

## S-Plus Commands for Survival Estimation

```
> t_c(1,1,2,2,3,4,4,5,5,8,8,8,8,11,11,12,12,15,17,22,23)

> surv.fit(t,status=rep(1,21))
95 percent confidence interval is of type "log"
time n.risk n.event survival std.dev lower 95% CI upper 95% CI
 1     21      2 0.90476190 0.06405645 0.78753505 1.00000000
 2     19      2 0.80952381 0.08568909 0.65785306 0.9961629
 3     17      1 0.76190476 0.09294286 0.59988048 0.9676909
 4     16      2 0.66666667 0.10286890 0.49268063 0.9020944
 5     14      2 0.57142857 0.10798985 0.39454812 0.8276066
 8     12      4 0.38095238 0.10597117 0.22084536 0.6571327
11      8      2 0.28571429 0.09858079 0.14529127 0.5618552
12      6      2 0.19047619 0.08568909 0.07887014 0.4600116
15      4      1 0.14285714 0.07636035 0.05010898 0.4072755
17      3      1 0.09523810 0.06405645 0.02548583 0.3558956
22      2      1 0.04761905 0.04647143 0.00703223 0.3224544
23      1      1 0.00000000          NA          NA          NA
```

## Estimating the Survival Function

### One-sample nonparametric methods:

We will consider three methods for estimating a survivorship function

$$S(t) = Pr(T \geq t)$$

without resorting to parametric methods:

(1) **Kaplan-Meier**

(2) **Life-table** (Actuarial Estimator)

(3) via the **Cumulative hazard estimator**

## (1) The Kaplan-Meier Estimator

The Kaplan-Meier (or KM) estimator is probably the most popular approach. It can be justified from several perspectives:

- product limit estimator
- likelihood justification
- redistribute to the right estimator

We will start with an intuitive motivation based on conditional probabilities, then review some of the other justifications.

### Motivation:

First, consider an example where there is no censoring.

The following are times of remission (weeks) for 21 leukemia patients receiving control treatment (Table 1.1 of Cox & Oakes):

1, 1, 2, 2, 3, 4, 4, 5, 5, 8, 8, 8, 8, 11, 11, 12, 12, 15, 17, 22, 23

How would we estimate  $S(10)$ , the probability that an individual survives to time 10 or later?

What about  $\tilde{S}(8)$ ? Is it  $\frac{12}{21}$  or  $\frac{8}{21}$ ?

Let's construct a table of  $\tilde{S}(t)$ :

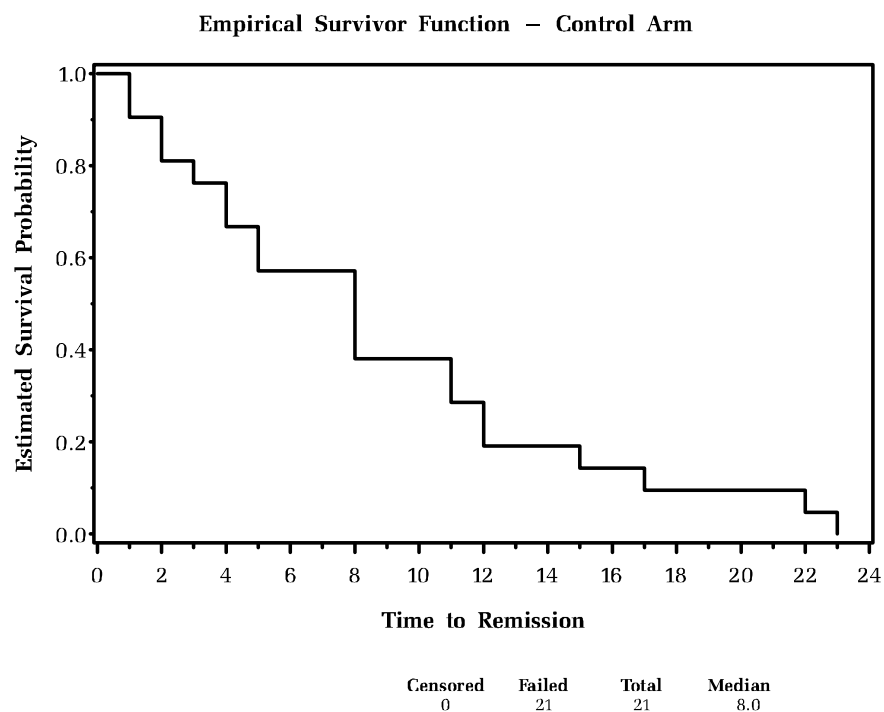
Values of t	$\hat{S}(t)$
$t \leq 1$	$21/21=1.000$
$1 < t \leq 2$	$19/21=0.905$
$2 < t \leq 3$	$17/21=0.809$
$3 < t \leq 4$	
$4 < t \leq 5$	
$5 < t \leq 8$	
$8 < t \leq 11$	
$11 < t \leq 12$	
$12 < t \leq 15$	
$15 < t \leq 17$	
$17 < t \leq 22$	
$22 < t \leq 23$	

## Empirical Survival Function:

When there is no censoring, the general formula is:

$$\tilde{S}(t) = \frac{\# \text{ individuals with } T \geq t}{\text{total sample size}}$$

Example for leukemia data (control arm):



## What if there is censoring?

Consider the treated group from Table 1.1 of Cox and Oakes:

6<sup>+</sup>, 6, 6, 6, 7, 9<sup>+</sup>, 10<sup>+</sup>, 10, 11<sup>+</sup>, 13, 16, 17<sup>+</sup>

19<sup>+</sup>, 20<sup>+</sup>, 22, 23, 25<sup>+</sup>, 32<sup>+</sup>, 32<sup>+</sup>, 34<sup>+</sup>, 35<sup>+</sup>

[Note: times with <sup>+</sup> are right censored]

We know  $S(6) = 21/21$ , because everyone survived at least until time 6 or greater. But, we can't say  $S(7) = 17/21$ , because we don't know the status of the person who was censored at time 6.

In a 1958 paper in the *Journal of the American Statistical Association*, Kaplan and Meier proposed a way to nonparametrically estimate  $S(t)$ , even in the presence of censoring. The method is based on the ideas of **conditional probability**.

**A quick review of conditional probability:**

**Conditional Probability:** Suppose A and B are two events. Then,

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

**Multiplication law of probability:** can be obtained from the above relationship, by multiplying both sides by  $P(B)$ :

$$P(A \cap B) = P(A|B) P(B)$$

**Extension to more than 2 events:**

Suppose  $A_1, A_2 \dots A_k$  are k different events. Then, the probability of all k events happening together can be written as a product of conditional probabilities:

$$\begin{aligned} P(A_1 \cap A_2 \dots \cap A_k) &= P(A_k | A_{k-1} \cap \dots \cap A_1) \times \\ &\quad \times P(A_{k-1} | A_{k-2} \cap \dots \cap A_1) \\ &\quad \dots \\ &\quad \times P(A_2 | A_1) \\ &\quad \times P(A_1) \end{aligned}$$

**Now, let's apply these ideas to estimate  $S(t)$ :**

Suppose  $a_k < t \leq a_{k+1}$ . Then

$$\begin{aligned} S(t) &= P(T \geq a_{k+1}) \\ &= P(T \geq a_1, T \geq a_2, \dots, T \geq a_{k+1}) \\ &= P(T \geq a_1) \times \prod_{j=1}^k P(T \geq a_{j+1} | T \geq a_j) \\ &= \prod_{j=1}^k [1 - P(T = a_j | T \geq a_j)] \\ &= \prod_{j=1}^k [1 - \lambda_j] \end{aligned}$$

$$\begin{aligned} \text{so } \hat{S}(t) &\cong \prod_{j=1}^k \left(1 - \frac{d_j}{r_j}\right) \\ &= \prod_{j:a_j < t} \left(1 - \frac{d_j}{r_j}\right) \end{aligned}$$

$d_j$  is the number of deaths at  $a_j$   
 $r_j$  is the number at risk at  $a_j$

## Intuition behind the Kaplan-Meier Estimator

Think of dividing the observed timespan of the study into a series of fine intervals so that there is a separate interval for each time of death or censoring:



Using the law of conditional probability,

$$Pr(T \geq t) = \prod_j Pr(\text{survive } j\text{-th interval } I_j \mid \text{survived to start of } I_j)$$

where the product is taken over all the intervals including or preceding time  $t$ .

