



Max Planck Odense Center on the Biodemography of Aging
University of Southern Denmark

Non-parametric Survival Analysis

Adam Lenart

MaxO - SDU

The survival function $S(t)$

The survival function at time t is the probability of surviving longer than time t or, alternatively, the probability of failing after time t .

$$S(t) = P(T > t)$$

where T is the time to event of interest (for example, time until death after cancer diagnosis).

The Kaplan-Meier estimator $\hat{S}(t)$

The Kaplan-Meier estimator is a non-parametric estimate of the survival function from the observed data, using information about censoring.

Example:

A cohort of 12 cancer patients followed up from diagnosis until death: 7 of them die, 5 drop out of the study (so their observations are right censored).

Ordered times: 5, 6*, 7, 8, 9*, 12*, 14, 15, 16, 20*, 22*, 23

The Kaplan-Meier estimator $\hat{S}(t)$

id	time	status
1	5	1
2	6	0
3	7	1
4	8	1
5	9	0
6	12	0
7	14	1
8	15	1
9	16	1
10	20	0
11	22	0
12	23	1

The Kaplan-Meier estimator $\hat{S}(t)$

$$S(12)=P(\text{survive in } [0,12])$$

$$=P(\text{survive in } [0,5]) \cdot$$

$$P(\text{survive in } [5,6] | \text{survive}[0,5]) \cdot$$

$$P(\text{survive in } [6,7] | \text{survive}[5,6]) \cdot$$

$$P(\text{survive in } [7,8] | \text{survive}[6,7]) \cdot$$

$$P(\text{survive in } [8,9] | \text{survive}[7,8]) \cdot$$

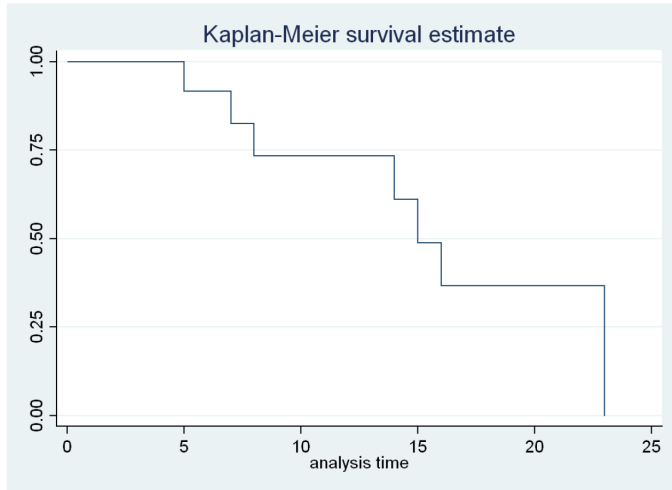
$$P(\text{survive in } [9,12] | \text{survive}[8,9])$$

Probabilities of surviving in an interval conditional on having survived until that interval.

The Kaplan-Meier estimator $\hat{S}(t)$

$(t-1, t)$	at risk	failed	censored	$p(t-1, t)$	$\hat{S}(t)$
(0,5]	12	1	0	$1 - \frac{1}{12}$	$\frac{11}{12} = 0.92$
(5,6]	11	0	1	1	$\frac{11}{12} \cdot 1 = 0.92$
(6,7]	10	1	0	$1 - \frac{1}{10}$	$\frac{11}{12} \cdot 1 \cdot \frac{9}{10} = 0.82$
(7,8]	9	1	0	$1 - \frac{1}{9}$	$\frac{11}{12} \cdot 1 \cdot \frac{9}{10} \cdot \frac{8}{9} = 0.73$
(8,9]	8	0	1	1	$\frac{11}{12} \cdot 1 \cdot \frac{9}{10} \cdot \frac{8}{9} \cdot 1 = 0.73$
(9,12]	7	0	1	1	$\frac{11}{12} \cdot 1 \cdot \frac{9}{10} \cdot \frac{8}{9} \cdot 1 \cdot 1 = 0.73$
(12,14]	6	1	0
(14,15]	5	1	0
(15,16]	4	1	0
(16,20]	3	0	1		
(20,22]	2	0	1		
(22,23]	1	1	0		

The Kaplan-Meier estimator $\hat{S}(t)$



The Kaplan-Meier estimator $\hat{S}(t)$

Kaplan-Meier formula:

$$\hat{S}(t) = \prod_{t_i \leq t} \left(1 - \frac{1}{Y(t_i)} \right)$$

with death times t_1, \dots, t_d and $Y(t_i)$ = number at risk just before t_i .

For multiple deaths at the same time (also called **ties**):

$$\hat{S}(t) = \prod_{t_i \leq t} \left(1 - \frac{m_i}{Y(t_i)} \right)$$

with m_i = number of deaths at time t_i .

Kaplan-Meier is also known as the **product limit** estimate.

Kaplan-Meier in Stata

The first step for any type of survival analysis in Stata (so not only for Kaplan-Meier) is to **stset** the data, to define them as **survival time data**.

