



Survival Analysis

Module II: Leaves and Trees

Dr. Mark Williamson
DaCCoTA
University of North Dakota

Introduction

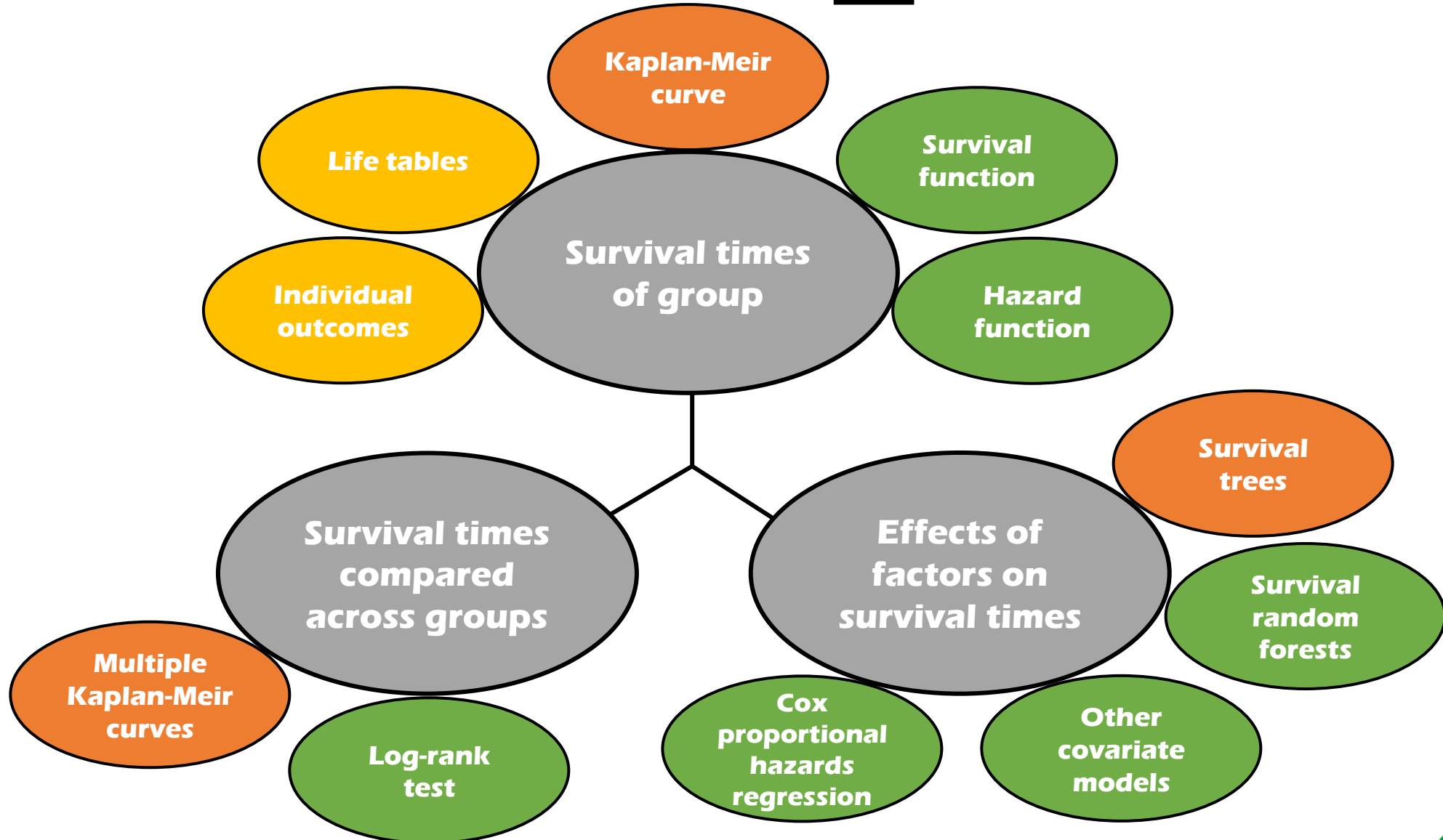


- Last time, we soared over a general overview of survival analysis and played around with some basic code examples.

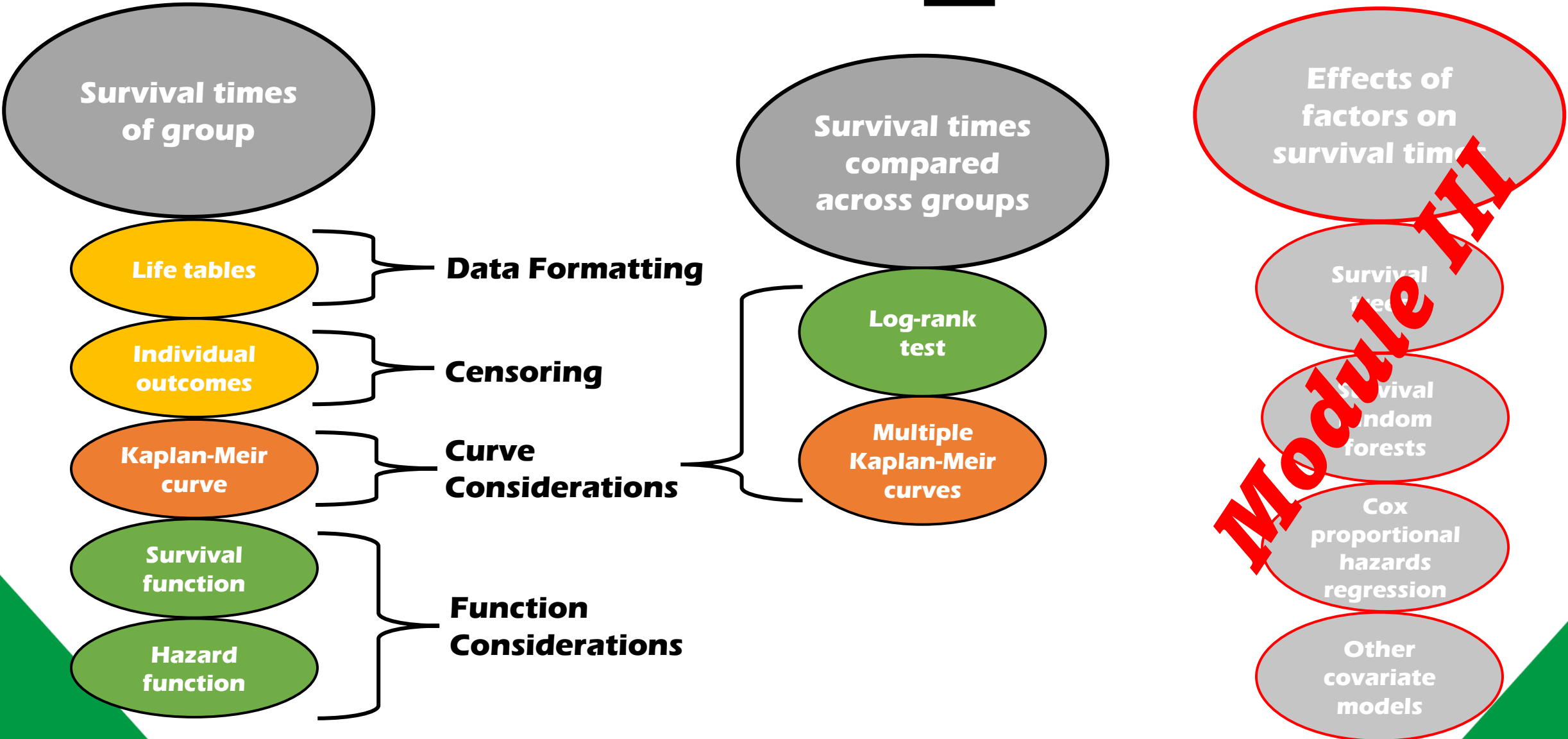
- Today, we'll dig our roots in deeper into the details.
- We will also look at more detailed examples in R and SAS.



Scope



Scope



Data Formatting



Basic information needed

- Time
- Status (Event or 'Censoring')
- Group
- Covariates

Getting data in the right format

- Change dates to time
- Clarify status (0/1 vs 1/2, etc.)

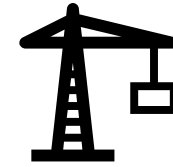
Calculating values for life tables

- N -> # of subjects at risk
- 1-D/N -> proportion survived at interval
- S(t) -> survival function

$$= S(t_{i-1}) * (1 - D_{i-1}/N_{i-1})$$

Name	ID	Time	Status	Group	Weight
Blue Anklo	1	26	1	H	120
Red Anklo	2	26	1	H	150
Green Anklo	3	31	1	H	200
Green Tricer	4	18	1	H	180
Black Iguan	5	38	1	H	88
Grey Hadr	6	13	1	H	73
Green Sauro	7	30	1	H	334
Orange Pter	8	26	1	C	12
Orange Carno	9	50	1	C	56
Blue_Red_Carno	10	27	1	C	93
Grey_Carno	11	16	1	C	86
Light_Blue_Carno	12	41	1	C	102

Data Formatting



Basic information needed

- Time
- Status (Event or 'Censoring')
- Group
- Covariates

Getting data in the right format

- Change dates to time
- Clarify status (0/1 vs 1/2, etc.)

Calculating values for life tables

- N -> # of subjects at risk
- 1-D/N -> proportion survived at interval
- S(t) -> survival function

$$= S(t_{i-1}) * (1 - D_{i-1}/N_{i-1})$$

Time	Dead	N	1-D/N	S(t)
0		12		1
13	1	12	0.916667	0.916667
16	1	11	0.909091	0.833333
18	1	10	0.9	0.75
26	3	9	0.666667	0.5
27	1	6	0.833333	0.416667
30	1	5	0.8	0.333333
31	1	4	0.75	0.25
38	1	3	0.666667	0.166667
41	1	2	0.5	0.083333
50	1	1	0	0

Censoring

What is censoring?