4 possibilities for each interval:

- (1) **No events (death or censoring)** - conditional probability of surviving the interval is 1
- (2) **Censoring** - assume they survive to the end of the interval, so that the conditional probability of surviving the interval is 1
- (3) **Death, but no censoring** - conditional probability of *not* surviving the interval is # deaths (d) divided by # 'at risk' (r) at the beginning of the interval. So the conditional probability of surviving the interval is  $1 - (d/r)$ .
- (4) **Tied deaths and censoring** - assume censorings last to the end of the interval, so that conditional probability of surviving the interval is still  $1 - (d/r)$

### General Formula for $j$ th interval:

It turns out we can write a general formula for the conditional probability of surviving the  $j$ -th interval that holds for all 4 cases:

$$1 - \frac{d_j}{r_j}$$

We could use the same approach by grouping the event times into intervals (say, one interval for each month), and then counting up the number of deaths (events) in each to estimate the probability of surviving the interval (this is called the *lifetable estimate*).

However, the assumption that those censored last until the end of the interval wouldn't be quite accurate, so we would end up with a cruder approximation.

As the intervals get finer and finer, the approximations made in estimating the probabilities of getting through each interval become smaller and smaller, so that the estimator converges to the true  $S(t)$ .

This intuition clarifies why an alternative name for the KM is the product limit estimator.

**The Kaplan-Meier estimator of the survivorship function (or survival probability)  $S(t) = Pr(T \geq t)$  is:**

$$\begin{aligned}\hat{S}(t) &= \prod_{j:\tau_j < t} \frac{r_j - d_j}{r_j} \\ &= \prod_{j:\tau_j < t} \left(1 - \frac{d_j}{r_j}\right)\end{aligned}$$

where

- $\tau_1, \dots, \tau_K$  is the set of  $K$  distinct death times observed in the sample
- $d_j$  is the number of deaths at  $\tau_j$
- $r_j$  is the number of individuals “at risk” right before the  $j$ -th death time (everyone dead or censored at or after that time).
- $c_j$  is the number of censored observations between the  $j$ -th and  $(j + 1)$ -st death times. Censorings tied at  $\tau_j$  are included in  $c_j$

**Note: two useful formulas are:**

$$(1) \quad r_j = r_{j-1} - d_{j-1} - c_{j-1}$$

$$(2) \quad r_j = \sum_{l \geq j} (c_l + d_l)$$

## Calculating the KM - Cox and Oakes example

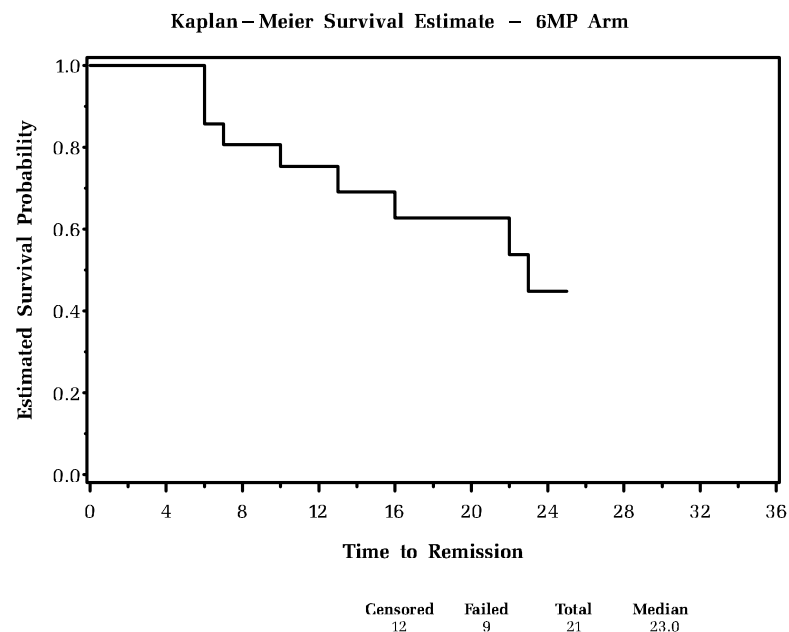
Make a table with a row for every death or censoring time:

$\tau_j$	$d_j$	$c_j$	$r_j$	$1 - (d_j/r_j)$	$\hat{S}(\tau_j^+)$
6	3	1	21	$\frac{18}{21} = 0.857$	
7	1	0	17		
9	0	1	16		
10					
11					
13					
16					
17					
19					
20					
22					
23					

**Note that:**

- $\hat{S}(t^+)$  only changes at death (failure) times
- $\hat{S}(t^+)$  is 1 up to the first death time
- $\hat{S}(t^+)$  only goes to 0 if the last event is a death

## KM plot for treated leukemia patients



**Note:** most statistical software packages summarize the KM survival function at  $\tau_j^+$ , i.e., *just after* the time of the  $j$ -th failure.

**In other words, they provide  $\hat{S}(\tau_j^+)$ .**

When there is no censoring, the empirical survival estimate would then be:

$$\tilde{S}(t^+) = \frac{\# \text{ individuals with } T > t}{\text{total sample size}}$$

## Output from STATA KM Estimator:

failure time: weeks  
failure/censor: remiss

Time	Beg. Total	Fail	Net Lost	Survivor Function	Std. Error	[95% Conf. Int.]	
6	21	3	1	0.8571	0.0764	0.6197	0.9516
7	17	1	0	0.8067	0.0869	0.5631	0.9228
9	16	0	1	0.8067	0.0869	0.5631	0.9228
10	15	1	1	0.7529	0.0963	0.5032	0.8894
11	13	0	1	0.7529	0.0963	0.5032	0.8894
13	12	1	0	0.6902	0.1068	0.4316	0.8491
16	11	1	0	0.6275	0.1141	0.3675	0.8049
17	10	0	1	0.6275	0.1141	0.3675	0.8049
19	9	0	1	0.6275	0.1141	0.3675	0.8049
20	8	0	1	0.6275	0.1141	0.3675	0.8049
22	7	1	0	0.5378	0.1282	0.2678	0.7468
23	6	1	0	0.4482	0.1346	0.1881	0.6801
25	5	0	1	0.4482	0.1346	0.1881	0.6801
32	4	0	2	0.4482	0.1346	0.1881	0.6801
34	2	0	1	0.4482	0.1346	0.1881	0.6801
35	1	0	1	0.4482	0.1346	0.1881	0.6801

## Two Other Justifications for KM Estimator

### I. Likelihood-based derivation (Cox and Oakes)

For a discrete failure time variable, define:

- $d_j$  number of failures at  $a_j$
- $r_j$  number of individuals at risk at  $a_j$   
(including those censored at  $a_j$ ).
- $\lambda_j$  Pr(death) in  $j$ -th interval  
(conditional on survival to start of interval)

The likelihood is that of  $g$  independent binomials:

$$L(\boldsymbol{\lambda}) = \prod_{j=1}^g \lambda_j^{d_j} (1 - \lambda_j)^{r_j - d_j}$$

Therefore, the **maximum likelihood estimator** of  $\lambda_j$  is:

$$\hat{\lambda}_j = d_j / r_j$$

Now we plug in the MLE's of  $\lambda$  to estimate  $S(t)$ :

$$\begin{aligned} \hat{S}(t) &= \prod_{j:a_j < t} (1 - \hat{\lambda}_j) \\ &= \prod_{j:a_j < t} \left(1 - \frac{d_j}{r_j}\right) \end{aligned}$$



## II. Redistribute to the right justification

(Efron, 1967)

In the absence of censoring,  $\hat{S}(t)$  is just the proportion of individuals with  $T \geq t$ . The idea behind Efron's approach is to spread the contributions of censored observations out over all the possible times to their right.

### Algorithm:

- Step (1): arrange the  $n$  observed times (deaths or censorings) in increasing order. If there are ties, put censored after deaths.
- Step (2): Assign weight  $(1/n)$  to each time.
- Step (3): Moving from left to right, each time you encounter a censored observation, distribute its mass to all times to its right.
- Step (4): Calculate  $\hat{S}_j$  by subtracting the final weight for time  $j$  from  $\hat{S}_{j-1}$

## Example of “redistribute to the right” algorithm

Consider the following event times:

2, 2.5+, 3, 3, 4, 4.5+, 5, 6, 7

The algorithm goes as follows:

(Step 1) Times	Step 2	Step 3a	Step 3b	(Step 4) $\hat{S}(\tau_j)$
2	1/9=0.11			0.889
2.5+	1/9=0.11	0		0.889
3	2/9=0.22	0.25		0.635
4	1/9=0.11	0.13		0.508
4.5+	1/9=0.11	0.13	0	0.508
5	1/9=0.11	0.13	0.17	0.339
6	1/9=0.11	0.13	0.17	0.169
7	1/9=0.11	0.13	0.17	0.000

This comes out the same as the product limit approach.

## Properties of the KM estimator

### In the case of no censoring:

$$\hat{S}(t) = \tilde{S}(t) = \frac{\# \text{ deaths at } t \text{ or greater}}{n}$$

where  $n$  is the number of individuals in the study.