Basic syntax of the command: `stset time var, failure var`

```
. list
```

	id	status	time
1.	1	1	5
2.	2	0	6
3.	3	1	7
.	.	.	.
11.	11	0	22
12.	12	1	23

```
. stset time, failure(status=1)
```

The `failure` option allows us to specify if an observation is right censored.

Kaplan-Meier in Stata

Once you have prepared the data for survival analysis with `stset`, the syntax for the Kaplan-Meier command in Stata is very easy:

for one curve for the whole population:

`sts list` estimates the survival function

`sts graph` plots the survival function

for curves by subgroups:

`sts list, by(varname)`

`sts graph, by(varname)`

There is no need to specify the data, because it always refers to the last survival-time dataset available in the memory.

Kaplan-Meier in Stata

You can also use the drop-down menu:

(Users\zarull\Documents\Teaching\SurvivalAnalysis\Data\heroindata.dta - [Results])

Graphics Statistics User Window Help

Statistics

- Summaries, tables, and tests
- Linear models and related
- Binary outcomes
- Ordinal outcomes
- Categorical outcomes
- Count outcomes
- Generalized linear models
- Treatment effects
- Endogenous covariates
- Sample-selection models
- Exact statistics
- Nonparametric analysis
- Time series
- Multivariate time series
- Longitudinal/panel data
- Multilevel mixed-effects models
- Survival analysis
 - Setup and utilities
 - Regression models
 - Summary statistics, tests, and tables
 - Graphs
 - Power and sample size
- Epidemiology and related
- SEM (structural equation modeling)
- Survey data analysis
- Multiple imputation
- Multivariate analysis
- Power and sample size
- Resampling
- Postestimation
- Other

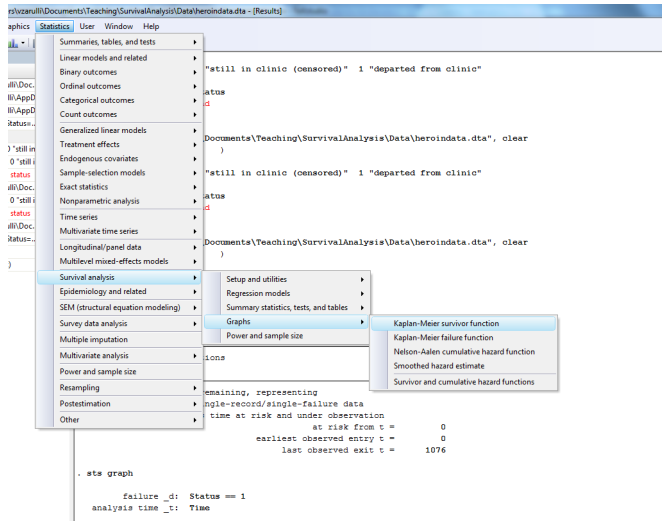
Summary statistics, tests, and tables

- Summarize survival-time data
- Describe survival-time data
- Report incidence-rate comparison
- Tabulate Mantel-Haenszel rate ratios
- Tabulate Mantel-Cox rate ratios
- Person-time, incidence rates, and SMR
- Tabulate failure rates and rate ratios
- Create survivor, hazard, and other variables
- List survivor and cumulative hazard functions
- Test equality of survivor functions
- Life tables for survival data
- CIs for means and percentiles of survival time

	Survivor	Std.	[95% Conf. Int.]	
Function	Error			
2	1.0000	-	-	-
0	0.9958	0.0042	0.9703	0.9994
0	0.9915	0.0060	0.9665	0.9979
0	0.9873	0.0073	0.9611	0.9959
0	0.9831	0.0084	0.9555	0.9936
0	0.9788	0.0094	0.9499	0.9911
2	0.9788	0.0094	0.9499	0.9911
0	0.9745	0.0103	0.9442	0.9885
0	0.9703	0.0111	0.9386	0.9857
0	0.9660	0.0118	0.9331	0.9828
			0.9223	0.9769
			0.9170	0.9738
			n. 0.9170	n. 0.9738

Kaplan-Meier in Stata

You can also use the drop-down menu:



Exercise 1

Exercise 1:

Heroindata.dta contain information about times, in days, that heroin addicts spend in a clinic. The variables are:

- ▶ ID (individual identification number)
- ▶ Clinic (1 or 2)
- ▶ Status (0 = still in clinic at end of study (censored) or 1 = departed from clinic)
- ▶ Time (days spent in clinic)
- ▶ Prison (1 = prison record or 0 = no record)
- ▶ Dose (methadone dosage (mg/day))

Exercise 1

Tasks:

1. Open the data "heroindata.dta".
2. Prepare the data for the analysis by declaring them survival-time data.
3. Estimate and plot the Kaplan-Meier curve for the whole population of drug addicts.
4. Estimate and plot the clinic specific Kaplan-Meier curve.
5. Add censoring time to the previous plot with the `censored(number)` option.

Logrank test

To compare the survival curves of groups by testing their equality (H_0 : Equality of the survival functions).

Comparison of the overall survival curve by comparing, at each failure time, the expected versus observed number of failures for each group and then combining the comparisons over all observed failure times.

