

Analiza cenzuriranih podataka i krivulje preživljjenja

Mladen Petrovečki

Katedra za medicinsku informatiku, Medicinski fakultet Sveučilišta u Rijeci
i Klinička bolnica Dubrava, Zagreb

Analiza preživljjenja

1. temeljni pojmovi
2. izračun vjerojatnosti preživljjenja
 - a) tablice preživljjenja
 - b) Kaplan-Meierov postupak
3. rizik umiranja
4. programska potpora
5. usporedba podataka o preživljjenju
6. statističko zaključivanje
7. regresijska analiza cenzuriranih podataka

Obrada podataka o preživljjenju bolesnika

- analiza preživljjenja
- *survival analysis*
- ponekad
 - analiza tablica preživljjenja
 - analiza osiguravateljskih (aktuarskih) podataka
 - *actuarial analysis*

Analiza preživljjenja

- Edmund Halley, 17. st
- engleski astronom, geofizičar, matematičar, meteorolog i fizičar
- http://en.wikipedia.org/wiki/Edmond_Halley

www.aktuari.hr

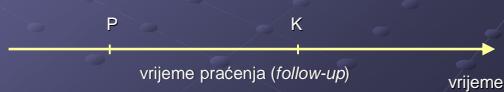
- aktuar – stručnjak koji se bavi problemima financijske neizvjesnosti i rizika koristeći matematičke metode teorije vjerojatnosti, statistike i financijske matematike
- posao – analiza podataka iz prošlosti, procjenu postojećih rizika i razvoj modela za projekciju budućih događaja
- zaposlenje – osiguranje i mirovinsko osiguranje
- znanja – matematika, ekonomija, praksa i zakoni države u kojoj radi, demografska i financijska kretanja, vještina komunikacije

Analiza preživljjenja

- psihijatrija – 1%
- patologija – 1%
- kirurgija – 12%
- onkologija – 14%
- izvorni znanstveni radovi u *The New England Journal of Medicine* – 32%
- podaci 1986.-2001., Dawson Saunders & Trapp, Basic and Clinical Biostatistics

Analiza preživljjenja

- analiza podataka vezanih uz vremensko praćenje događaja
- dve točke praćenja:
 - početak (P) (*time origin*)
 - kraj (K) (*end point*)



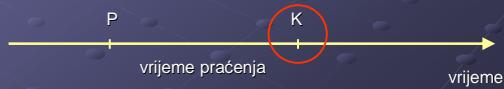
Početak praćenja

- rođenje
- pojava znaka bolesti
- postavljanje dijagnoze
- početak liječenja
- dan operativnog zahvata



Kraj praćenja

- smrt od osnovne bolesti
- smrt (svi ostali mogući uzroci)
- ponovno javljanje bolesti
- postizanje učinka liječenja
- gubitak iz uzorka (ispitne skupine)



Kraj praćenja

- smrt od osnovne bolesti
- smrt (svi ostali mogući uzroci)

uskladeno
preživljivanje
engl. *adjusted survival rate*

ukupno preživljivanje
engl. *observed survival rate*

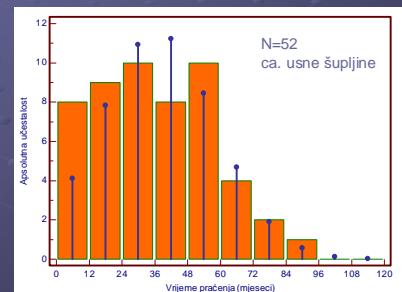
Manual for Staging of Cancer
3rd ed., AJCC

Vrijeme praćenja

- raspodjela u pravilu nije simetrična
- podaci su nepotpuni, praćenje je nepotpuno, "cenzurirano" (*censored data*)
- podaci za primjere:
 - istraživanje karcinoma usne šupljine
 - MFK KBD
 - dr. Ivica Lukšić
 - n = 52; 1. siječnja 2000. – 31. prosinca 2004.
 - reprezentativni probrani uzorak
 - dio populacije tog razdoblja
 - prva dg. karcinoma, bez regionalnih metastaza, itd.