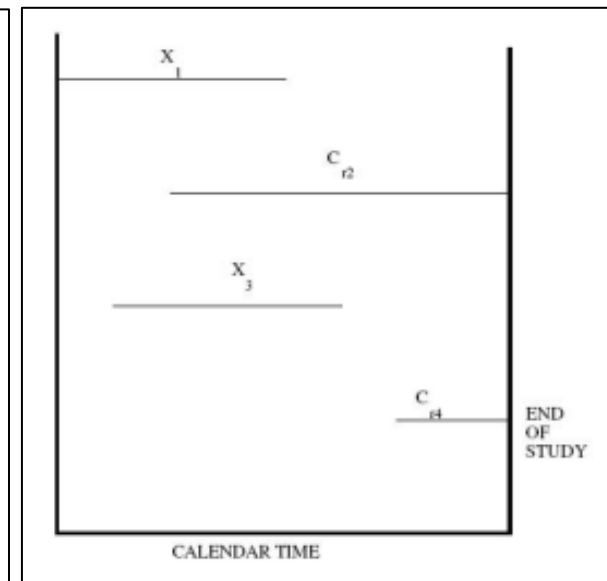
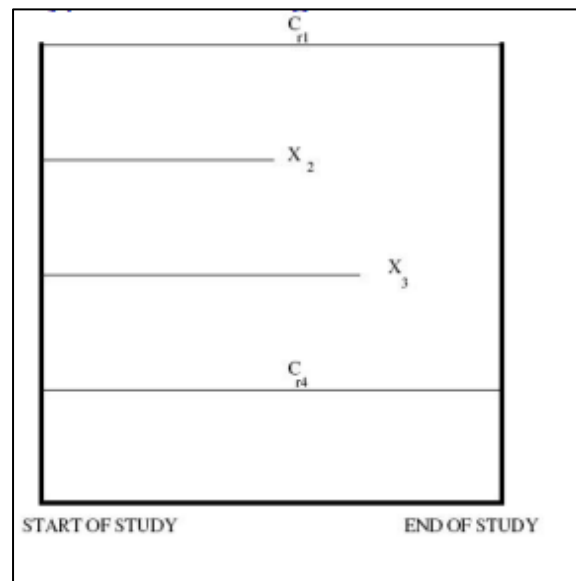
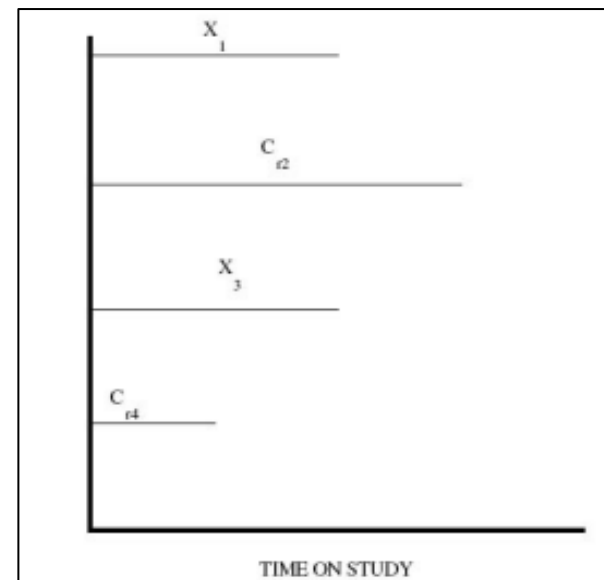
- Full information on exact event time not present [1]
- AKA, event is not recorded during the time period
- Unrealistic to have enough time to capture all events
- Typically **right censored**
- After 1st patient censored, survival becomes **estimate** [2]
- Informativeness of censoring needs to be considered [3,4]

Reasons for censoring [3]

- Loss to follow-up
- Withdrawal from study
- No event by end of study period

Types of Censoring [1]

- Fixed type I censoring
- Random type I censoring
- Type II censoring
- Censoring can also be generalized [5]



Function Considerations

Generating survival and hazard estimates

- When discrete, equations are iterative and used to produce non-parametric curves
- Each can be derived from the others
- In the absence of censoring, survival data could be analyzed with an empirical cumulative distribution function
- Otherwise, Kaplan-Meier estimator is common -> non-parametric and 'choppy'
- Can also use maximum likelihood (ML) to estimate via a distribution -> parametric and 'smooth'

Comparing survival estimates

- AKA, telling the difference between curves
- Standard is Log Rank Test [6]
- Others include Alternative Log Rank, Wilcoxon, and Taron-Ware Test

Type	Discrete function
Survival	$S(t_i) = S(t_{i-1}) * (1 - D_{i-1}/N_{i-1})$
Hazard	$D_i=1: h(t_i) = D_i/N_i$ $D_i>1: h(t_i) = 1 / (D_i - N_i + 1)$
Cumulative Hazard	$H(t_i) = -\ln(1 - D_i/N_i)$

$$h(t) = -\frac{\partial \log(S(t))}{\partial t}$$

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

$$S(t) = \exp(-H(t))$$

Function Considerations

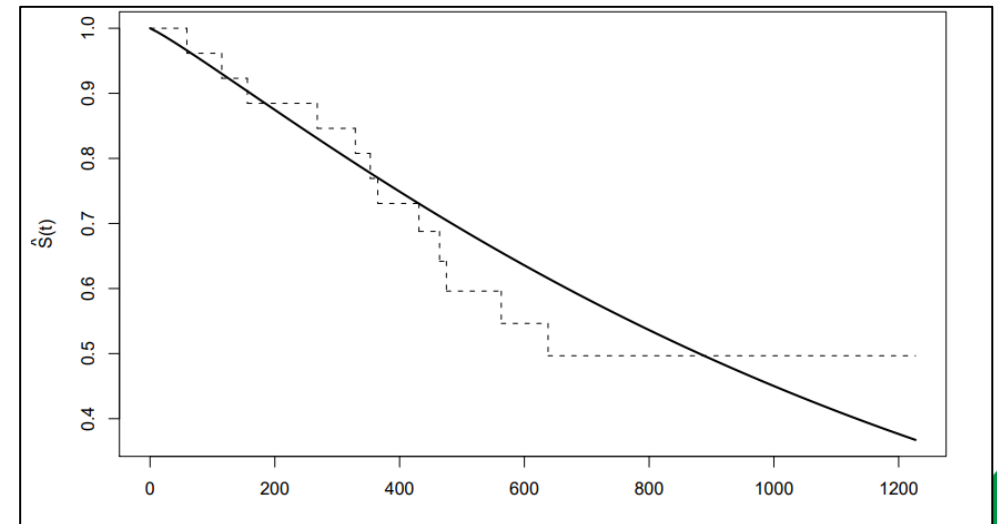
Generating survival and hazard estimates

- When discrete, equations are iterative and used to produce non-parametric curves
- Each can be derived from the others
- In the absence of censoring, survival data could be analyzed with an empirical cumulative distribution function
- Otherwise, Kaplan-Meier estimator is common -> non-parametric and 'choppy'
- Can also use maximum likelihood (ML) to estimate via a distribution -> parametric and 'smooth'

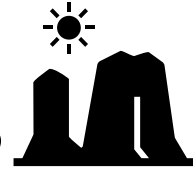
Comparing survival estimates

- AKA, telling the difference between curves
- Standard is Log Rank Test [6]
- Others include Alternative Log Rank, Wilcoxon, and Taron-Ware Test

Type	Discrete function
Survival	$S(t_i) = S(t_{i-1}) * (1 - D_{i-1}/N_{i-1})$
Hazard	$D_i=1: h(t_i) = D_i/N_i$ $D_i>1: h(t_i) = 1 / (D_i - N_i + 1)$
Cumulative Hazard	$H(t_i) = -\ln(1 - D_i/N_i)$



Function Considerations

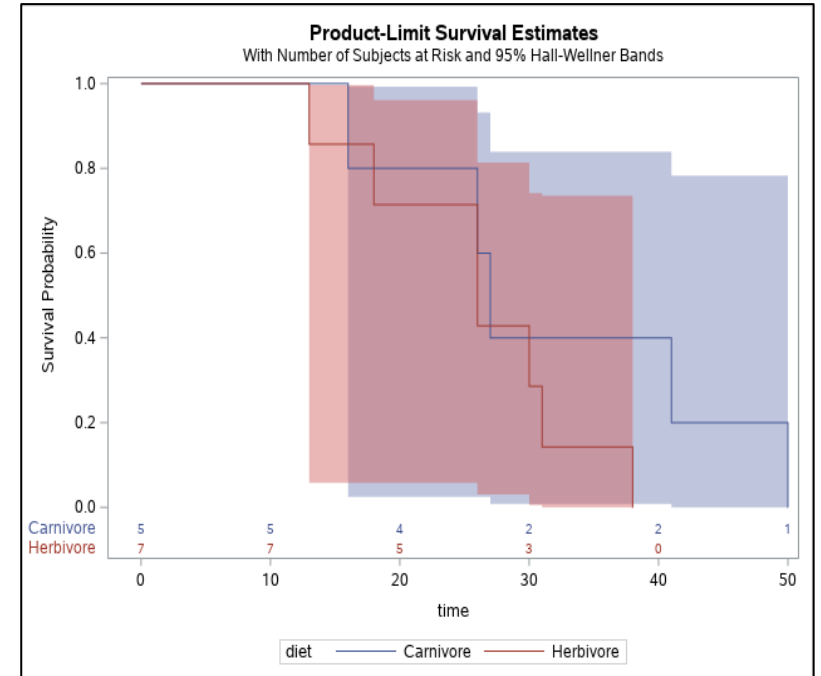


Generating survival and hazard estimates

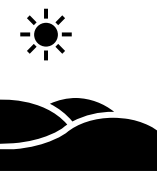
- When discrete, equations are iterative and used to produce non-parametric curves
- Each can be derived from the others
- In the absence of censoring, survival data could be analyzed with an empirical cumulative distribution function
- Otherwise, Kaplan-Meier estimator is common -> non-parametric and 'choppy'
- Can also use maximum likelihood (ML) to estimate via a distribution -> parametric and 'smooth'

Comparing survival estimates

- AKA, telling the difference between curves
- Standard is Log Rank Test [6]
- Others include Alternative Log Rank, Wilcoxon, and Taron-Ware Test