This is just like an estimated probability from a binomial distribution, so we have:

$$\hat{S}(t) \simeq \mathcal{N}(S(t), S(t)[1 - S(t)]/n)$$

### How does censoring affect this?

- $\hat{S}(t)$  is still approximately normal
- The mean of  $\hat{S}(t)$  converges to the true  $S(t)$
- The variance is a bit more complicated (since the denominator  $n$  includes some censored observations).

Once we get the variance, then we can construct (pointwise)  $(1 - \alpha)\%$  confidence bands about  $\hat{S}(t)$ :

$$\hat{S}(t) \pm z_{1-\alpha/2} se[\hat{S}(t)]$$

## Greenwood's formula (Collett 2.1.3)

We can think of the KM estimator as

$$\hat{S}(t) = \prod_{j:\tau_j < t} (1 - \hat{\lambda}_j)$$

where  $\hat{\lambda}_j = d_j/r_j$ .

Since the  $\hat{\lambda}_j$ 's are just binomial proportions, we can apply standard likelihood theory to show that each  $\hat{\lambda}_j$  is approximately normal, with mean the true  $\lambda_j$ , and

$$var(\hat{\lambda}_j) \approx \frac{\hat{\lambda}_j(1 - \hat{\lambda}_j)}{r_j}$$

Also, the  $\hat{\lambda}_j$ 's are independent in large enough samples.

Since  $\hat{S}(t)$  is a function of the  $\lambda_j$ 's, we can estimate its variance using the **delta method**:

**Delta method:** If  $Y$  is normal with mean  $\mu$  and variance  $\sigma^2$ , then  $g(Y)$  is approximately normally distributed with mean  $g(\mu)$  and variance  $[g'(\mu)]^2\sigma^2$ .

## Two specific examples of the delta method:

(A)  $Z = \log(Y)$

$$\text{then } Z \sim N \left[ \log(\mu), \left( \frac{1}{\mu} \right)^2 \sigma^2 \right]$$

(B)  $Z = \exp(Y)$

$$\text{then } Z \sim N \left[ e^\mu, [e^\mu]^2 \sigma^2 \right]$$

The examples above use the following results from calculus:

$$\frac{d}{dx} \log u = \frac{1}{u} \left( \frac{du}{dx} \right)$$

$$\frac{d}{dx} e^u = e^u \left( \frac{du}{dx} \right)$$

## Greenwood's formula (continued)

Instead of dealing with  $\hat{S}(t)$  directly, we will look at its log:

$$\log[\hat{S}(t)] = \sum_{j:\tau_j < t} \log(1 - \hat{\lambda}_j)$$

Thus, by approximate independence of the  $\hat{\lambda}_j$ 's,

$$\text{var}(\log[\hat{S}(t)]) = \sum_{j:\tau_j < t} \text{var}[\log(1 - \hat{\lambda}_j)]$$

$$\begin{aligned} \text{by (A)} \quad &= \sum_{j:\tau_j < t} \left( \frac{1}{1 - \hat{\lambda}_j} \right)^2 \text{var}(\hat{\lambda}_j) \\ &= \sum_{j:\tau_j < t} \left( \frac{1}{1 - \hat{\lambda}_j} \right)^2 \hat{\lambda}_j(1 - \hat{\lambda}_j)/r_j \\ &= \sum_{j:\tau_j < t} \frac{\hat{\lambda}_j}{(1 - \hat{\lambda}_j)r_j} \\ &= \sum_{j:\tau_j < t} \frac{d_j}{(r_j - d_j)r_j} \end{aligned}$$

Now,  $\hat{S}(t) = \exp[\log[\hat{S}(t)]]$ . Thus by (B),

$$\text{var}(\hat{S}(t)) = [\hat{S}(t)]^2 \text{var}[\log[\hat{S}(t)]]$$

### Greenwood's Formula:

$$\text{var}(\hat{S}(t)) = [\hat{S}(t)]^2 \sum_{j:\tau_j < t} \frac{d_j}{(r_j - d_j)r_j}$$

## Back to confidence intervals

For a 95% confidence interval, we could use

$$\hat{S}(t) \pm z_{1-\alpha/2} se[\hat{S}(t)]$$

where  $se[\hat{S}(t)]$  is calculated using Greenwood's formula.

**Problem:** This approach can yield values  $> 1$  or  $< 0$ .

**Better approach:** Get a 95% confidence interval for

$$L(t) = \log(-\log(S(t)))$$

Since this quantity is unrestricted, the confidence interval will be in the proper range when we transform back.

**To see why this works, note the following:**

- Since  $\hat{S}(t)$  is an estimated probability

$$0 \leq \hat{S}(t) \leq 1$$

- Taking the log of  $\hat{S}(t)$  has bounds:

$$-\infty \leq \log[\hat{S}(t)] \leq 0$$

- Taking the opposite:

$$0 \leq -\log[\hat{S}(t)] \leq \infty$$

- Taking the log again:

$$-\infty \leq \log[-\log[\hat{S}(t)]] \leq \infty$$

To transform back, reverse steps with  $S(t) = \exp(-\exp(L(t)))$

## Log-log Approach for Confidence Intervals:

(1) Define  $L(t) = \log(-\log(S(t)))$

(2) Form a 95% confidence interval for  $L(t)$  based on  $\hat{L}(t)$ , yielding  $[\hat{L}(t) - A, \hat{L}(t) + A]$

(3) Since  $S(t) = \exp(-\exp(L(t)))$ , the confidence bounds for the 95% CI on  $S(t)$  are:

$$[\exp(-e^{(\hat{L}(t)+A)}), \exp(-e^{(\hat{L}(t)-A)})]$$

(note that the upper and lower bounds switch)

(4) Substituting  $\hat{L}(t) = \log(-\log(\hat{S}(t)))$  back into the above bounds, we get confidence bounds of

$$([\hat{S}(t)]^{e^A}, [\hat{S}(t)]^{e^{-A}})$$

## What is A?

- $A$  is  $1.96 \text{ se}(\hat{L}(t))$

- To calculate this, we need to calculate

$$\text{var}(\hat{L}(t)) = \text{var}[\log(-\log(\hat{S}(t)))]$$

- From our previous calculations, we know

$$\text{var}(\log[\hat{S}(t)]) = \sum_{j:\tau_j < t} \frac{d_j}{(r_j - d_j)r_j}$$

- Applying the delta method as in example (A), we get:

$$\begin{aligned} \text{var}(\hat{L}(t)) &= \text{var}(\log(-\log[\hat{S}(t)])) \\ &= \frac{1}{[\log \hat{S}(t)]^2} \sum_{j:\tau_j < t} \frac{d_j}{(r_j - d_j)r_j} \end{aligned}$$

- We take the square root of the above to get  $\text{se}(\hat{L}(t))$ , and then form the confidence intervals as:

$$\hat{S}(t) e^{\pm 1.96 \text{ se}(\hat{L}(t))}$$

- This is the approach that Stata uses. Splus gives an option to calculate these bounds (use `conf.type='log-log'` in `surv.fit`).

## Summary of Confidence Intervals on $S(t)$

- Calculate  $\hat{S}(t) \pm 1.96 \text{ se}[\hat{S}(t)]$  where  $\text{se}[\hat{S}(t)]$  is calculated using Greenwood's formula, and replace negative lower bounds by 0 and upper bounds greater than 1 by 1.
  - Recommended by Collett
  - This is the default using SAS
  - not very satisfactory
- Use a log transformation to stabilize the variance and allow for non-symmetric confidence intervals. This is what is normally done for the confidence interval of an estimated odds ratio.
  - Use  $\text{var}[\log(\hat{S}(t))] = \sum_{j:\tau_j < t} \frac{d_j}{(r_j - d_j)r_j}$  already calculated as part of Greenwood's formula
  - This is the default in Splus
- Use the log-log transformation just described
  - Somewhat complicated, but always yields proper bounds
  - This is the default in Stata.

## Software for Kaplan-Meier Curves

- Stata - stset and sts commands
- SAS - PROC LIFETEST
- Splus - surv.fit(time,status)

## Defaults for Confidence Interval Calculations

- Stata - “log-log”  $\Rightarrow \hat{L}(t) \pm 1.96 se[\hat{L}(t)]$   
where  $L(t) = \log[-\log(S(t))]$
- SAS - “plain”  $\Rightarrow \hat{S}(t) \pm 1.96 se[\hat{S}(t)]$
- Splus - “log”  $\Rightarrow \log S(t) \pm 1.96 se[\log(\hat{S}(t))]$

but Splus will also give either of the other two options if you request them.