Vrijeme praćenja (1)

- raspodjela u pravilu nije simetrična



Vrijeme praćenja (2)

- potpuni podaci (potpuno praćenje)



Vrijeme praćenja (2)

- podaci su nepotpuni, praćenje je nepotpuno, "cenzurirano"
 - cenzurirano vrijeme praćenja = jedinka tijekom praćenja ne dostiže očekivani događaj



Vrijeme praćenja (3)

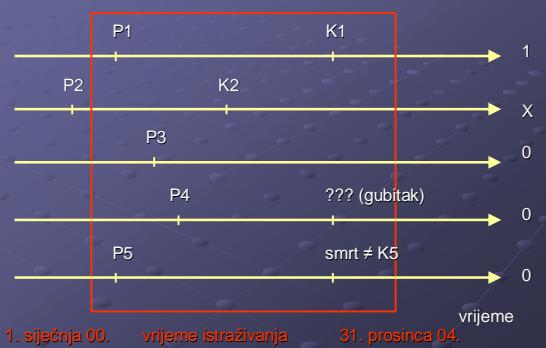
- podaci su nepotpuni, praćenje je nepotpuno, "cenzurirano"
 - cenzurirano vrijeme praćenja = jedinka tijekom praćenja ne dostiže očekivani događaj



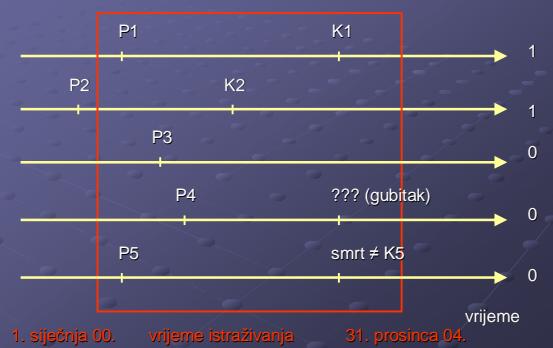
Cenzuriranje

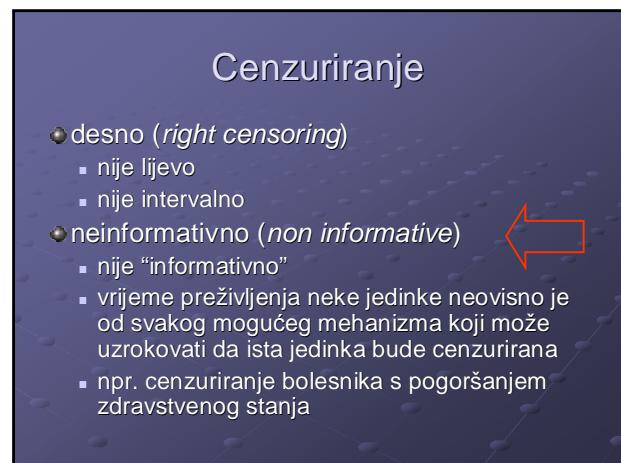
- događaj se ostvaruje = 1
- sve ostalo = 0 (cenzurirani podaci)
 - kraj istraživanja (*end of the study*)
 - gubitak iz praćenja (*lost to follow-up*)
 - ostali događaji

Cenzuriranje: bolesnici s postavljenom dijagnozom (P) u zadanih pet godina



Cenzuriranje: bolesnici liječeni u zadanim petogodišnjem razdoblju





A sada – veselje!

$$R(t) = P\{T > t\} = \int_t^{\infty} f(u) du = 1 - F(t).$$

- funkcija preživljivanja
 - biomedicina
 - survival function*
- funkcija pouzdanosti
 - inženjerstvo
 - reliability function*

S(t) ili R(t):

- vjerojatnost da će jedinka preživjeti ili točno doživjeti vrijeme od t jedinica praćenja, ili
- vjerojatnost preživljivanja jedinke u rasponu od početka praćenja do trenutka praćenja t

A sada – još veće veselje!