Test of Equality over Strata			
Test	Chi-Square	DF	Pr > Chi-Square
Log-Rank	1.2420	1	0.2651
Wilcoxon	0.5226	1	0.4697
Tarone	0.7964	1	0.3722
Peto	0.5061	1	0.4768
Modified Peto	0.4459	1	0.5043
Fleming(1)	0.5226	1	0.4697



Curve Considerations

Kaplan-Meier

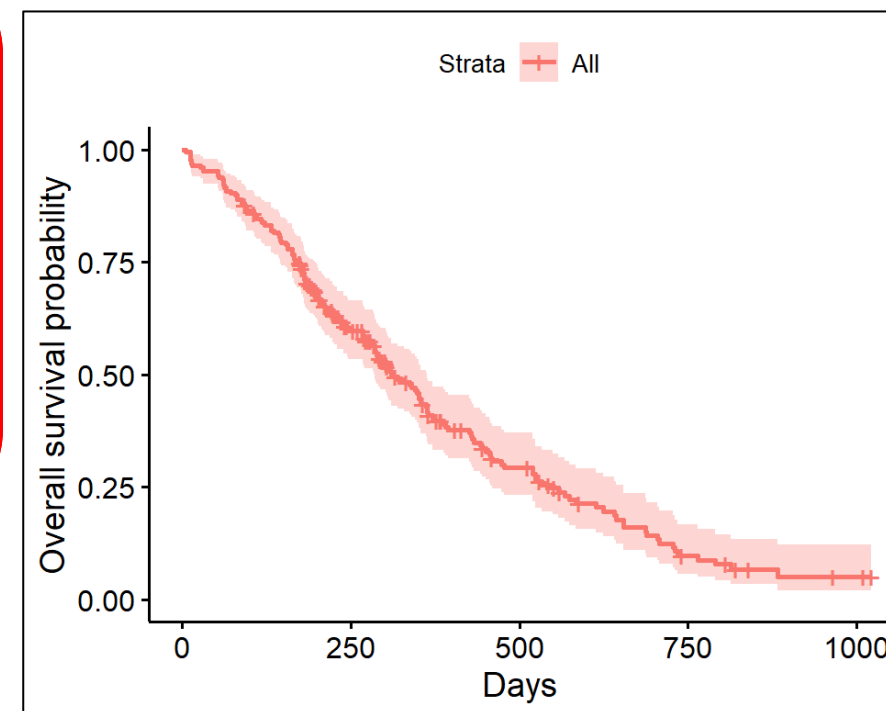
- Non-parametric estimation of the 'true' survival function [7]
- Deals with differing survival times where censoring is present [2]
- Probability of surviving in a given length of time while considering time as intervals [8]
- Different types of curves include overall survival, disease-free survival, progression-free survival, disease-specific survival [2]
- Tick marks are censored data
- Test difference between 2 curves? -> Log-Rank Test
- More than 2 curves or covariates? ->Cox regression [9]

Nelson-Aalen

- Non-parametric estimation of the 'true' cumulative hazard function [10]
- Used to estimate the cumulative number of expected **events** within a certain period of time [7]

Parametric Curves

- Will cover in part III
- Suffice to say they attempt to fit the survival curve smoothly





Curve Considerations

Kaplan-Meier

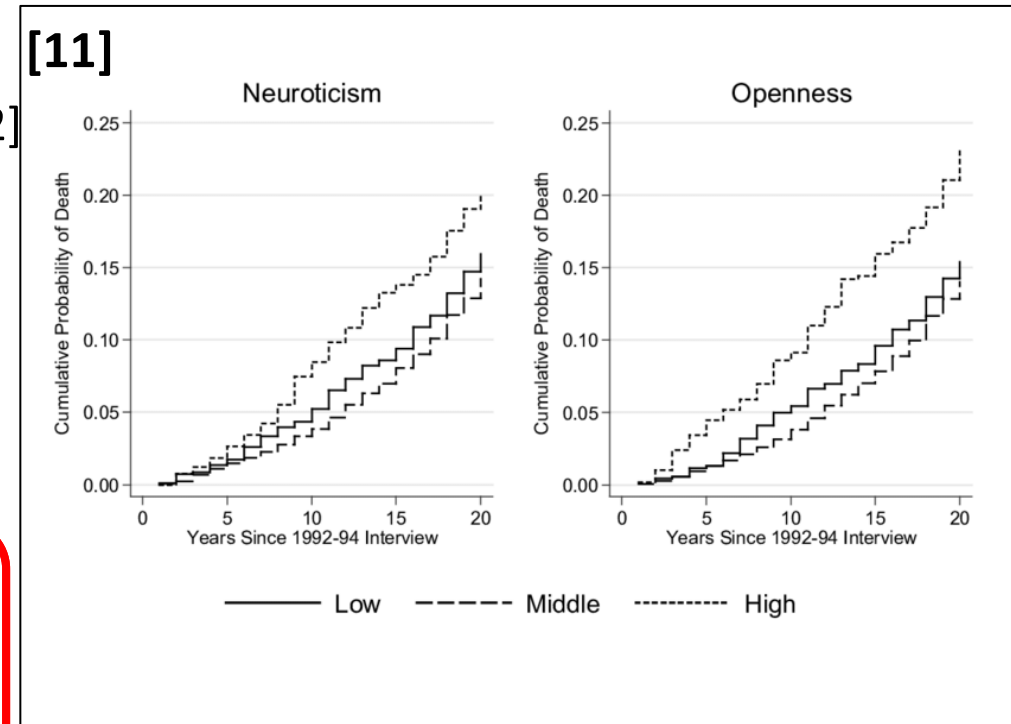
- Non-parametric estimation of the 'true' survival function [7]
- Deals with differing survival times where censoring is present [2]
- Probability of surviving in a given length of time while considering time as intervals [8]
- Different types of curves include overall survival, disease-free survival, progression-free survival, disease-specific survival [2]
- Tick marks are censored data
- Test difference between 2 curves? -> Log-Rank Test
- More than 2 curves or covariates? ->Cox regression [9]

Nelson-Aalen

- Non-parametric estimation of the 'true' cumulative hazard function [10]
- Used to estimate the cumulative number of expected **events** within a certain period of time [7]

Parametric Curves

- Will cover in part III
- Suffice to say they attempt to fit the survival curve smoothly





Curve Considerations

Kaplan-Meier

- Non-parametric estimation of the 'true' survival function [7]
- Deals with differing survival times where censoring is present [2]
- Probability of surviving in a given length of time while considering time as intervals [8]
- Different types of curves include overall survival, disease-free survival, progression-free survival, disease-specific survival [2]
- Tick marks are censored data
- Test difference between 2 curves? -> Log-Rank Test
- More than 2 curves or covariates? ->Cox regression [9]

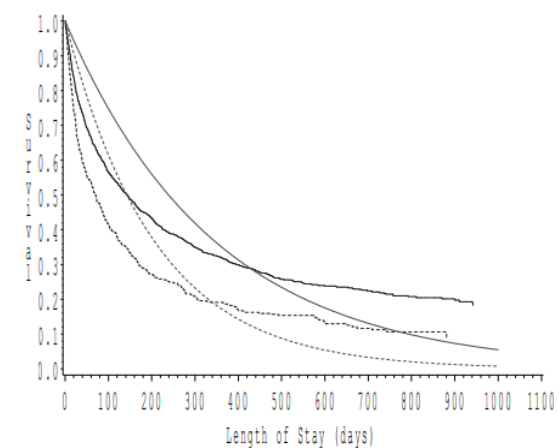
Nelson-Aalen

- Non-parametric estimation of the 'true' cumulative hazard function [10]
- Used to estimate the cumulative number of expected **events** within a certain period of time [7]

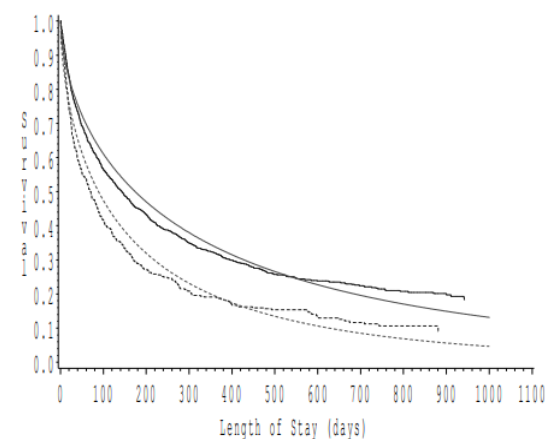
Parametric Curves

- Will cover in part III
- Suffice to say they attempt to fit the survival curve smoothly