## Stata Commands

Create a file called “leukemia.dat” with the raw data, with a column for treatment, weeks to relapse (i.e., duration of remission), and relapse status:

```
.infile trt remiss status using leukemia.dat

.stset remiss status      (sets up a failure time dataset,
                          with failtime status in that order,
                          type help stset to get details)

.sts list                (estimated S(t), se[S(t)], and 95% CI)

.sts graph, saving(kmtrt) (creates a Kaplan-Meier plot, and
                          saves the plot in file kmtrt.gph,
                          type ‘help gphdot’ to get some
                          printing instructions)

.graph using kmtrt      (redisplays the graph at any later time)
```

If the dataset has already been created and loaded into Stata, then you can substitute the following commands for initializing the data:

```
.use leukem              (finds Stata dataset leukem.dta)

.describe                (provides a description of the dataset)

.stset remiss status     (declares data to be failure type)

.stdes                   (gives a description of the survival dataset)
```

## STATA Output for Treated Leukemia Patients:

```
.use leukem
.stset remiss status if trt==1
.sts list

    failure time:  remiss
failure/censor:  status
```

Time	Beg. Total	Fail	Net Lost	Survivor Function	Std. Error	[95% Conf. Int.]	
6	21	3	1	0.8571	0.0764	0.6197	0.9516
7	17	1	0	0.8067	0.0869	0.5631	0.9228
9	16	0	1	0.8067	0.0869	0.5631	0.9228
10	15	1	1	0.7529	0.0963	0.5032	0.8894
11	13	0	1	0.7529	0.0963	0.5032	0.8894
13	12	1	0	0.6902	0.1068	0.4316	0.8491
16	11	1	0	0.6275	0.1141	0.3675	0.8049
17	10	0	1	0.6275	0.1141	0.3675	0.8049
19	9	0	1	0.6275	0.1141	0.3675	0.8049
20	8	0	1	0.6275	0.1141	0.3675	0.8049
22	7	1	0	0.5378	0.1282	0.2678	0.7468
23	6	1	0	0.4482	0.1346	0.1881	0.6801
25	5	0	1	0.4482	0.1346	0.1881	0.6801
32	4	0	2	0.4482	0.1346	0.1881	0.6801
34	2	0	1	0.4482	0.1346	0.1881	0.6801
35	1	0	1	0.4482	0.1346	0.1881	0.6801

## SAS Commands for Kaplan Meier Estimator - PROC LIFETEST

The SAS command for the Kaplan-Meier estimate is:

```
time failtime*censor(1);
or time failtime*failind(0);
```

The first variable is the failure time, and the second is the failure or censoring indicator. In parentheses you need to put the specific numeric value that corresponds to censoring.

The upper and lower confidence limits on  $\hat{S}(t)$  are included in the data set "OUTSURV" when specified. The upper and lower limits are called: **sdf\_ucl**, **sdf\_lcl**.

```
data leukemia;
    input weeks remiss;
    label weeks='Time to Remission (in weeks)'
          remiss='Remission indicator (1=yes,0=no)';
cards;
6 1
6 1
..... ( lines edited out here)
34 0
35 0
;

proc lifetest data=leukemia outsurv=confint;
    time weeks*remiss(0);
    title 'Leukemia data from Table 1.1 of Cox and Oakes';
run;

proc print data=confint;
title '95% Confidence Intervals for Estimated Survival';
```

## Output from SAS PROC LIFETEST

Note: this information is not printed if you use NOPRINT.

Leukemia data from Table 1.1 of Cox and Oakes

The LIFETEST Procedure

Product-Limit Survival Estimates

WEEKS	Survival	Failure	Survival Standard Error	Number Failed	Number Left
0.0000	1.0000	0	0	0	21
6.0000	.	.	.	1	20
6.0000	.	.	.	2	19
6.0000	0.8571	0.1429	0.0764	3	18
6.0000*	.	.	.	3	17
7.0000	0.8067	0.1933	0.0869	4	16
9.0000*	.	.	.	4	15
10.0000	0.7529	0.2471	0.0963	5	14
10.0000*	.	.	.	5	13
11.0000*	.	.	.	5	12
13.0000	0.6902	0.3098	0.1068	6	11
16.0000	0.6275	0.3725	0.1141	7	10
17.0000*	.	.	.	7	9
19.0000*	.	.	.	7	8
20.0000*	.	.	.	7	7
22.0000	0.5378	0.4622	0.1282	8	6
23.0000	0.4482	0.5518	0.1346	9	5
25.0000*	.	.	.	9	4
32.0000*	.	.	.	9	3
32.0000*	.	.	.	9	2
34.0000*	.	.	.	9	1
35.0000*	.	.	.	9	0

\* Censored Observation

## Output from printing the CONFINT file

95% Confidence Intervals for Estimated Survival

OBS	WEEKS	_CENSOR_	SURVIVAL	SDF_LCL	SDF_UCL
1	0	0	1.00000	1.00000	1.00000
2	6	0	0.85714	0.70748	1.00000
3	6	1	0.85714	.	.
4	7	0	0.80672	0.63633	0.97711
5	9	1	0.80672	.	.
6	10	0	0.75294	0.56410	0.94178
7	10	1	0.75294	.	.
8	11	1	0.75294	.	.
9	13	0	0.69020	0.48084	0.89955
10	16	0	0.62745	0.40391	0.85099
11	17	1	0.62745	.	.
12	19	1	0.62745	.	.
13	20	1	0.62745	.	.
14	22	0	0.53782	0.28648	0.78915
15	23	0	0.44818	0.18439	0.71197
16	25	1	.	.	.
17	32	1	.	.	.
18	32	1	.	.	.
19	34	1	.	.	.
20	35	1	.	.	.

The output dataset will have one observation for each unique combination of WEEKS and \_CENSOR\_. It will also add an observation for failure time equal to 0.



## Splus Commands

Create a file called “leukemia.dat” with the variables names in the first row, as follows:

```
t      c
6      1
6      1
etc ...
```

In Splus, type

```
y_read.table('leukemia.dat',header=T)
surv.fit(y$t,y$c)
plot(surv.fit(y$t,y$c))
```

(the plot command will also yield 95% confidence intervals)

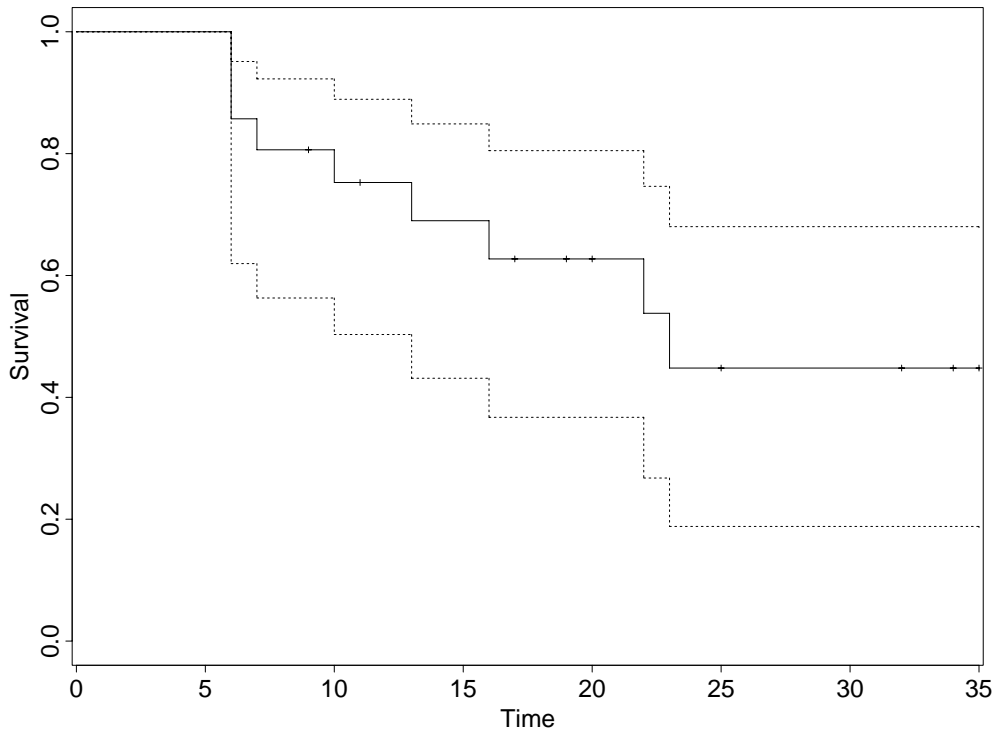
To specify the type of confidence intervals, use the `conf.type=` option in the `surv.fit` statements: e.g. `conf.type=“log-log”` or `conf.type=“plain”`

```
>surv.fit(y$t,y$c)
95 percent confidence interval is of type "log"
time n.risk n.event survival std.dev lower 95% CI upper 95% CI
 6      21      3 0.8571429 0.07636035 0.7198171 1.0000000
 7      17      1 0.8067227 0.08693529 0.6531242 0.9964437
10      15      1 0.7529412 0.09634965 0.5859190 0.9675748
13      12      1 0.6901961 0.10681471 0.5096131 0.9347692
16      11      1 0.6274510 0.11405387 0.4393939 0.8959949
22      7       1 0.5378151 0.12823375 0.3370366 0.8582008
23      6       1 0.4481793 0.13459146 0.2487882 0.8073720
```

```
> surv.fit(y$t,y$c,conf.type="log-log")
95 percent confidence interval is of type "log-log"
time n.risk n.event survival std.dev lower 95% CI upper 95% CI
 6      21      3 0.8571429 0.07636035 0.6197180 0.9515517
 7      17      1 0.8067227 0.08693529 0.5631466 0.9228090
10      15      1 0.7529412 0.09634965 0.5031995 0.8893618
13      12      1 0.6901961 0.10681471 0.4316102 0.8490660
16      11      1 0.6274510 0.11405387 0.3675109 0.8049122
22      7       1 0.5378151 0.12823375 0.2677789 0.7467907
23      6       1 0.4481793 0.13459146 0.1880520 0.6801426
```

```
> surv.fit(y$t,y$c,conf.type="plain")
95 percent confidence interval is of type "plain"
time n.risk n.event survival std.dev lower 95% CI upper 95% CI
 6      21      3 0.8571429 0.07636035 0.7074793 1.0000000
 7      17      1 0.8067227 0.08693529 0.6363327 0.9771127
10      15      1 0.7529412 0.09634965 0.5640993 0.9417830
13      12      1 0.6901961 0.10681471 0.4808431 0.8995491
16      11      1 0.6274510 0.11405387 0.4039095 0.8509924
22      7       1 0.5378151 0.12823375 0.2864816 0.7891487
23      6       1 0.4481793 0.13459146 0.1843849 0.7119737
```

