

NTNU
Det skapende universitet

Survival analysis (Levetidsanalyse)
Rosner 14.8-14.11

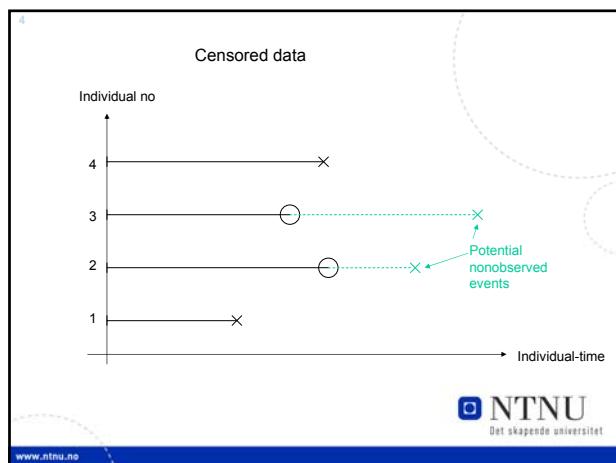
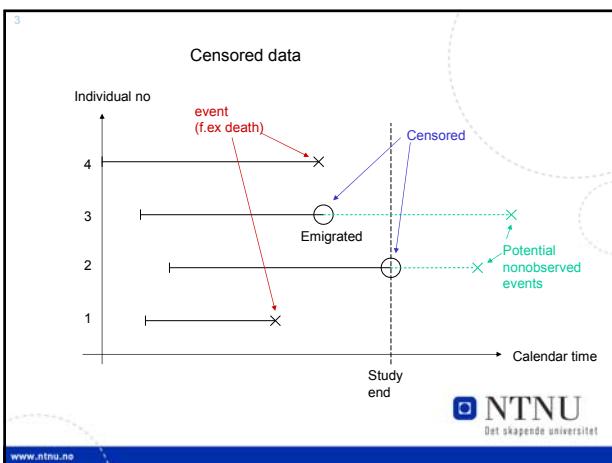
By Stian Lydersen
Lecture 21 april 2010

www.ntnu.no

Lifetime (duration of a state)

- Examples:
 - Time to death (from birth, from diagnosis, from treatment start)
 - Time to a diagnosis
 - Time from start of treatment until declared disease-free
- Lifetimes are
 - Always greater than 0
 - Often censored
- $S(t) = \text{Probability of surviving } t$

NTNU
Det skapende universitet



Lifetime analysis – most used methods

- Actuarial calculations (approximate calculations within time intervals)
- Kaplan-Meier plot (empirical survival probability)
- Log-rank test: Non-parametric test for difference between groups
- Regression analysis:
 - Semi-parametric (Cox regression = Proportional hazard regression)
 - Parametric (Assume a certain shape of the distribution, f.ex. Weibull-distribution)

NTNU
Det skapende universitet

ISI-search 7 april 2010

COX DR
REGRESSION MODELS AND LIFE-TABLES
JOURNAL OF THE ROYAL STATISTICAL SOCIETY SERIES B-STATISTICAL METHODOLOGY 34 (2):
187& 1972
Times Cited: 24441

Title: NONPARAMETRIC-ESTIMATION FROM INCOMPLETE OBSERVATIONS
Author(s): KAPLAN EL, MEIER P
Source: JOURNAL OF THE AMERICAN STATISTICAL ASSOCIATION 53 (282): 457-481 1958
Times Cited: 33937

NTNU
Det skapende universitet

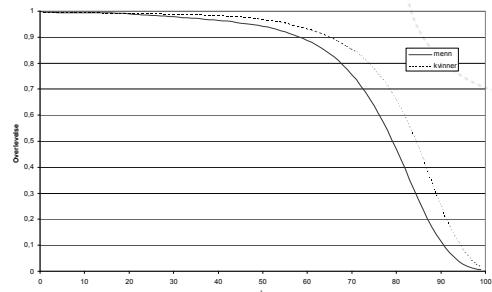
Methods for lifetime analysis assume non-informative censoring

- For each individual there exist a potential lifetime T and a censoring time C .
- We observe only $\min(T, C)$ and D , where
 - $D=0$ if censored
 - $D=1$ if event (f.ex death)
- The methods for lifetime analysis assume T, D to be independent
- Not possible to check independence from data



www.ntnu.no

Norge 2000. Basert på dødelighetstall fra www.ssb.no



www.ntnu.no

Kaplan-Meier (produkt-grense) estimatoren for overlevelsessannsynligheten

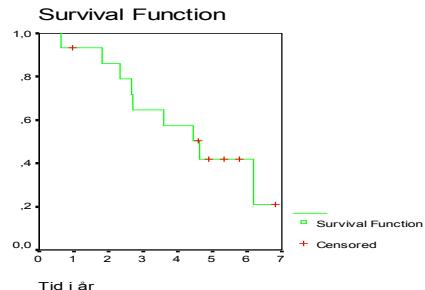
Data om overlevelse for 15 pasienter med melanom (fra Aalen, 1998)