Predicted Survival for Exponential model vs Kaplan–Meier



Predicted Survival for Weibull model vs Kaplan–Meier



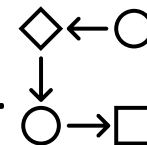
Assessment 1

qualtrics^{XM}



https://und.qualtrics.com/jfe/form/SV_6PznuRrv40l2lFk

Step-by-step Example 1.1



Survival Analysis in R: NCCTG Lung Cancer

`#Packages`

`library(ggplot2)`

`library(survival)`

`library(survminer)`

`library(broom)`

`library(knitr)`

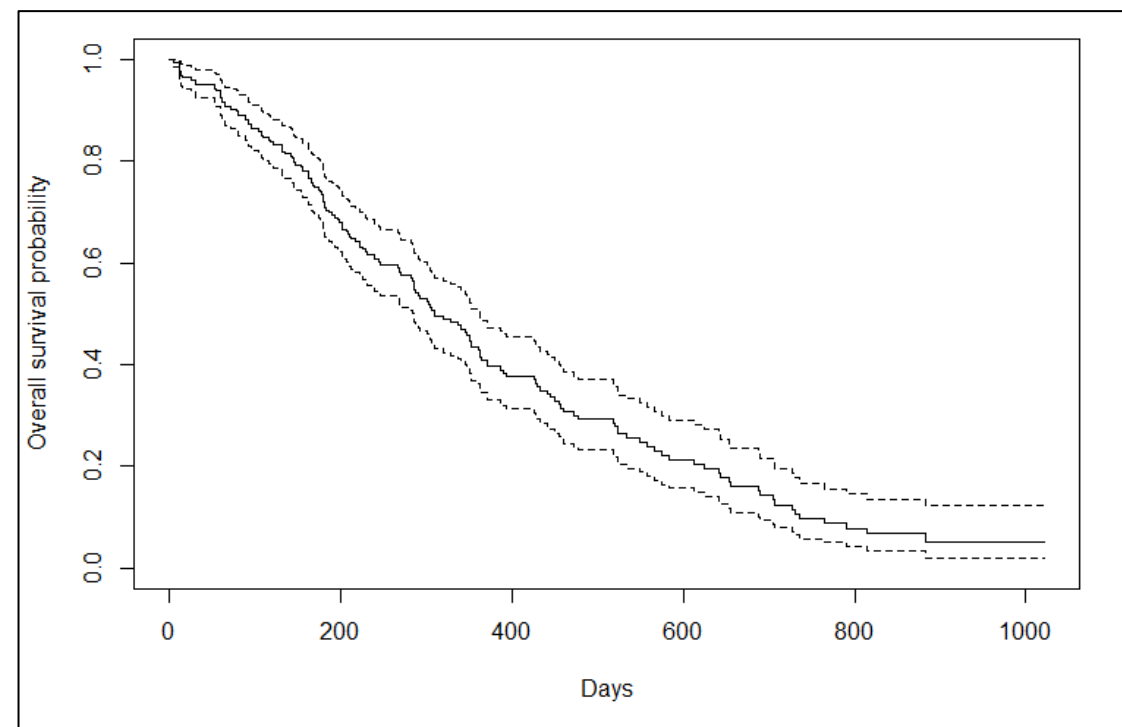
`head(lung)`

`#Basic Plot`

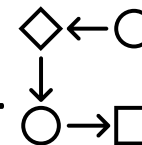
```
fit1 <- survfit(Surv(time, status)~1,  
               data=lung)
```

```
plot(fit1, xlab="Days", ylab="Overall survival probability")
```

	inst	time	status	age	sex	ph.ecog	ph.karno	pat.karno	meal.cal	wt.loss
1	3	306	2	74	1	1	90	100	1175	NA
2	3	455	2	68	1	0	90	90	1225	15
3	3	1010	1	56	1	0	90	90	NA	15
4	5	210	2	57	1	1	90	60	1150	11
5	1	883	2	60	1	0	100	90	NA	0
6	12	1022	1	74	1	1	50	80	513	0



Step-by-step Example 1.1



#Nicer Plot

```
fit2 <- ggsurvplot(
  fit = survfit(Surv(time, status)~ 1, data=lung),
  xlab="Days",
  ylab="Overall survival probability",
  palette = "orange"
)
```

fit2

#Median survival time

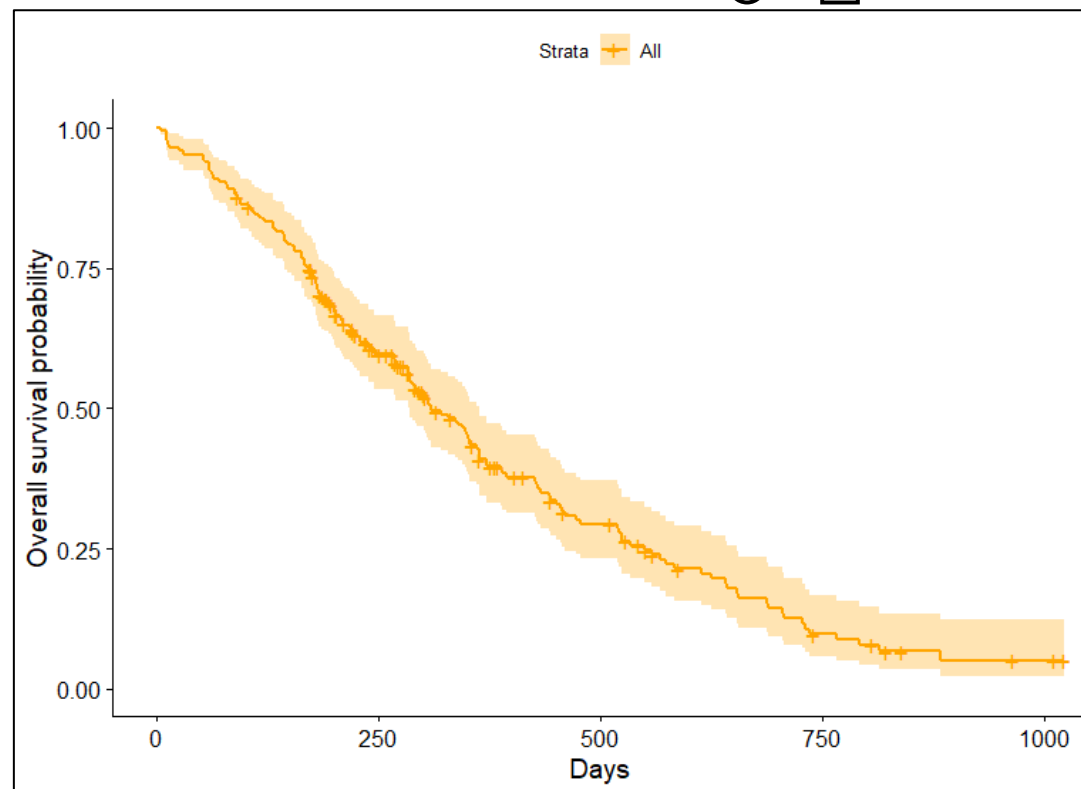
fit1

#Survival time based on years

```
summary(fit1, times=182.625) #six months
```

```
summary(fit1, times=365.25) #one year
```

```
summary(fit1, times=730.50) #two years
```

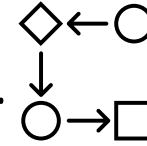


```
Call: survfit(formula = surv(time, status) ~ 1, data = lung)
```

n	events	median	0.95LCL	0.95UCL
228	165	310	285	363

Six months	One year	Two years
0.708	0.409	0.116

Step-by-step Example 1.2



Survival Analysis in R: Rat treatment

`head(rats)`

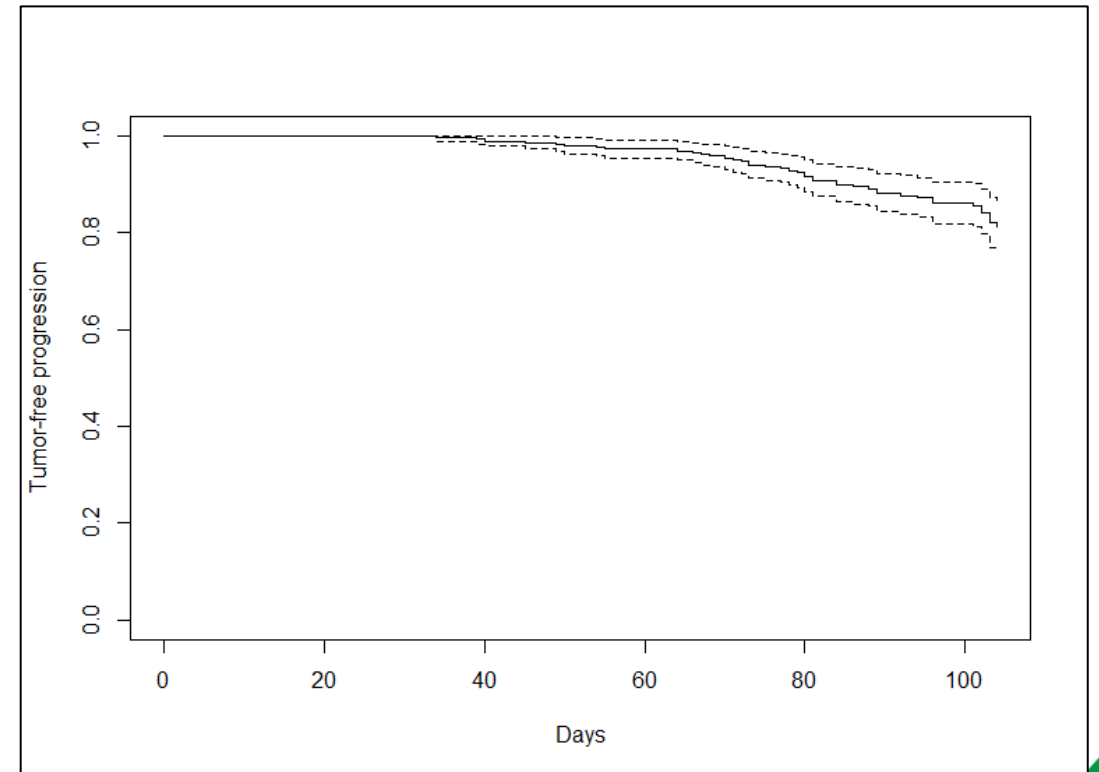
#Basic Plot

```
fit3 <- survfit(Surv(time, status),  
               data=rats)  
plot(fit3, xlab="Days", ylab="Tumor-free progression")
```

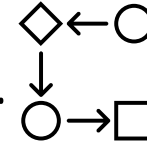
#Males and Female Plot

```
surv_obj1 <- Surv(time=rats$time, event=rats$status)  
fit4 <- survfit(surv_obj1 ~ sex, data=rats)  
ggsurvplot(fit4, data=rats, pval=TRUE)
```

	litter	rx	time	status	sex
1	1	1	101	0	f
2	1	0	49	1	f
3	1	0	104	0	f
4	2	1	91	0	m
5	2	0	104	0	m
6	2	0	102	0	m