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{R(t) - R(t + \Delta t)}{\Delta t \cdot R(t)}$$

$$F(t) = \int_0^t f(x) dx$$

- funkcija rizika
 - hazard function*
- kumulativna funkcija rizika umiranja
 - $H(t) = -\log S(t)$

$h(t)$:

- vjerojatnost da će jedinka umrijeti u trenutku t , uz uvjet da je preživjela do toga trenutka, i uz uvjet da je
- $F(t)$ funkcija gustoće

Zaključak: podaci o preživljivanju

- vjerojatnost preživljivanja
 - $S(t)$
- rizik umiranja
 - $H(t) = -\log S(t)$

Izračunavanje preživljivanja

- neparametrijski postupci
 - Cutler-Edererov postupak (tablice preživljivanja)
 - Kaplan-Meireov postupak
- parametrijski

I. Tablice preživljenja

- osiguravateljske tablice
- tablice smrtnosti

vrijeme praćenja	vjerojatnost smrti $q=d/(n-w/2)$	vjerojatnost preživljivanja $p=1-q$	kumulativna vjerojatnost preživljavanja $S(t)=\prod_p$
0-12 mj.	0,11	0,89	0,89
13-24 mj.	0,27	0,73	0,65
25-36 mj.	0	1	0,65
37-48 mj.	0,4	0,6	0,39
49-60 mj.	0	1	0,39

Kako do preživljenja?

- upis podataka
- preuređenje podataka
- izračun podataka

1. Upis podataka, Excel®

	pacijent	datumop	datumkraj	cenzus	mjeseci
1	23456	23.6.2000	15.4.2007	0	81,8
2	24485	15.10.2003	8.11.2005	0	24,8
3	23080	25.7.2000	29.8.2004	0	49,2
4	23511	28.12.2001	15.2.2007	0	61,6
5	24188	20.2.2002	29.10.2004	0	32,3
6	22701	17.12.2003	8.6.2005	1	17,7
7	22701	17.12.2003	8.6.2005	1	17,7
8	24241	17.7.2002	29.4.2007	0	57,4
9	23480	15.5.2003	20.8.2007	0	51,2
10	22823	5.10.2000	26.9.2002	1	23,7

2. Preuređenje podataka

A	B	C	D	E	
1	pacijent	datumop	datumkraj	cenzus	mjeseci
2	24485	15.10.2003	8.11.2005	0	24,8
3	23080	25.7.2000	29.8.2004	0	49,2
4	24188	20.2.2002	29.10.2004	0	32,3
5	22701	17.12.2003	8.6.2005	1	17,7
6	23998	27.4.2000	2.2.2004	1	45,2
7	24544	9.1.2002	29.9.2003	1	20,6
8	23869	10.10.2000	16.11.2003	0	37,2
9	22819	1.3.2001	6.2.2002	1	11,2
10	22921	26.4.2004	7.12.2004	0	7,4
11	23309	9.7.2004	18.7.2006	0	24,3

vrijeme praćenja	živi na početku intervala	smrtni ishod u intervalu	cenzurirani u intervalu
0-12 mj.	10	1	1
13-24 mj.	8	2	1
25-36 mj.	5	0	2
37-48 mj.	3	1	1
49-60 mj.	1	0	1

3. Izračun podataka

vrijeme praćenja	vjerojatnost smrti $q=d/(n-w/2)$	vjerojatnost preživljivanja $p=1-q$	kumulativna vjerojatnost preživljavanja $S(t)=\prod_p$
0-12 mj.	0,11	0,89	0,89
13-24 mj.	0,27	0,73	0,65
25-36 mj.	0	1	0,65
37-48 mj.	0,4	0,6	0,39
49-60 mj.	0	1	0,39