## KM Survival Estimate and Confidence intervals (SPlus)



## Means, Medians, Quantiles based on the KM

- **Mean:**  $\sum_{j=1}^k \tau_j Pr(T = \tau_j)$
- **Median** - by definition, this is the time,  $\tau$ , such that  $S(\tau) = 0.5$ . However, in practice, it is defined as the smallest time such that  $\hat{S}(\tau) \leq 0.5$ . The median is more appropriate for censored survival data than the mean.

For the treated leukemia patients, we find:

$$\hat{S}(22) = 0.5378$$

$$\hat{S}(23) = 0.4482$$

The median is thus 23. This can also be seen visually on the graph to the left.

- **Lower quartile (25<sup>th</sup> percentile):**  
the smallest time (LQ) such that  $\hat{S}(LQ) \leq 0.75$
- **Upper quartile (75<sup>th</sup> percentile):**  
the smallest time (UQ) such that  $\hat{S}(UQ) \leq 0.25$

## The (2) Lifetable Estimator of Survival:

We said that we would consider the following three methods for estimating a survivorship function

$$S(t) = Pr(T \geq t)$$

without resorting to parametric methods:

(1) ✓ **Kaplan-Meier**

(2)  $\implies$  **Life-table** (Actuarial Estimator)

(3)  $\implies$  **Cumulative hazard estimator**

## (2) The Lifetable or Actuarial Estimator

- one of the oldest techniques around
- used by actuaries, demographers, etc.
- **applies when the data are grouped**

Our goal is still to estimate the survival function, hazard, and density function, but this is complicated by the fact that we don't know exactly when during each time interval an event occurs.

Lee (section 4.2) provides a good description of lifetable methods, and distinguishes several types according to the data sources:

### POPULATION LIFE TABLES

- **cohort life table** - describes the mortality experience from birth to death for a particular cohort of people born at about the same time. People at risk at the start of the interval are those who survived the previous interval.
- **current life table** - constructed from (1) census information on the number of individuals alive at each age, for a given year and (2) vital statistics on the number of deaths or failures in a given year, by age. This type of lifetable is often reported in terms of a hypothetical cohort of 100,000 people.

Generally, censoring is not an issue for Population Life Tables.

CLINICAL LIFE TABLES - applies to grouped survival data from studies in patients with specific diseases. Because patients can enter the study at different times, or be lost to follow-up, censoring must be allowed.

### **Notation**

- the  $j$ -th time interval is  $[t_{j-1}, t_j)$
- $c_j$  - the number of censorings in the  $j$ -th interval
- $d_j$  - the number of failures in the  $j$ -th interval
- $r_j$  is the number entering the interval

**Example:** 2418 Males with Angina Pectoris (Lee, p.91)

Year after Diagnosis	$j$	$d_j$	$c_j$	$r_j$	$r'_j = r_j - c_j/2$
[0, 1)	1	456	0	2418	2418.0
[1, 2)	2	226	39	1962	1942.5 (1962 - $\frac{39}{2}$ )
[2, 3)	3	152	22	1697	1686.0
[3, 4)	4	171	23	1523	1511.5
[4, 5)	5	135	24	1329	1317.0
[5, 6)	6	125	107	1170	1116.5
[6, 7)	7	83	133	938	871.5
etc..					

## Estimating the survivorship function

We could apply the K-M formula directly to the numbers in the table on the previous page, estimating  $S(t)$  as

$$\hat{S}(t) = \prod_{j:\tau_j < t} \left(1 - \frac{d_j}{r_j}\right)$$

However, this approach is unsatisfactory for grouped data.... it treats the problem as though it were in discrete time, with events happening only at 1 yr, 2 yr, etc. In fact, what we are trying to calculate here is the conditional probability of dying within the interval, given survival to the beginning of it.

## What should we do with the censored people?

### We can assume that censorings occur:

- at the beginning of each interval:  $r'_j = r_j - c_j$
- at the end of each interval:  $r'_j = r_j$
- on average halfway through the interval:

$$r'_j = r_j - c_j/2$$

The last assumption yields the Actuarial Estimator. It is appropriate if censorings occur uniformly throughout the interval.

## Constructing the lifetable

First, some additional notation for the  $j$ -th interval,  $[t_{j-1}, t_j)$ :

- **Midpoint** ( $t_{mj}$ ) - useful for plotting the density and the hazard function
- **Width** ( $b_j = t_j - t_{j-1}$ ) needed for calculating the hazard in the  $j$ -th interval

### Quantities estimated:

- Conditional probability of dying

$$\hat{q}_j = d_j/r'_j$$

- Conditional probability of surviving

$$\hat{p}_j = 1 - \hat{q}_j$$

- Cumulative probability of surviving at  $t_j$ :

$$\begin{aligned}\hat{S}(t_j) &= \prod_{\ell \leq j} \hat{p}_\ell \\ &= \prod_{\ell \leq j} \left(1 - \frac{d_\ell}{r'_\ell}\right)\end{aligned}$$

### Some important points to note:

- Because the intervals are defined as  $[t_{j-1}, t_j)$ , the first interval typically starts with  $t_0 = 0$ .
- Stata estimates the survival function at the right-hand endpoint of each interval, i.e.,  $S(t_j)$
- However, SAS estimates the survival function at the left-hand endpoint,  $S(t_{j-1})$ .
- The implication in SAS is that  $\hat{S}(t_0) = 1$  and  $\hat{S}(t_1) = p_1$

### Other quantities estimated at the midpoint of the $j$ -th interval:

- **Hazard** in the  $j$ -th interval:

$$\begin{aligned}\hat{\lambda}(t_{mj}) &= \frac{d_j}{b_j(r'_j - d_j/2)} \\ &= \frac{\hat{q}_j}{b_j(1 - \hat{q}_j/2)}\end{aligned}$$

the number of deaths in the interval divided by the average number of survivors at the midpoint

- **density** at the midpoint of the  $j$ -th interval:

$$\begin{aligned}\hat{f}(t_{mj}) &= \frac{\hat{S}(t_{j-1}) - \hat{S}(t_j)}{b_j} \\ &= \frac{\hat{S}(t_{j-1}) \hat{q}_j}{b_j}\end{aligned}$$

Note: Another way to get this is:

$$\begin{aligned}\hat{f}(t_{mj}) &= \hat{\lambda}(t_{mj})\hat{S}(t_{mj}) \\ &= \hat{\lambda}(t_{mj})[\hat{S}(t_j) + \hat{S}(t_{j-1})]/2\end{aligned}$$

## Constructing the Lifetable using Stata

Uses the `ltable` command.