Step-by-step Example 1.2



Survival Analysis in R: Rat treatment

```
head(rats)
```

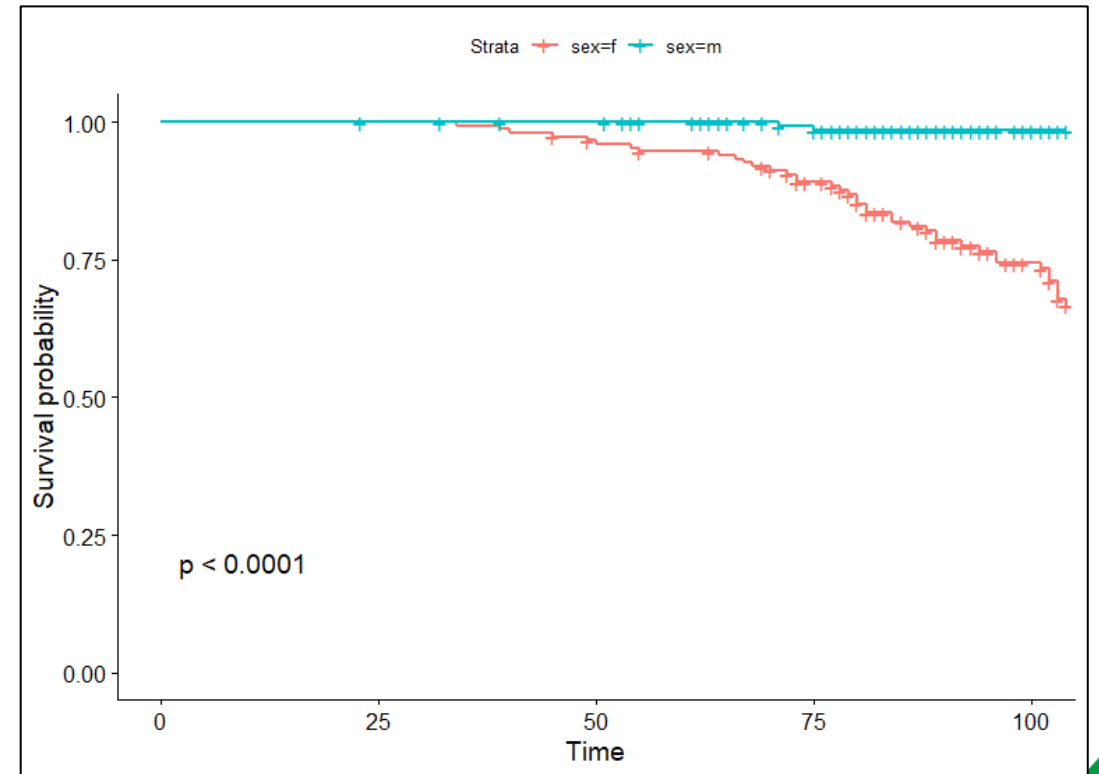
```
#Basic Plot
```

```
fit3 <- survfit(Surv(time, status),  
               data=rats)  
plot(fit3, xlab="Days", ylab="Tumor-free progression")
```

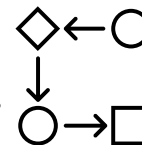
```
#Males and Female Plot
```

```
surv_obj1 <- Surv(time=rats$time, event=rats$status)  
fit4 <- survfit(surv_obj1 ~ sex, data=rats)  
ggsurvplot(fit4, data=rats, pval=TRUE)
```

	litter	rx	time	status	sex
1	1	1	101	0	f
2	1	0	49	1	f
3	1	0	104	0	f
4	2	1	91	0	m
5	2	0	104	0	m
6	2	0	102	0	m



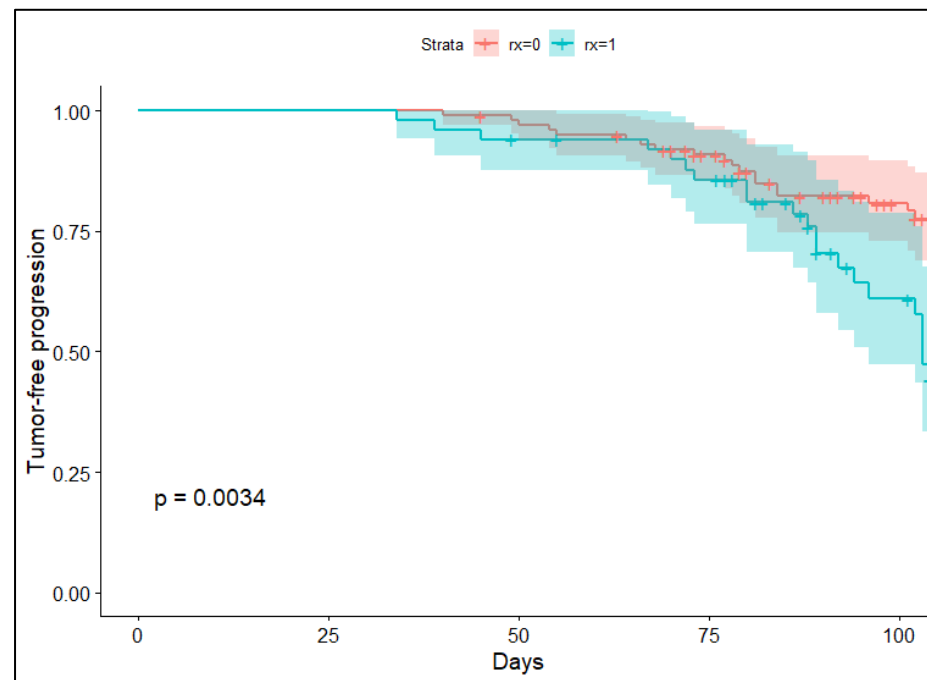
Step-by-step Example 1.2



```
#Treatment and Control Plot -> females only
f.rats <-rats[rats$sex=="f",]
surv_obj2 <- Surv(time=f.rats$time, event=f.rats$status)
fit5 <-survfit(surv_obj2 ~ rx, data=f.rats)
ggsurvplot(fit5, data=f.rats, pval=TRUE, conf.int=TRUE,
  xlab="Days",
  ylab="Tumor-free progression",)
```

```
#Rank test
survdif(surv_obj2~rx, data=f.rats)
```

```
#Hazards ratio
fit6 <- coxph(surv_obj2~ rx, data=f.rats)
broom::tidy(fit6, exp=TRUE) %>% kable()
```



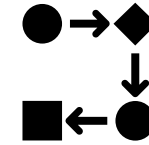
```
Call:
survdif(formula = surv_obj2 ~ rx, data = f.rats)

      N observed Expected (O-E)^2/E (O-E)^2/V
rx=0  100      19      27.5      2.65      8.61
rx=1   50      21      12.5      5.87      8.61

Chisq= 8.6 on 1 degrees of freedom, p= 0.003
```

term	estimate	std.error	statistic	p.value
rx	2.471277	0.3175104	2.849466	0.0043793

Step-by-step Example 2.1



Survival Analysis in SAS: NCCTG Lung Cancer Remix

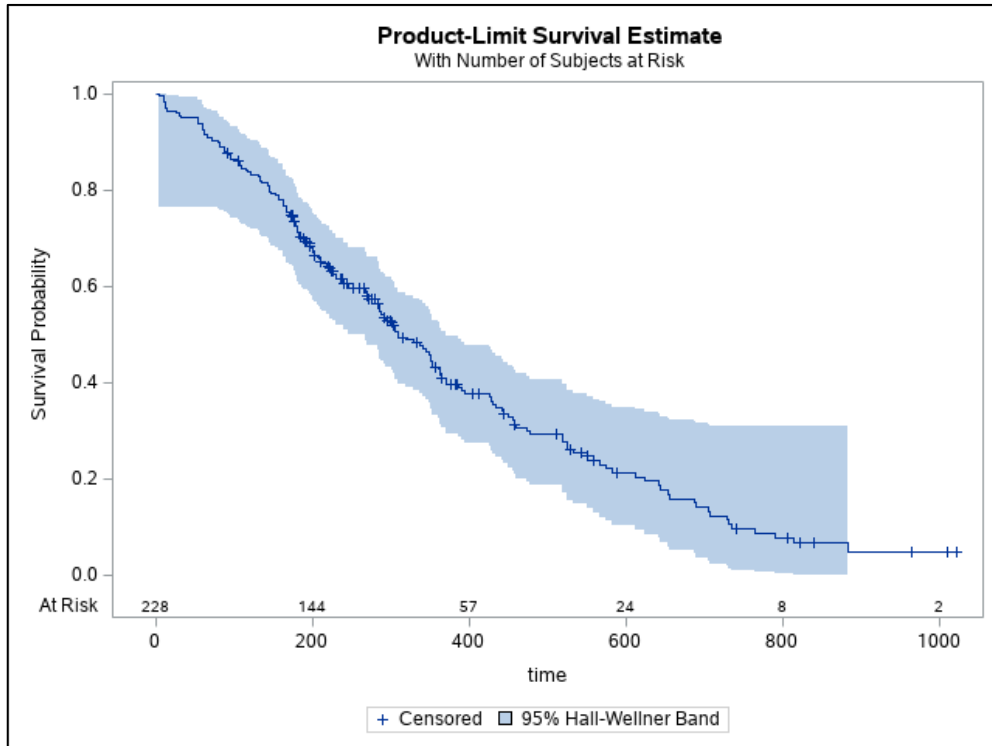
```

*Load Data;
DATA lung;
  input time status;
  cards;
306 2
455 2
1010 1
210 2
883 2
1022 1
310 2
361 2
218 2
166 2
...
;
PROC PRINT data=lung(obs=10);
  
```

Obs	time	status
1	306	2
2	455	2
3	1010	1
4	210	2
5	883	2
6	1022	1
7	310	2
8	361	2
9	218	2
10	166	2

Step-by-step Example 2.1

***Survival Plot and median survival time;
 PROC LIFETEST data=lung plots=survival(atrisk cb);
 time time*status(1);**

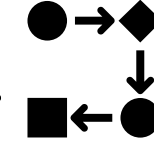


***Survival time based on years;
 PROC LIFETEST data=lung timelist=182.625 365.25 730.50;
 time time*status(1);**

Product-Limit Survival Estimates						
Timelist	time	Survival	Failure	Survival Standard Error	Number Failed	Number Left
182.625	182.00	0.7081	0.2919	0.0303	66	156
365.250	364.00	0.4092	0.5908	0.0358	121	65
730.500	728.00	0.1157	0.8843	0.0283	159	13