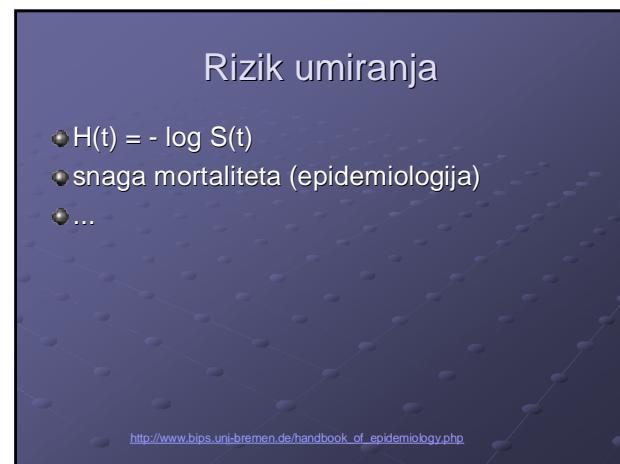
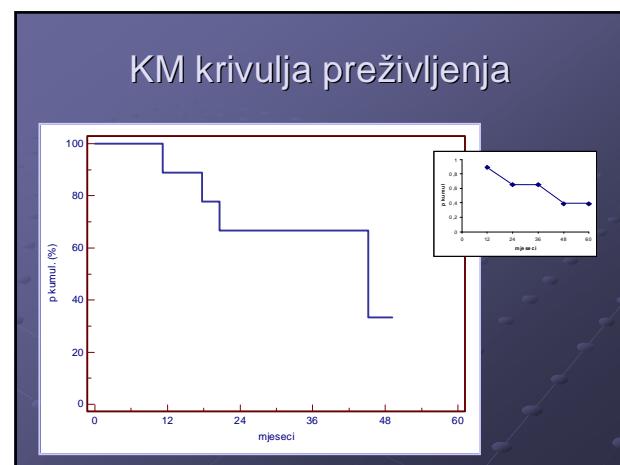
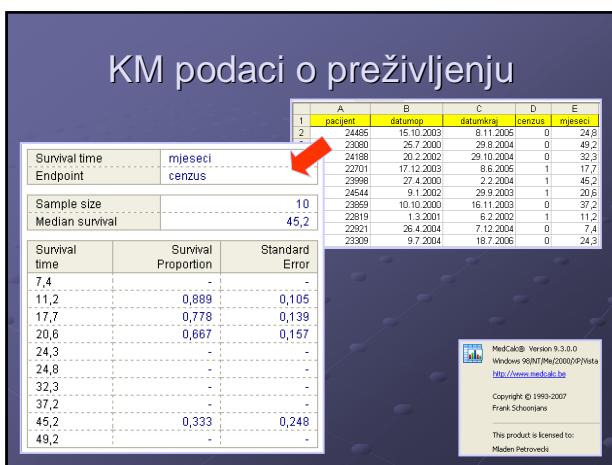
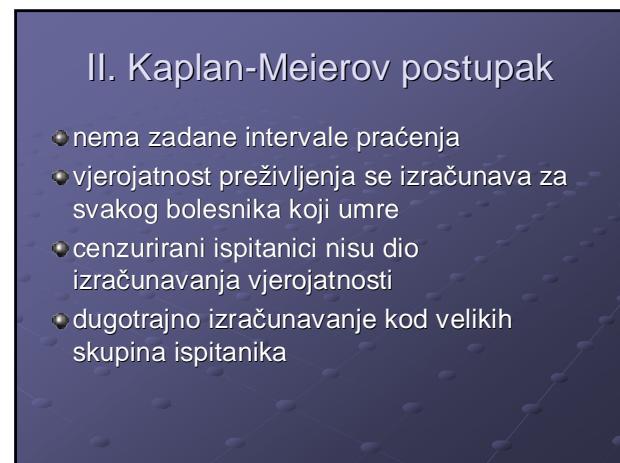
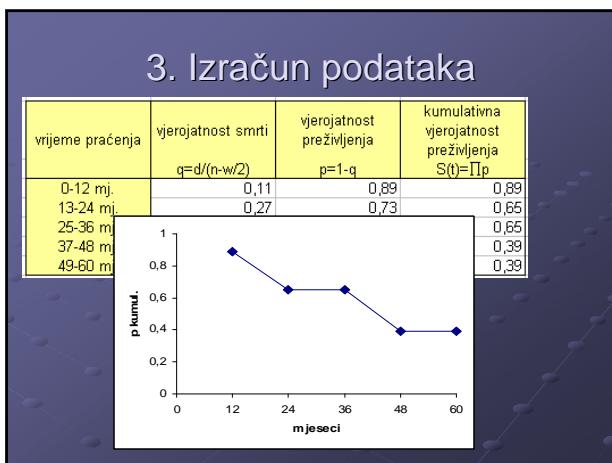
vrijeme praćenja	živi na početku intervala	smrtni ishod u intervalu	cenzurirani u intervalu
0-12 mj.	10	1	1
13-24 mj.	8	2	1
25-36 mj.	5	0	2
37-48 mj.	3	1	1
49-60 mj.	1	0	1