If the raw data are already grouped, then the `freq` statement must be used when reading the data.

```
. infile years status count using angina.dat
(32 observations read)
```

```
. ltable years status [freq=count]
```

Interval	Beg. Total	Deaths	Lost	Survival	Std. Error	[95% Conf. Int.]
0	1	2418	456	0	0.8114	0.0080 0.7952 0.8264
1	2	1962	226	39	0.7170	0.0092 0.6986 0.7346
2	3	1697	152	22	0.6524	0.0097 0.6329 0.6711
3	4	1523	171	23	0.5786	0.0101 0.5584 0.5981
4	5	1329	135	24	0.5193	0.0103 0.4989 0.5392
5	6	1170	125	107	0.4611	0.0104 0.4407 0.4813
6	7	938	83	133	0.4172	0.0105 0.3967 0.4376
7	8	722	74	102	0.3712	0.0106 0.3505 0.3919
8	9	546	51	68	0.3342	0.0107 0.3133 0.3553
9	10	427	42	64	0.2987	0.0109 0.2775 0.3201
10	11	321	43	45	0.2557	0.0111 0.2341 0.2777
11	12	233	34	53	0.2136	0.0114 0.1917 0.2363
12	13	146	18	33	0.1839	0.0118 0.1614 0.2075
13	14	95	9	27	0.1636	0.0123 0.1404 0.1884
14	15	59	6	23	0.1429	0.0133 0.1180 0.1701
15	16	30	0	30	0.1429	0.0133 0.1180 0.1701

It is also possible to get estimates of the hazard function,  $\hat{\lambda}_j$ , and its standard error using the “`hazard`” option:

```
. ltable years status [freq=count], hazard
```

Interval	Beg. Total	Cum. Failure	Std. Error	Hazard	Std. Error	[95% Conf Int]
0	1	2418	0.1886	0.0080	0.2082	0.0097 0.1892 0.2272
1	2	1962	0.2830	0.0092	0.1235	0.0082 0.1075 0.1396
2	3	1697	0.3476	0.0097	0.0944	0.0076 0.0794 0.1094
3	4	1523	0.4214	0.0101	0.1199	0.0092 0.1020 0.1379
4	5	1329	0.4807	0.0103	0.1080	0.0093 0.0898 0.1262
5	6	1170	0.5389	0.0104	0.1186	0.0106 0.0978 0.1393
6	7	938	0.5828	0.0105	0.1000	0.0110 0.0785 0.1215
7	8	722	0.6288	0.0106	0.1167	0.0135 0.0902 0.1433
8	9	546	0.6658	0.0107	0.1048	0.0147 0.0761 0.1336
9	10	427	0.7013	0.0109	0.1123	0.0173 0.0784 0.1462
10	11	321	0.7443	0.0111	0.1552	0.0236 0.1090 0.2015
11	12	233	0.7864	0.0114	0.1794	0.0306 0.1194 0.2395
12	13	146	0.8161	0.0118	0.1494	0.0351 0.0806 0.2182
13	14	95	0.8364	0.0123	0.1169	0.0389 0.0407 0.1931
14	15	59	0.8571	0.0133	0.1348	0.0549 0.0272 0.2425
15	16	30	0.8571	0.0133	0.0000	. . .

There is also a “`failure`” option which gives the number of failures (like the default), and also provides a 95% confidence interval on the cumulative failure probability.

## Constructing the lifetable using SAS

If the raw data are already grouped, then the FREQ statement must be used when reading the data.

SAS requires that the interval endpoints be specified, using one of the following (see SAS manual or online help for more detail):

- **intervals** - specify the the interval endpoints
- **width** - specify the width of each interval
- **ninterval** - specify the number of intervals

```
Title 'Actuarial Estimator for Angina Pectoris Example';
data angina;
  input years status count;
cards;
0.5 1 456
1.5 1 226
2.5 1 152          /* angina cases */
3.5 1 171
4.5 1 135
5.5 1 125
.
.
0.5 0 0
1.5 0 39
2.5 0 22          /* censored */
3.5 0 23
4.5 0 24
5.5 0 107
.
.
proc lifetest data=angina outsurv=survres intervals=0 to 15 by 1 method=act;
  time years*status(0);
  freq count;
```

## SAS output:

Actuarial Estimator for Angina Pectoris Example

The LIFETEST Procedure

Life Table Survival Estimates

Interval [Lower, Upper)	Number Failed	Number Censored	Effective Sample Size	Conditional Probability of Failure	Conditional Probability Standard Error	
0	1	456	0	2418.0	0.1886	0.00796
1	2	226	39	1942.5	0.1163	0.00728
2	3	152	22	1686.0	0.0902	0.00698
3	4	171	23	1511.5	0.1131	0.00815
4	5	135	24	1317.0	0.1025	0.00836
5	6	125	107	1116.5	0.1120	0.00944
6	7	83	133	871.5	0.0952	0.00994
7	8	74	102	671.0	0.1103	0.0121
8	9	51	68	512.0	0.0996	0.0132
9	10	42	64	395.0	0.1063	0.0155
10	11	43	45	298.5	0.1441	0.0203
11	12	34	53	206.5	0.1646	0.0258
12	13	18	33	129.5	0.1390	0.0304
13	14	9	27	81.5	0.1104	0.0347
14	15	6	23	47.5	0.1263	0.0482
15	.	0	30	15.0	0	0

Interval [Lower, Upper)	Survival	Failure	Survival Standard Error	Median Residual Lifetime	Median Standard Error
0	1	1.0000	0	5.3313	0.1749
1	2	0.8114	0.1886	6.2499	0.2001
2	3	0.7170	0.2830	6.3432	0.2361
3	4	0.6524	0.3476	6.2262	0.2361
4	5	0.5786	0.4214	6.2185	0.1853
5	6	0.5193	0.4807	5.9077	0.1806
6	7	0.4611	0.5389	5.5962	0.1855
7	8	0.4172	0.5828	5.1671	0.2713
8	9	0.3712	0.6288	4.9421	0.2763
9	10	0.3342	0.6658	4.8258	0.4141
10	11	0.2987	0.7013	4.6888	0.4183
11	12	0.2557	0.7443	.	.
12	13	0.2136	0.7864	.	.
13	14	0.1839	0.8161	.	.
14	15	0.1636	0.8364	.	.
15	.	0.1429	0.8571	.	.



**more SAS output:** (estimated density  $\hat{f}_j$  and hazard  $\hat{\lambda}_j$ )

Evaluated at the Midpoint of the Interval

Interval [Lower, Upper)		PDF	PDF Standard Error	Hazard	Hazard Standard Error
0	1	0.1886	0.00796	0.208219	0.009698
1	2	0.0944	0.00598	0.123531	0.008201
2	3	0.0646	0.00507	0.09441	0.007649
3	4	0.0738	0.00543	0.119916	0.009154
4	5	0.0593	0.00495	0.108043	0.009285
5	6	0.0581	0.00503	0.118596	0.010589
6	7	0.0439	0.00469	0.1	0.010963
7	8	0.0460	0.00518	0.116719	0.013545
8	9	0.0370	0.00502	0.10483	0.014659
9	10	0.0355	0.00531	0.112299	0.017301
10	11	0.0430	0.00627	0.155235	0.023602
11	12	0.0421	0.00685	0.17942	0.030646
12	13	0.0297	0.00668	0.149378	0.03511
13	14	0.0203	0.00651	0.116883	0.038894
14	15	0.0207	0.00804	0.134831	0.054919
15	.	.	.	.	.

Summary of the Number of Censored and Uncensored Values

Total	Failed	Censored	%Censored
2418	1625	793	32.7957

Suppose we wish to use the actuarial method, but the data do not come grouped.