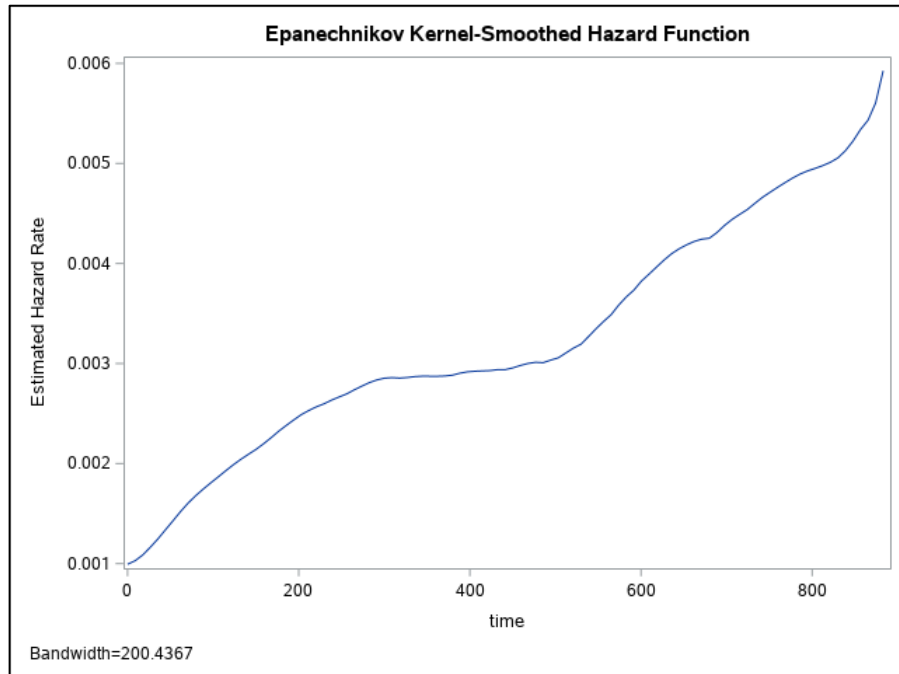
Quartile Estimates				
Percent	Point Estimate	95% Confidence Interval		
		Transform	[Lower	Upper)
75	550.00	LOGLOG	457.00	643.00
50	310.00	LOGLOG	284.00	361.00
25	170.00	LOGLOG	144.00	194.00

Step-by-step Example 2.1



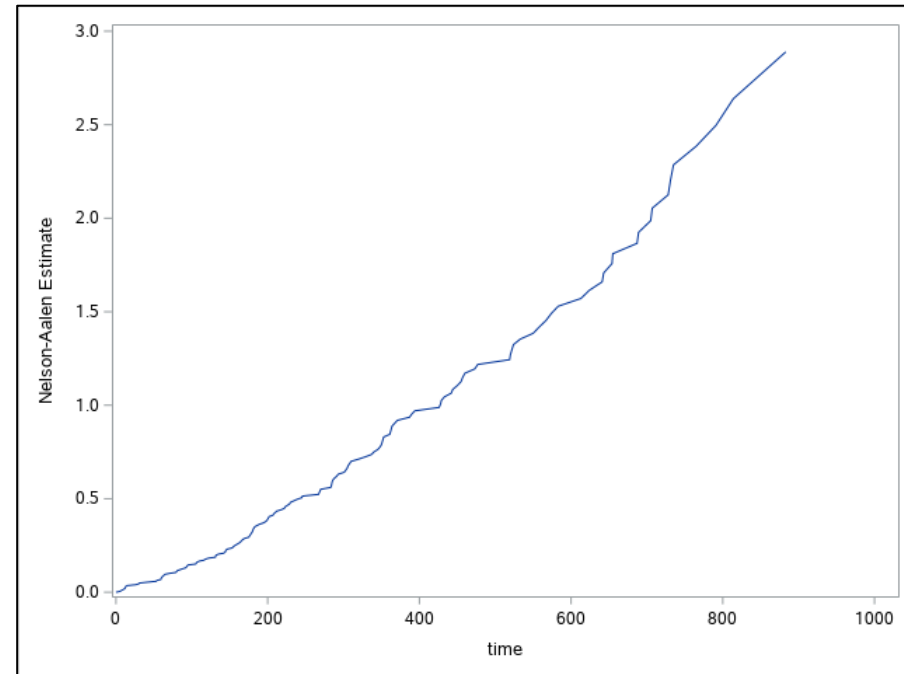
***Hazard Function;**

```
PROC LIFETEST data=lung plots=hazard;  
time time*status(1);
```



***Culmulative Hazard Function;**

```
PROC LIFETEST data=lung nelson method=pl;  
time time*status(1);  
ods output ProductLimitEstimates=lung_ple;  
PROC PRINT data=lung_ple(obs=25);  
PROC SGPLOT data=lung_ple;  
series x=time y=CumHaz;
```



Step-by-step Example 2.2

Survival Analysis in SAS: Bone Marrow Transplant

***Get data;**

DATA BMT;
set sashelp.BMT;

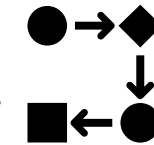
PROC PRINT data=BMT(obs=15);

PROC FREQ data=BMT;
tables Group*Status;

Obs	Group	T	Status
1	ALL	2081	0
2	ALL	1602	0
3	ALL	1496	0
4	ALL	1462	0
5	ALL	1433	0
6	ALL	1377	0
7	ALL	1330	0
8	ALL	996	0
9	ALL	226	0
10	ALL	1199	0
11	ALL	1111	0
12	ALL	530	0
13	ALL	1182	0
14	ALL	1167	0
15	ALL	418	1

Table of Group by Status			
Group(Disease Group)	Status(Event Indicator: 1=Event 0=Censored)		
	0	1	Total
ALL	14 10.22 36.84 25.93	24 17.52 63.16 28.92	38 27.74
AML-High Risk	11 8.03 24.44 20.37	34 24.82 75.56 40.96	45 32.85
AML-Low Risk	29 21.17 53.70 53.70	25 18.25 46.30 30.12	54 39.42
Total	54 39.42	83 60.58	137 100.00

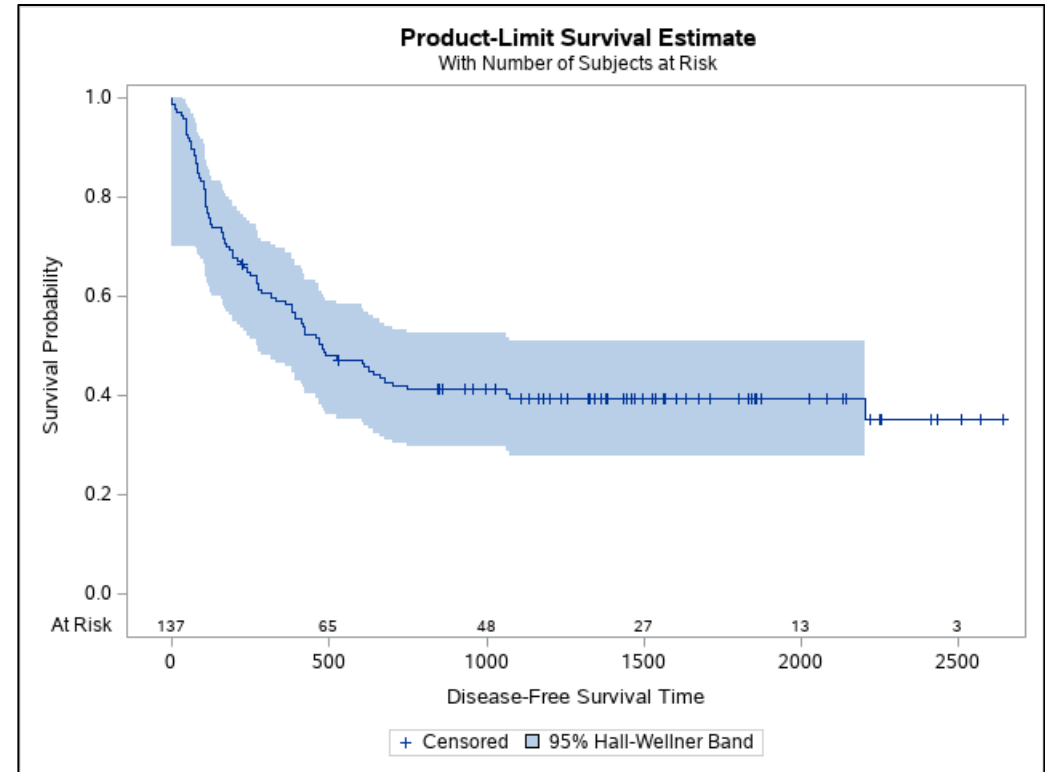
Step-by-step Example 2.2



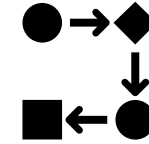
***Overall Survival Plot;**
PROC LIFETEST data=BMT plots=survival(atrisk cb);
time T*Status(0);

***Survival Function (low and high risk) and rank test;**
PROC LIFETEST data=BMT plots=survival(atrisk cb);
where Group not in ('ALL');
time T*Status(0);
strata Group / test=logrank;

Test of Equality over Strata			
Test	Chi-Square	DF	Pr > Chi-Square
Log-Rank	13.4456	1	0.0002



