_{T_j \leq t} \exp\left\{-\frac{1}{Y(T_j)}\right\} \end{aligned}$$

For practical purposes there is little difference between the two estimators

But from a theoretical point of view, the Kaplan-Meier estimator is the natural one (and it may be generalized to Markov models)

10

### Kaplan-Meier estimator: Properties

$$A^*(t) = \int_0^t J(s) \alpha(s) ds \approx A(t)$$

$$S^*(t) = \prod_{0 \leq s \leq t} (1 - dA^*(s)) = \exp\{-A^*(t)\} \approx S(t)$$

May show that (this is Duhamel's equation)

$$\frac{\hat{S}(t)}{S^*(t)} - 1 = - \underbrace{\int_0^t \frac{\hat{S}(s-)}{S^*(s)}}_{\approx 1} d(\hat{A} - A^*)(s) \approx -(\hat{A}(t) - A^*(t))$$

Asymptotically:  $\frac{\hat{S}(t)}{S(t)} - 1 \approx -(\hat{A}(t) - A(t))$

11

Thus:

$$\hat{S}(t) - S(t) \approx -S(t) \cdot (\hat{A}(t) - A(t))$$

The statistical properties for Kaplan-Meier may be derived from those of Nelson-Aalen:

- $\text{Var}\{\hat{S}(t)\} \approx \{S(t)\}^2 \cdot \text{Var}\{\hat{A}(t)\}$
- Variance estimator:  $\hat{\tau}^2(t) = \{\hat{S}(t)\}^2 \cdot \hat{\sigma}^2(t)$   
with  $\hat{\sigma}^2(t) = \int_0^t \{Y(s)\}^{-2} dN(s)$
- $\hat{S}(t)$  is as normally distributed around  $S(t)$

12

Usually the variance is estimated by *Greenwood's formula*:

$$\tilde{\tau}^2(t) = [\hat{S}(t)]^2 \cdot \tilde{\sigma}^2(t)$$

$$\text{with } \tilde{\sigma}^2(t) = \int_0^t [Y(s)\{Y(s) - \Delta N(s)\}]^{-1} dN(s)$$

Only minor difference between the two variance estimators

Pointwise 95% confidence limits for  $S(t)$

**Linear:**  $\hat{S}(t) \pm 1.96 \cdot \hat{S}(t) \cdot \hat{\sigma}(t)$

**Log-log-transformed:**  $\hat{S}(t)^{\exp\{\pm 1.96 \cdot \hat{\sigma}(t) / \log \hat{S}(t)\}}$

(cf. Exercise 3.6)

13

## Estimation of mean survival time

The mean survival time is given by (exercise 1.3)

$$E(T) = \int_0^{\infty} S(u) du$$

Due to censoring, this may usually not be estimated

We may consider the expected survival in  $[0, t]$ :

$$\mu_t = \int_0^t S(u) du$$

This may be estimated by

$$\mu_t = \int_0^t \hat{S}(u) du$$

14

## Estimation of median survival time and other fractiles of the survival distribution

The  $p$ -th fractile  $\xi_p$  of the survival distribution is given by (exercise 1.2)

$$F(\xi_p) = p \quad \text{or equivalently} \quad S(\xi_p) = 1 - p$$

It is estimated by

$$\hat{\xi}_p = \inf \{t : \hat{S}(t) \leq 1 - p\}$$

Confidence intervals may be found by "inverting" the confidence intervals for the survival function (exercise 3.8)

15

### Example 3.10: Median time between first and second births

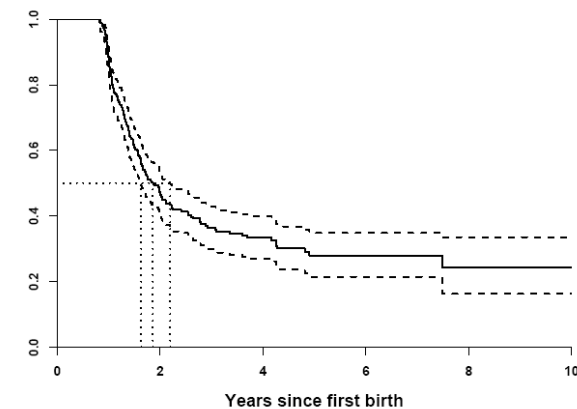


Fig. 3.13 Kaplan-Meier estimate with 95% log-log-transformed confidence intervals for the time between first and second birth for women who lost the first child within one year after birth. It is indicated at the figure how one may obtain the estimated median time with confidence limits.

16