Consider the treated nursing home patients, with length of stay (los) grouped into 100 day intervals:

```
.use nurshome

.drop if rx==0                (keep only the treated patients)
(881 observations deleted)

.stset los fail

.ltable los fail, intervals(100)
```

Interval	Beg.	Total Deaths	Lost	Survival	Std. Error	[95% Conf. Int.]
0	100	710	328	0	0.5380	0.0187 0.5006 0.5739
100	200	382	86	0	0.4169	0.0185 0.3805 0.4529
200	300	296	65	0	0.3254	0.0176 0.2911 0.3600
300	400	231	38	0	0.2718	0.0167 0.2396 0.3050
400	500	193	32	1	0.2266	0.0157 0.1966 0.2581
500	600	160	13	0	0.2082	0.0152 0.1792 0.2388
600	700	147	13	0	0.1898	0.0147 0.1619 0.2195
700	800	134	10	30	0.1739	0.0143 0.1468 0.2029
800	900	94	4	29	0.1651	0.0143 0.1383 0.1941
900	1000	61	4	30	0.1508	0.0147 0.1233 0.1808
1000	1100	27	0	27	0.1508	0.0147 0.1233 0.1808

## SAS Commands for lifetable analysis - grouping data

```

Title 'Actuarial Estimator for nursing home data';
data morris ;
  infile 'ch12.dat' ;
  input los age trt gender marstat hltstat cens ;

data morristr;
  set morris;
  if trt=1;

proc lifetest data=morristr outsurv=survres
  intervals=0 to 1100 by 100 method=act;
  time los*cens(1);
run ;

proc print data=survres;
run;

```

## Actuarial estimator for treated nursing home patients

Actuarial Estimator for Nursing Home Patients

The LIFETEST Procedure

Life Table Survival Estimates

Interval [Lower, Upper)	Number Failed	Number Censored	Effective Sample Size	Conditional Probability of Failure
0 100	330	0	712.0	0.4635
100 200	86	0	382.0	0.2251
200 300	65	0	296.0	0.2196
300 400	38	0	231.0	0.1645
400 500	32	1	192.5	0.1662
500 600	13	0	160.0	0.0813
600 700	13	0	147.0	0.0884
700 800	10	30	119.0	0.0840
800 900	4	29	79.5	0.0503
900 1000	4	30	46.0	0.0870
1000 1100	0	27	13.5	0

Interval [Lower, Upper)	Conditional Probability Standard Error	Survival	Failure	Survival Standard Error	Median Residual Lifetime
0 100	0.0187	1.0000	0	0	130.2
100 200	0.0214	0.5365	0.4635	0.0187	306.2
200 300	0.0241	0.4157	0.5843	0.0185	398.8
300 400	0.0244	0.3244	0.6756	0.0175	617.0
400 500	0.0268	0.2711	0.7289	0.0167	.
500 600	0.0216	0.2260	0.7740	0.0157	.
600 700	0.0234	0.2076	0.7924	0.0152	.
700 800	0.0254	0.1893	0.8107	0.0147	.
800 900	0.0245	0.1734	0.8266	0.0143	.
900 1000	0.0415	0.1647	0.8353	0.0142	.
1000 1100	0	0.1503	0.8497	0.0147	.

Actuarial estimator for treated nursing home patients, cont'd

Evaluated at the Midpoint  
of the Interval

Interval [Lower, Upper)	Median Standard Error	PDF PDF	PDF Standard Error	Hazard Hazard	Hazard Standard Error
0	100	15.5136	0.00463	0.000187	0.006033
100	200	30.4597	0.00121	0.000122	0.002537
200	300	65.7947	0.000913	0.000108	0.002467
300	400	74.5466	0.000534	0.000084	0.001792
400	500	.	0.000451	0.000078	0.001813
500	600	.	0.000184	0.00005	0.000847
600	700	.	0.000184	0.00005	0.000925
700	800	.	0.000159	0.00005	0.000877
800	900	.	0.000087	0.000043	0.000516
900	1000	.	0.000143	0.00007	0.000909
1000	1100	.	0	.	0

Summary of the Number of Censored and Uncensored Values

Total	Failed	Censored	%Censored
712	595	117	16.4326

Actuarial estimator for treated nursing home patients, cont'd  
Output from SURVRES dataset

Actuarial Estimator for Nursing Home Patients

OBS	LOS	SURVIVAL	SDF_LCL	SDF_UCL	MIDPOINT	PDF
1	0	1.00000	1.00000	1.00000	50	.0046348
2	100	0.53652	0.49989	0.57315	150	.0012079
3	200	0.41573	0.37953	0.45193	250	.0009129
4	300	0.32444	0.29005	0.35883	350	.0005337
5	400	0.27107	0.23842	0.30372	450	.0004506
6	500	0.22601	0.19528	0.25674	550	.0001836
7	600	0.20764	0.17783	0.23745	650	.0001836
8	700	0.18928	0.16048	0.21808	750	.0001591
9	800	0.17337	0.14536	0.20139	850	.0000872
10	900	0.16465	0.13677	0.19253	950	.0001432
11	1000	0.15033	0.12157	0.17910	1050	.0000000

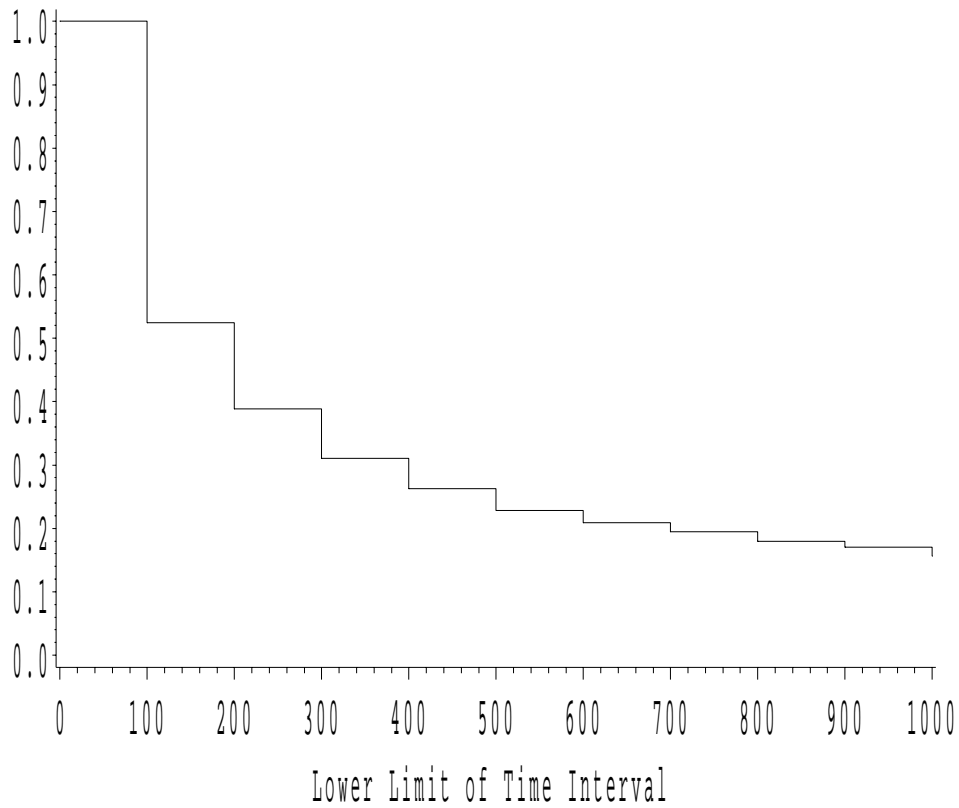
  

OBS	PDF_LCL	PDF_UCL	HAZARD	HAZ_LCL	HAZ_UCL
1	.0042685	.0050011	.0060329	.0054123	.0066535
2	.0009685	.0014472	.0025369	.0020050	.0030687
3	.0007014	.0011245	.0024668	.0018717	.0030619
4	.0003686	.0006988	.0017925	.0012248	.0023601
5	.0002981	.0006031	.0018130	.0011874	.0024386
6	.0000847	.0002825	.0008469	.0003869	.0013069
7	.0000847	.0002825	.0009253	.0004228	.0014277
8	.0000617	.0002565	.0008772	.0003340	.0014203
9	.0000027	.0001717	.0005161	.0000105	.0010218
10	.0000069	.0002794	.0009091	.0000191	.0017991
11	.	.	.0000000	.	.