Step-by-step Example 2.2



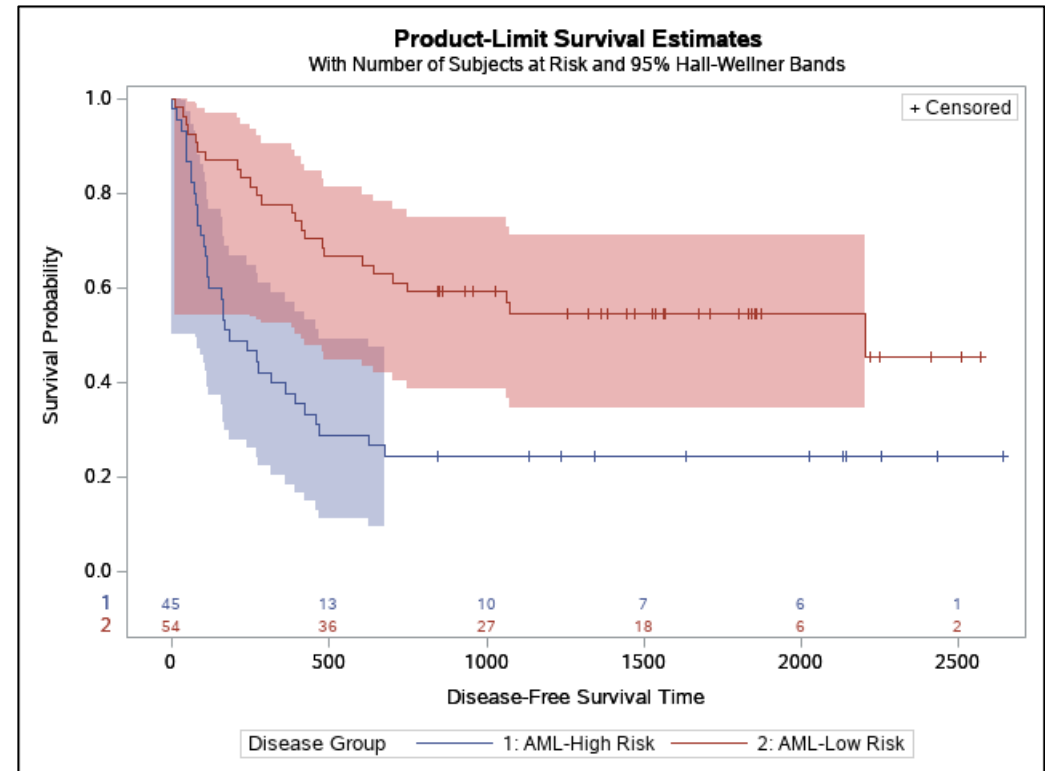
***Overall Survival Plot;**

**PROC LIFETEST data=BMT plots=survival(atrisk cb);
 time T*Status(0);**

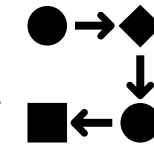
***Survival Function (low and high risk) and rank test;**

**PROC LIFETEST data=BMT plots=survival(atrisk cb);
 where Group not in ('ALL');
 time T*Status(0);
 strata Group / test=logrank;**

Test of Equality over Strata			
Test	Chi-Square	DF	Pr > Chi-Square
Log-Rank	13.4456	1	0.0002



Step-by-step Example 2.2

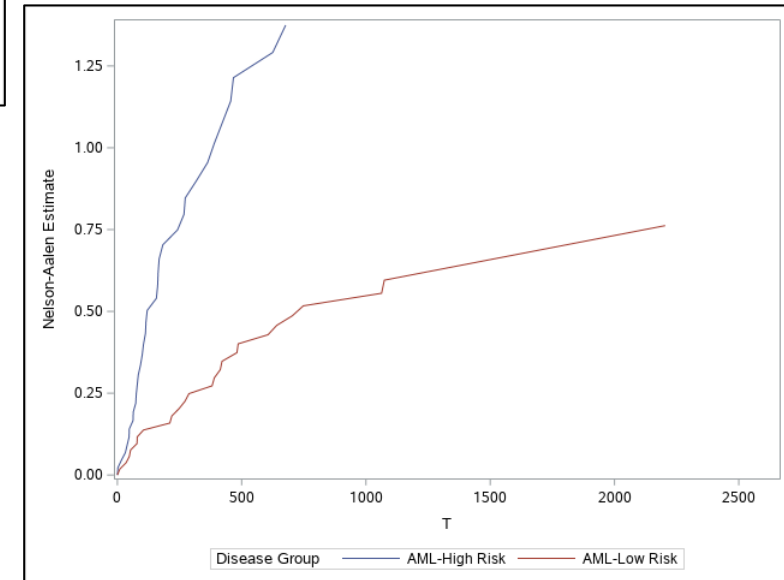
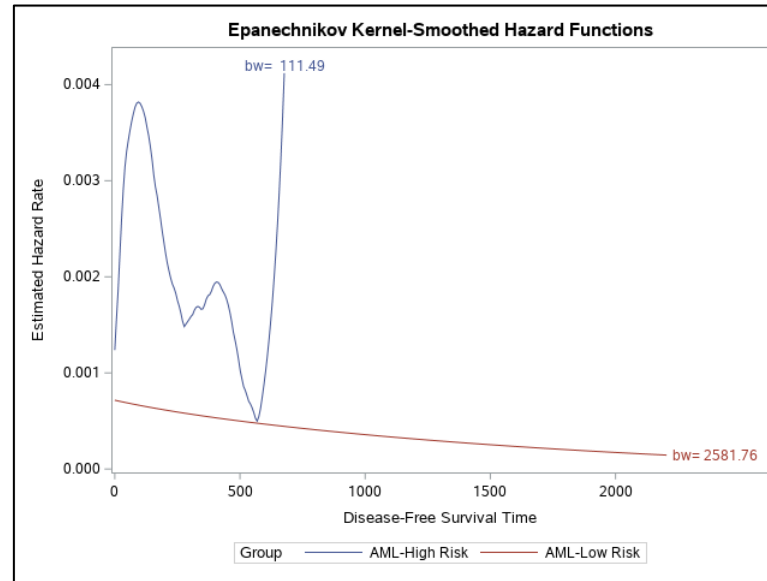


```

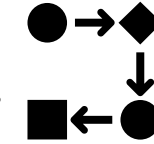
*Hazard Function;
PROC LIFETEST data=BMT plots=hazard;
  where Group not in ('ALL');
  time T*Status(0);
  strata Group;

*Culmulative Hazard Function;
PROC LIFETEST data=BMT nelson method=pl;
  where Group not in ('ALL');
  time T*Status(0);
  strata Group;
  ods output ProductLimitEstimates=BMT_ple;

PROC PRINT data=BMT_ple;
PROC SGPLOT data=BMT_ple;
  series x=T y=CumHaz/ group=Group;
  
```



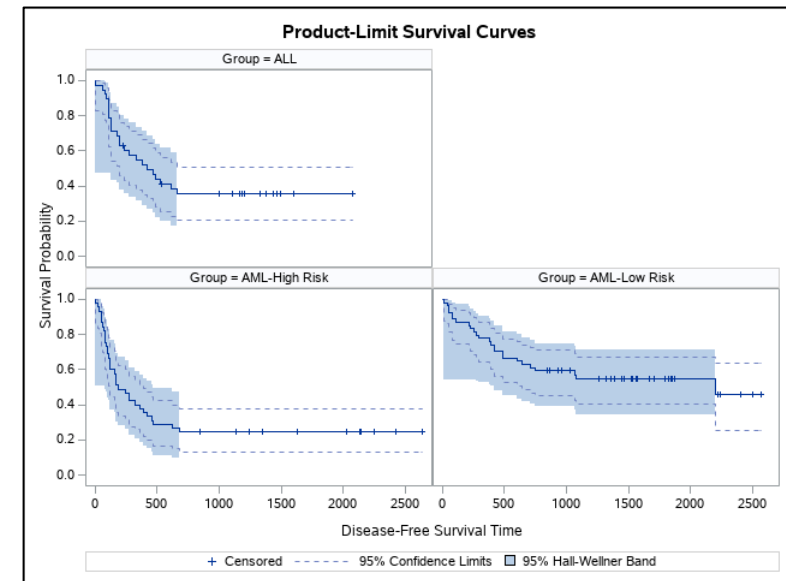
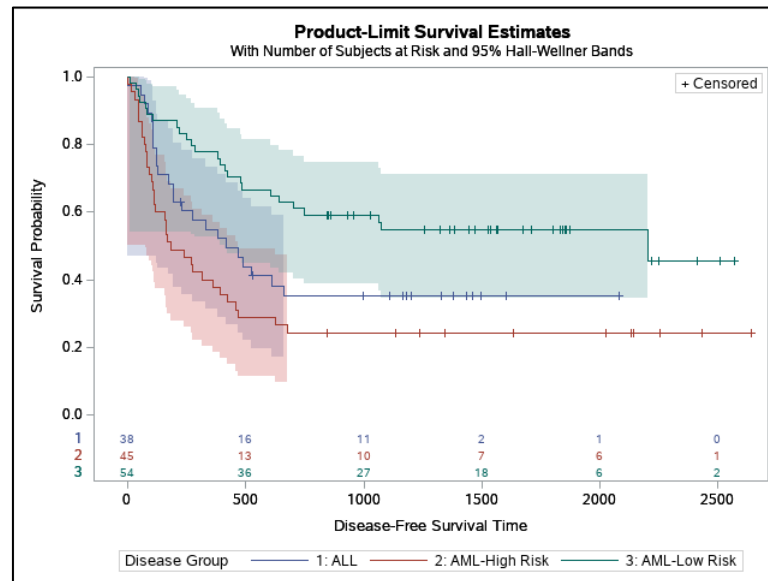
Step-by-step Example 2.2



***Adding a third group;**
PROC LIFETEST data=BMT plots=survival(atrisk cb);
time T*Status(0);
strata Group / test=logrank adjust=sidak;

***Plot curves separately;**
PROC LIFETEST data=BMT plots=survival(cl cb=hw
strata=panel);
time T*Status(0);
strata Group;

Adjustment for Multiple Comparisons for the Logrank Test				
Strata Comparison		Chi-Square	p-Values	
Group	Group		Raw	Sidak
ALL	AML-High Risk	2.6610	0.1028	0.2779
ALL	AML-Low Risk	5.1400	0.0234	0.0685
AML-High Risk	AML-Low Risk	13.8011	0.0002	0.0006



Assessment 2



qualtrics^{XM}[®]



UNIVERSITY OF NORTH DAKOTA



https://und.qualtrics.com/jfe/form/SV_6KXKzqGqDIh2BP8

Caveats and Concerns



- Too much or too few censoring
 - No censoring, standard regression could possibly be used [1]
 - Too many censors questions the validity of the study design [2]
- Truncation is also a thing [12]
- Checking curves and survival function to match
- Make sure assumptions are valid [13]
 - Random sampling*
 - Independent survival times
 - Measuring survival doesn't impact survival
 - Time to event is known with accuracy
 - Censoring mechanism is random