Logrank test

For every time failure t_j

time t_j	failed	survived	at risk
group 1	f_{j1}	$r_{j1} - f_{j1}$	r_{j1}
group 2	f_{j2}	$r_{j2} - f_{j2}$	r_{j2}
total	f_j	$r_j - f_j$	r_j

under the validity of H_0 : $E_{j1} = f_j \frac{r_{j1}}{r_j}$ and $V_j = \frac{r_{j1} r_{j2} f_j (r_j - f_j)}{r_j^2 (r_j - 1)}$

Test statistic:

$$LR = \frac{(O_1 - E_1)^2}{V} \sim \chi_1^2$$

where

$E_1 = \sum_j E_{j1}$ expected number of failures in group 1

$O_1 = \sum_j f_{j1}$ observed failures in group 1

$V = \sum_j V_j$ variance

Logrank test in Stata

Logrank test:

```
sts test varname
```

Stratified logrank test:

```
sts test varname, strata(varname) detail
```

When, while testing the equality of the survival curves of Group 1 and Group 2, we think that an additional variable, for example Sex, is worth to be taken into account (because males and females survival experiences differ...)

Logrank test in Stata

The interpretation of the test follows the usual hypothesis testing: the p-value tells us whether we can or can not reject the null hypothesis.

Log-rank test for equality of survivor functions

Clinic	Events observed	Events expected
1	122	90.91
2	28	59.09
Total	150	150.00
chi2 (1) =		27.89
Pr>chi2 =		0.0000

Exercise 2

Tasks:

1. Test the equality of the survival curves of drug addicts by Clinic.
2. It may be that those who have been in prison and those who have not been have different survival patterns. This may be a confounding factor in the analysis of the performance of the clinic. Plot Kaplan-Meier curves by Clinic and Prison. What is your impression?
3. Perform a logrank test between survival curves of the two clinic, stratified by the variable Prison, to confirm the Kaplan-Meier result.

The hazard and cumulative hazard functions

The hazard function, $h(t)$ indicates the risk of death at the instant t , given survival until t .

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

The cumulative hazard function, $H(t)$ is:

$$H(t) = \int_0^t h(x) dx$$

The survival function and the hazard functions are related by the following relation:

$$S(t) = e^{-\int_0^t h(x) dx} = e^{-H(t)}$$

The Nelson-Aalen estimator

In virtue of their relation, we could derive the cumulative hazard and the hazard from the survival function in the following way:

$$H(t) = -\ln S(t) \text{ and } h(t) = -\frac{d}{dt} \ln S(t)$$

However, when $S(t)$ is a step function and the sample size is rather small, a more appropriate method is the Nelson-Aalen estimator of the cumulative hazard $\hat{H}(t)$.

With death times t_1, \dots, t_d :

$$\hat{H}(t) = \sum_{t_i \leq t} \frac{m_i}{Y(t_i)}$$

where $Y(t_i)$ = number at risk just before t_i and m_i = number of deaths at time t_i .

The Nelson-Aalen estimator in Stata

Cumulative hazard estimate:

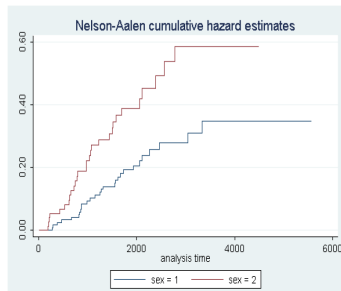
```
sts list, cumhaz
```

```
sts list, by(varname) cumhaz
```

Cumulative hazard plot:

```
sts graph, cumhaz
```

```
sta graph, by(varname) cumhaz
```



The Nelson-Aalen estimator in Stata

And for the hazard function?

When the survival is a step function it can not be directly differentiated to obtain the hazard.

A non-straightforward way to estimate it is by taking the steps of the Nelson-Aalen cumulative hazard and smoothing them with some smoothing algorithm.

In Stata you need to specify the option `hazard` in the `sts graph` and which kernel function (epanechnikov, gaussian...) and bandwidth to use for the smoothing.

```
sts graph, hazard kernel(...) width(...)
```

Exercise 3

Tasks:

1. Estimate the cumulative hazard of leaving the clinic for the drug addicts by clinic.
2. Plot the cumulative hazard by clinic, the survival function by clinic and save the plots in the memory.
(just add `saving(nameplot)` to the `sts graph` command)
3. Combine the plots in a single graph.
(use `gr combine "nameplot1" "nameplot2"`).

Exercise 3

Tasks:

4. Estimate the median survival time (time at which 50% of the individuals have failed and 50% is still surviving, also known as the 50th percentile) of the drug addicts in the two clinics (`stci, by(varname)`). Why the median survival time is not available for clinic 2?
5. Estimate the time at which $\frac{1}{4}$ of the patients have left the clinics (use `stci, by(varname) p(#)`, where `#` is the desired percentile).

Exercise 3

Tasks:

6. List the comparison of the survival function by clinic at 2 weeks, 1 month, 3 months, 6 months, 1 year, 2 years (`sts list, by(varname), at(#)`, where # is the desired survival time).