Antall i live og under observasjon	Status (0=sensurert, 1=død)	Observasjons- tid i år	Overlevels- sannsynlighet (Kaplan-Meier)
15	1	0,64	0,933
14	0	0,97	
13	1	1,81	0,862
12	1	2,35	0,79
11	1	2,65	0,718
10	1	2,69	0,646
9	1	3,59	0,574
8	1	4,44	0,503
7	0	4,6	
6	1	4,63	0,419
5	0	4,9	
4	0	5,32	
3	0	5,76	
2	1	6,18	0,209
1	0	6,83	



www.ntnu.no

Kaplan-Meier plott (empirisk overlevessannsynlighet)



Det skapende universitet

Comparing two groups:

H_0 : No group effect
 $S_1(t) = S_2(t)$ for all t

against

H_1 : lifetimes tend to be longer (alternatively shorter) in Group 1
 $S_1(t) \geq S_2(t)$ for all t with $>$ for some t (alternatively \leq og $<$)

Non-parametric test for complete or censored data:
Logrank test.

Approximately same result:

Test if $\beta=0$ in Cox model with x as group indicator.
(Exactly same result if no ties AND use of LR p-value in Cox model.)



www.ntnu.no

Kleinbaum & Klein (2005) s 17 Remisjonstider for 42 leukemi-pasienter (Freireich & al, Blood, 1963)

uker	hend.	grp	uker	hend.	grp	uker	hend.	grp
1	1	2	8	1	2	16	1	1
1	1	2	8	1	2	17	0	1
2	1	2	8	1	2	17	1	2
2	1	2	8	1	2	19	0	1
3	1	2	9	0	1	20	0	1
4	1	2	10	1	1	22	1	1
4	1	2	10	0	1	22	1	2
5	1	2	11	0	1	23	1	1
5	1	2	11	1	2	23	1	2
6	1	1	11	1	2	25	0	1
6	1	1	12	1	2	32	0	1
6	1	1	12	1	2	32	0	1
6	0	1	13	1	1	34	0	1
7	1	1	15	1	2	35	0	1

Hendelse
0 = sensur
1 = remisjon
Grp
1 = Behandling
2 = placebo



www.ntnu.no

13

Tid (uke)	hendelser				rest
	gr1 h1	gr2 h2	r1	r2	
1	0	2	21	21	
2	0	2	21	19	
3	0	1	21	17	
4	0	2	21	16	
5	0	2	21	14	
6	3	0	21	12	
7	1	0	17*	12	
8	0	4	16	12	
10	1	0	15	8	
11	0	2	13	8	
12	0	2	12	6	
13	1	0	12	4	
15	0	1	11	4	* Sensurerte tider kommer ikke frem, kursiv antyder uker med sensur.
16	1	0	11	3	
17	0	1	10	3	
22	1	1	7	2	
23	1	1	6	1	
Sum		9	21		



Det skapende universitet

www.ntnu.no

* Sensurerte tider
kommer ikke frem,
kursiv antyder
uker med sensur.

14

Log-rang testen

Likner på en kjikkvadrat "goodness of fit" test

Del opp i korte tidsintervall

(bare ett tidspunkt med hendelse i hvert intervall)

Forventet antall hendelser i intervall nr j:

Expected = andel i risikomengden \times antall hendelser totalt

$$e_{1j} = \left(\frac{r_{1j}}{r_{1j} + r_{2j}} \right) \times (h_{1j} + h_{2j}),$$

$$e_{2j} = \left(\frac{r_{2j}}{r_{1j} + r_{2j}} \right) \times (h_{1j} + h_{2j})$$