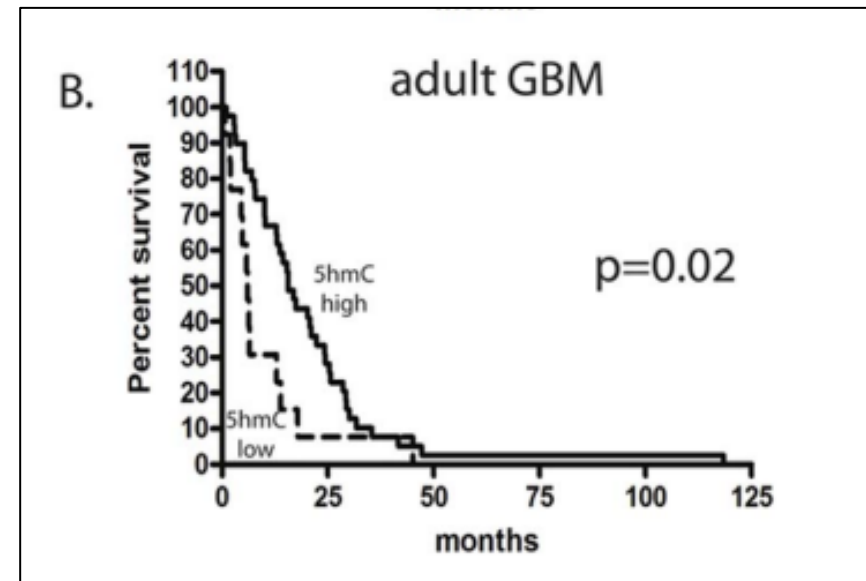
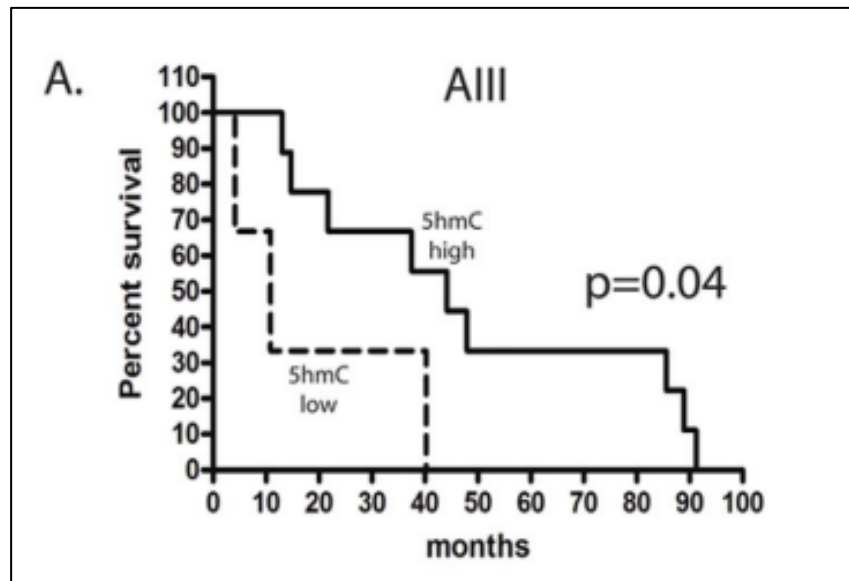
Real World Examples

Orr BA, Haffner MC, Nelson WG, Yegnasubramanian S, Eberhart CG (2012) Decreased 5-Hydroxymethylcytosine Is Associated with Neural Progenitor Phenotype in Normal Brain and Shorter Survival in Malignant Glioma. PLOS ONE 7(7): e41036.
<https://doi.org/10.1371/journal.pone.0041036>

[14]

Epigenetic modification of DNA by cytosine methylation to produce 5-methylcytosine (5mC) has become well-recognized as an important epigenetic process in human health and disease

...we identified a significant relationship between low levels of 5hmC and reduced survival in malignant glioma



Real World Examples



Mjørud M, Selbæk G, Bjertness E, Edwin TH, Engedal K, et al. (2020) Time from dementia diagnosis to nursing-home admission and death among persons with dementia: A multistate survival analysis. PLOS ONE 15(12): e0243513.

<https://doi.org/10.1371/journal.pone.0243513>

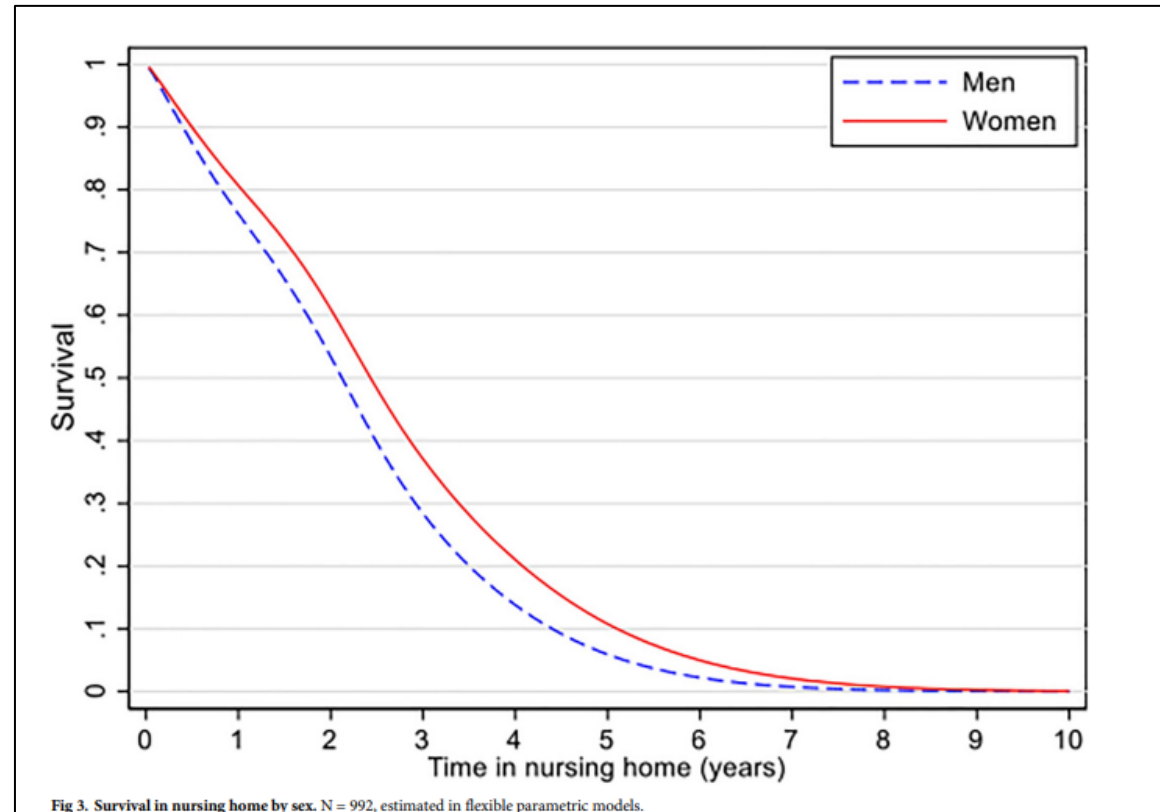
[15]

To estimate transition times from dementia diagnosis to nursing-home (NH) admission or death and to examine whether sex, education, marital status, level of cognitive impairment and dementia aetiology are associated with transition times.

Markov multistate survival analysis and flexible parametric models

The probability of NH admission was greater for women than men due to women's lower mortality rate.

Age, dementia aetiology and severity of cognitive impairment at time of diagnosis did not influence the probability of NH admission.



Summary and Conclusion

- Survival analysis provides estimates of survival for a group and can also determine if two groups are significantly different
- Censoring needs to be kept in mind
- The selection of functions, curves, and tests depends on the analysis
- Tune in next time for deeper details into survival analysis in Survival Analysis Module III: Deep Dive

References

- [1] <http://www.stat.columbia.edu/~madigan/W2025/notes/survival.pdf>
- [2] <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3932959/pdf/nihms549224.pdf>
- [3] [https://www.emilyzabor.com/tutorials/survival analysis in r tutorial.html](https://www.emilyzabor.com/tutorials/survival%20analysis%20in%20r%20tutorial.html)
- [4] <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3275994/pdf/PCR-3-40.pdf>
- [5] https://web.njit.edu/~wguo/Math%20659_2011/Math659_Chapter3.pdf
- [6] <https://www.math.ucsd.edu/~rxu/math284/slect3.pdf>
- [7] <https://towardsdatascience.com/introduction-to-survival-analysis-the-nelson-aalen-estimator-9780c63d549d>
- [8] <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3059453/pdf/IJAR-1-274.pdf>
- [9] <https://www.r-bloggers.com/2019/06/parametric-survival-modeling/>
- [10] <https://www.medicine.mcgill.ca/epidemiology/hanley/c609/material/NelsonAalenEstimator.pdf>
- [11] https://www.researchgate.net/publication/270593359_Who_Lives_Longer_and_Healthier_The_Role_of_Personality_Facial_Attractiveness_and_Intelligence
- [12] https://web.njit.edu/~wguo/Math%20659_2011/Math659_Chapter3.pdf
- [13] <https://www.jstor.org/stable/pdf/3802856.pdf?refreqid=excelsior%3Acc9cce2a2e3afe46b7af9801147e3e81>
- [14] <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0041036>
- [15] <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0243513>

Acknowledgements

- The DaCCoTA is supported by the National Institute of General Medical Sciences of the National Institutes of Health under Award Number U54GM128729.
- For the labs that use the Biostatistics, Epidemiology, and Research Design Core in any way, including this Module, please acknowledge us for publications. ***"Research reported in this publication was supported by DaCCoTA (the National Institute of General Medical Sciences of the National Institutes of Health under Award Number U54GM128729)".***

DaCCoTA
DAKOTA CANCER COLLABORATIVE
ON TRANSLATIONAL ACTIVITY