**Examples for Nursing home data:**

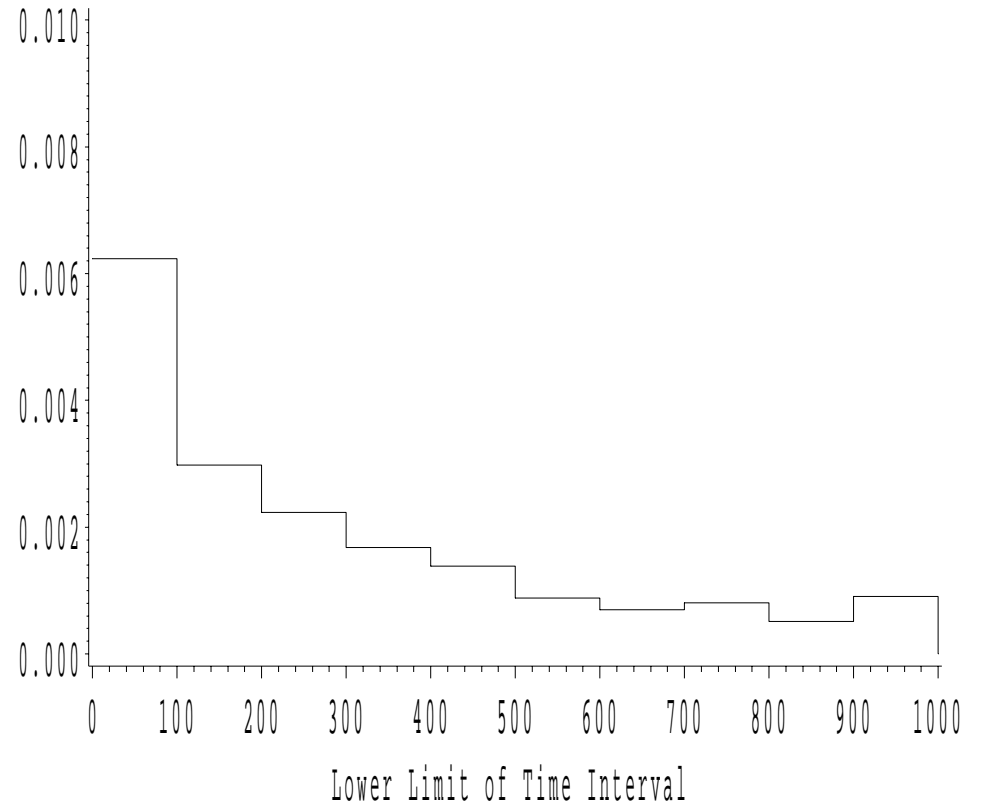
**Estimated Survival:**

**Duration of stay in nursing homes**  
Estimated Survival



**Estimated hazard:**

**Duration of stay in nursing homes**  
Estimated hazard



### (3) Estimating the cumulative hazard

(Nelson-Aalen estimator)

Suppose we want to estimate  $\Lambda(t) = \int_0^t \lambda(u) du$ , the cumulative hazard at time  $t$ .

Just as we did for the KM, think of dividing the observed timespan of the study into a series of fine intervals so that there is only one event per interval:

			D		C		C	D	D	D
--	--	--	---	--	---	--	---	---	---	---

$\Lambda(t)$  can then be approximated by a sum:

$$\hat{\Lambda}(t) = \sum_j \lambda_j \Delta$$

where the sum is over intervals,  $\lambda_j$  is the value of the hazard in the  $j$ -th interval and  $\Delta$  is the width of each interval. Since  $\hat{\lambda}\Delta$  is approximately the probability of dying in the interval, we can further approximate by

$$\hat{\Lambda}(t) = \sum_j d_j / r_j$$

It follows that  $\Lambda(t)$  will change only at death times, and hence we write the Nelson-Aalen estimator as:

$$\hat{\Lambda}_{NA}(t) = \sum_{j:\tau_j < t} d_j / r_j$$

			D		C		C	D	D	D
$r_j$	n	n	n	n-1	n-1	n-2	n-2	n-3	n-4	
$d_j$	0	0	1	0	0	0	0	1	1	
$c_j$	0	0	0	0	1	0	1	0	0	
$\hat{\lambda}(t_j)$	0	0	1/n	0	0	0	0	$\frac{1}{n-3}$	$\frac{1}{n-4}$	
$\hat{\Lambda}(t_j)$	0	0	1/n	1/n	1/n	1/n	1/n			

Once we have  $\hat{\Lambda}_{NA}(t)$ , we can also find another estimator of  $S(t)$  (Fleming-Harrington):

$$\hat{S}_{FH}(t) = \exp(-\hat{\Lambda}_{NA}(t))$$

In general, this estimator of the survival function will be close to the Kaplan-Meier estimator,  $\hat{S}_{KM}(t)$

We can also go the other way ... we can take the Kaplan-Meier estimate of  $S(t)$ , and use it to calculate an alternative estimate of the cumulative hazard function:

$$\hat{\Lambda}_{KM}(t) = -\log \hat{S}_{KM}(t)$$

## Stata commands for FH Survival Estimate

Say we want to obtain the Fleming-Harrington estimate of the survival function for married females, in the healthiest initial subgroup, who are randomized to the untreated group of the nursing home study.

First, we use the following commands to calculate the Nelson-Aalen cumulative hazard estimator:

```
. use nurshome

. keep if rx==0 & gender==0 & health==2 & married==1
(1579 observations deleted)

. sts list, na

      failure _d:  fail
analysis time _t:  los
```

Time	Beg. Total	Fail	Net Lost	Nelson-Aalen Cum. Haz.	Std. Error	[95% Conf. Int.]	
14	12	1	0	0.0833	0.0833	0.0117	0.5916
24	11	1	0	0.1742	0.1233	0.0435	0.6976
25	10	1	0	0.2742	0.1588	0.0882	0.8530
38	9	1	0	0.3854	0.1938	0.1438	1.0326
64	8	1	0	0.5104	0.2306	0.2105	1.2374
89	7	1	0	0.6532	0.2713	0.2894	1.4742
113	6	1	0	0.8199	0.3184	0.3830	1.7551
123	5	1	0	1.0199	0.3760	0.4952	2.1006
149	4	1	0	1.2699	0.4515	0.6326	2.5493
168	3	1	0	1.6032	0.5612	0.8073	3.1840
185	2	1	0	2.1032	0.7516	1.0439	4.2373
234	1	1	0	3.1032	1.2510	1.4082	6.8384

After generating the Nelson-Aalen estimator, we manually have to create a variable for the survival estimate:

```
. sts gen nelson=na

. gen sfh=exp(-nelson)

. list sfh
```

```
              sfh
1.  .9200444
2.  .8400932
3.  .7601478
4.  .6802101
5.  .6002833
6.  .5203723
7.  .4404857
8.  .3606392
9.  .2808661
10. .2012493
11. .1220639
12. .0449048
```

Additional built-in functions can be used to generate 95% confidence intervals on the FH survival estimate.

We can compare the Fleming-Harrington survival estimate to the KM estimate by rerunning the `sts list` command:

```
. sts list

. sts gen skm=s

. list skm sfh
```

	skm	sfh
1.	.91666667	.9200444
2.	.83333333	.8400932
3.	.75	.7601478
4.	.66666667	.6802101
5.	.58333333	.6002833
6.	.5	.5203723
7.	.41666667	.4404857
8.	.33333333	.3606392
9.	.25	.2808661
10.	.16666667	.2012493
11.	.08333333	.1220639
12.	0	.0449048

In this example, it looks like the Fleming-Harrington estimator is slightly higher than the KM at every time point, but with larger datasets the two will typically be much closer.

## Splus Commands for Fleming-Harrington Estimator:

(Nursing home data: females, untreated, married, healthy)

### Fleming-Harrington:

```
>fh<-surv.fit(los,cens,type="f",conf.type="log-log")
>fh
```

```
95 percent confidence interval is of type "log-log"
time n.risk n.event survival std.dev lower 95% CI upper 95% CI
 14    12      1 0.9200444 0.08007959 0.5244209125 0.9892988
 24    11      1 0.8400932 0.10845557 0.4750041174 0.9600371
 25    10      1 0.7601478 0.12669130 0.4055610500 0.9200425
 38     9      1 0.6802101 0.13884731 0.3367907188 0.8724502
 64     8      1 0.6002833 0.14645413 0.2718422278 0.8187596
 89     7      1 0.5203723 0.15021856 0.2115701242 0.7597900
113     6      1 0.4404857 0.15045450 0.1564397006 0.6960354
123     5      1 0.3606392 0.14723033 0.1069925657 0.6278888
149     4      1 0.2808661 0.14043303 0.0640979523 0.5560134
168     3      1 0.2012493 0.12990589 0.0293208029 0.4827590
185     2      1 0.1220639 0.11686728 0.0058990525 0.4224087
234     1      1 0.0449048 0.06216787 0.0005874321 0.2740658
```

### Kaplan-Meier:

```
>km<-surv.fit(los,cens,conf.type="log-log")
>km
```

```
95 percent confidence interval is of type "log-log"
time n.risk n.event survival std.dev lower 95% CI upper 95% CI
 14    12      1 0.91666667 0.07978559 0.538977181 0.9878256
 24    11      1 0.83333333 0.10758287 0.481714942 0.9555094
 25    10      1 0.75000000 0.12500000 0.408415913 0.9117204
 38     9      1 0.66666667 0.13608276 0.337018933 0.8597118
 64     8      1 0.58333333 0.14231876 0.270138924 0.8009402
 89     7      1 0.50000000 0.14433757 0.208477143 0.7360731
113     6      1 0.41666667 0.14231876 0.152471264 0.6653015
123     5      1 0.33333333 0.13608276 0.102703980 0.5884189
149     4      1 0.25000000 0.12500000 0.060144556 0.5047588
168     3      1 0.16666667 0.10758287 0.026510427 0.4129803
185     2      1 0.08333333 0.07978559 0.005052835 0.3110704
234     1      1 0.00000000          NA          NA          NA
```