Det skapende universitet

15

Tid (uke)	hendelser				expected	Gr2**	Obs-Exp
	gr1 h1	gr2 h2	gr1 r1	gr2 r2			
1	0	2	21	21	21/42 *2	21/42 *2	1.00
2	0	2	21	19	21/40 *2	19/40 *2	1.05
3	0	1	21	17	21/38 *1	17/38 *1	0.55
4	0	2	21	16	21/37 *2	16/37 *2	1.14
5	0	2	21	14	21/35 *2	14/35 *2	1.20
6	3	0	21	12	21/33 *3	12/33 *3	-1.09
7	1	0	17*	12	17/29 *2	12/29 *2	-0.41
8	0	4	16	12	16/28 *4	12/28 *4	2.29
10	1	0	15	8	15/23 *1	8/23 *1	-0.35
11	0	2	13	8	13/21 *2	8/21 *2	1.24
12	0	2	12	6	12/18 *2	6/18 *2	1.33
13	1	0	12	4	12/16 *1	4/16 *1	-0.25
15	0	1	11	4	11/15 *1	4/15 *1	0.73
16	1	0	11	3	11/14 *1	3/14 *1	-0.21
17	0	1	10	3	10/13 *1	3/13 *1	0.77
22	1	1	7	2	7/9 *2	2/9 *2	0.56
23	1	1	6	1	6/7 *2	1/7 *2	0.71
Sum		9	21		19.26	10.74	10.3



Det skapende universitet

www.ntnu.no

sensurerte kommer ikke frem, kursiv antyder uker med sensur

*Observerd-expected for gruppe 1 = -(Observerd-expected) for gruppe 2

16

Log-rank test

- log-rank observator: $(\sum_{failure times} obs-forv)^2 / var (\sum obs-forv)$

$$var (\sum obs-forv) = \sum \frac{r_1 * r_2 (h_1 + h_2) * (r_1 + r_2 - h_1 - h_2)}{(r_1 + r_2)^2 * (r_1 + r_2 - 1)} = 6.3$$

- Når to grupper blir sammenliknet brukes kun resultater for en gruppe; (obs-forv) for gruppe1 = -(obs-forv) for gruppe 2

$$\text{For gruppe 2: } (\sum_{failure times} obs-forv)^2 = (10.3)^2 = 106.09$$

$$\text{log-rank: } \chi^2 = 16.8, df = 1 \text{ gir p-verdi } < 0.0001$$

- Mer komplisert hvis det er mer enn 2 grupper som sammenliknes



Det skapende universitet

17

Survival probability, survival function :
("overlevelsessansynlighet")

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

Hazard rate, force of mortality ("dødelighet"): $z(t)$

The probability that an individual of age t
dies during the next time interval of duration Δt :

$$P(T \leq t + \Delta t | T > t) \approx z(t) \Delta t$$

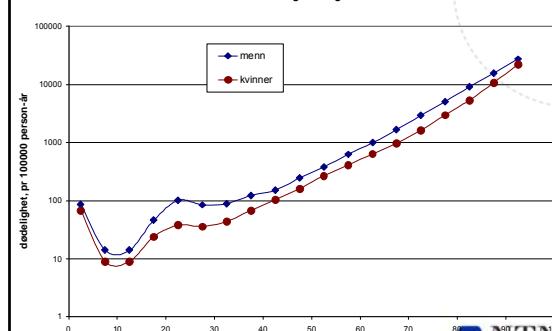


Det skapende universitet

www.ntnu.no

18

Dødelighet Norge 2005



Det skapende universitet

19

Formal definition of the hazard rate:

$$z(t) = \lim_{\Delta t \rightarrow 0} \frac{P(T \leq t + \Delta t | T > t)}{\Delta t}$$

Which is equal to

$$\begin{aligned} z(t) &= \lim_{\Delta t \rightarrow 0} \frac{P((T \leq t + \Delta t) \cap (T > t))}{\Delta t} = \lim_{\Delta t \rightarrow 0} \frac{P(t < T \leq t + \Delta t)}{\Delta t} \cdot \frac{1}{P(T > t)} \\ &= \lim_{\Delta t \rightarrow 0} \frac{F(t + \Delta t) - F(t)}{\Delta t} \cdot \frac{1}{P(T > t)} = \frac{f(t)}{S(t)} \end{aligned}$$

And also:

$$z(t) = -\frac{d}{dt}(\ln(S(t)))$$



www.ntnu.no

Sammenheng mellom $S(t)$ og $h(t)$:

$$S(t) = e^{-\int_0^t h(u) du} = e^{-H(t)}$$

hvor

$$H(t) = \int_0^t h(u) du$$

kalles kumulativ hazard.

Det er lettere å plotte $H(t)$ enn $h(t)$.

Hazardraten $h(t) = H'(t)$ er stigningsstallet til $H(t)$.



21

Cox Proporsjonal Hazard (PH) modell

$$h(t; \underline{x}) = h_0(t) \exp(\beta \underline{x})$$

$$\begin{aligned} h(t; x_1, \dots, x_p) &= h_0(t) \exp(\beta_1 x_1 + \dots + \beta_p x_p) \\ &= h_0(t) \exp(\beta_1 x_1) \dots \exp(\beta_p x_p) \end{aligned}$$

ekvivalent:

$$S(t; \underline{x}) = S_0(t)^{\exp(\beta \underline{x})}$$

hvor $h_0(t)$, $S_0(t)$ kalles "base line hazard rate", "base line overlevelsessannsynlighet"

$h_0(t)$ er ukjent og helt uspesifisert
 β_1, \dots, β_p er ukjente parametere



www.ntnu.no

22

An interpretation of β

- If two individuals (or groups) have respectively $x_1=0$ and $x_1=1$ and other covariates are equal: The ratio of the hazard rates equals $\exp(\beta_1)$, independent of "age" t .
- The quantity $\exp(\beta_1)$ is called the hazard ratio (HR).
- Relevant only if the covariate is not included in an interaction term!



23

Cox, D. R. (1972):

- A method for estimation of β_1, \dots, β_p
- and non-parametric estimation of $S_0(t)$ (Similar to the Kaplan-Meier estimate)
- for complete and censored data sets
- Time dependent $\beta_i(t)$ is possible.
- (Alternatively: Parametric model, for example lognormal or Weibull distribution for $S(t)$.)



www.ntnu.no

Partial likelihood (Cox):

$$l_p(\underline{\beta}) = \prod_{i=1}^m \frac{h_0(t_i) \exp(x_{(i)} \underline{\beta})}{\sum_{j \in R(t_{(i)})} h_0(t_j) \exp(x_{(j)} \underline{\beta})} = \prod_{i=1}^m \frac{\exp(x_{(i)} \underline{\beta})}{\sum_{j \in R(t_{(i)})} \exp(x_{(j)} \underline{\beta})}$$

where $t_{(1)} < t_{(2)} < \dots < t_{(m)}$ are the m distinct (uncensored) lifetimes.
 $R(t_{(j)})$ is the risk set at $t_{(j)}$,
that is, individuals alive and uncensored at $t_{(j)}$
 $x_{(j)}$ are the covariates for the individual with lifetime $t_{(j)}$.



www.ntnu.no

25

The partial likelihood does not depend on $h_0(t)$!

Estimates for $\underline{\beta}$ as well as variances and covariances for the are obtained by treating $l_p(\underline{\beta})$ as an ordinary likelihood.



www.ntnu.no

With tied observations, a more general definition of $l_p(\underline{\beta})$ is used.

Alternatives:

- Exakt. Very time consuming.
- Breslow (1974). Quick.
- Efron (1977) Quick. Almost identical results as exact.

SPSS: Only Breslow is available

Stata: All 3 available. Breslow is default.
Use option "Efron" when tied observations.



Det skapende universitet

27

Checking the PH assumption: Log - log plot

We have

$$S(t, \underline{x}) = S_0(t)^{\exp(\underline{x}' \underline{\beta})}$$

so

$$\log(-\log(S(t, \underline{x}))) = \underline{x}' \underline{\beta} + \log(-\log(S_0(t)))$$

Plot for different \underline{x} ought to be parallel under the PH assumption.

Difficult to judge parallelity when few observations.
More advanced methods exist.



www.ntnu.no

28

Software for Cox PH regression

- SPSS – has some
- MINITAB – has some
- STATA – has much
- S-plus or R- has much
- SAS - has much



Det skapende universitet

29

BMJ 2003;326:822

Parametric survival models may be more accurate than Kaplan-Meier estimates

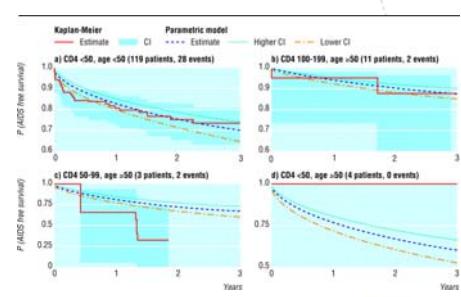
EDITOR---Lundin et al use Kaplan-Meier estimates of survival probabilities in their system for survival estimation in breast cancer ... They claim that researchers can obtain survival estimates based on actual data, rather than inferential estimates generated by a regression formula. However, any regression formula is based on actual data. More importantly, survival estimates from a regression model may be substantially more precise than Kaplan-Meier estimates when there are few patients in particular strata.



www.ntnu.no

30

Kaplan-Meier versus Weibull distribution



May, M. et al. BMJ 2003;326:822



BMJ

References

- Hosmer, D. W., Lemeshow, S. D., May, S.: "Applied Survival Analysis. Regression Modeling of Time to Event Data." 2nd Ed. Wiley, 2008.
 - Comprehensive and well written book.
- Kleinbaum, D. G., and Klein, M.: "Survival Analysis: A Self-Learning Text". 2nd ed. Springer, 2005.
 - Introductory book, easy to read