3. Izračun podataka

vrijeme praćenja	vjerojatnost smrti $q=d/(n-w/2)$	vjerojatnost preživljivanja $p=1-q$	kumulativna vjerojatnost preživljavanja $S(t)=\prod_p$
0-12 mj.	0,11	0,89	0,89
13-24 mj.	0,27	0,73	0,65
25-36 mj.	0	1	0,65
37-48 mj.	0,4	0,6	0,39
49-60 mj.	0	1	0,39

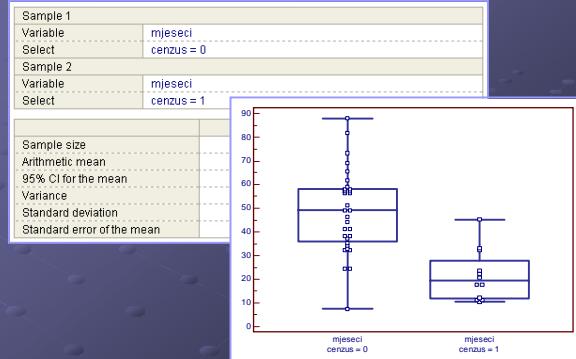
$$q = \frac{d + \frac{1}{2}wq}{n}$$

d – smrtni ishod u intervalu
n – živi na početku intervala
w – izgubljeni u intervalu

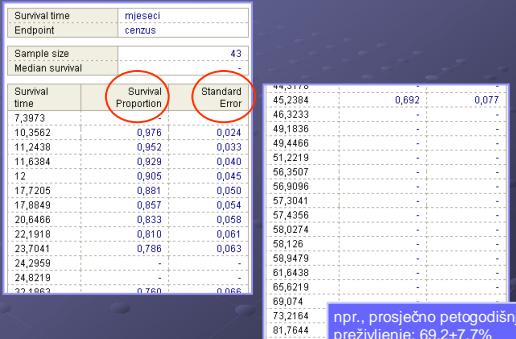


Programska potpora

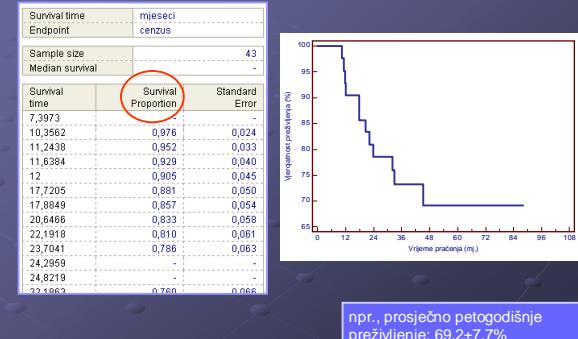
Primjer...



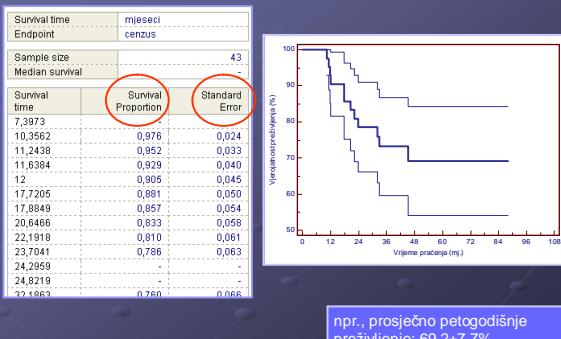
Preživljenje, MedCalc®



Krivilja preživljenja, MedCalc®



Granice pouzdanosti



Izračun granica pouzdanosti

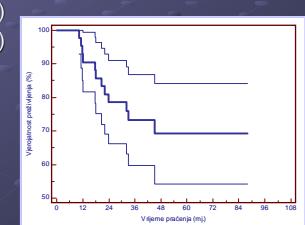
• granice pouzdanosti (*Confidence Intervals*)

• $CI = x \pm z \cdot SE(x)$

- 95%CI = $x \pm 1,96 \cdot SE(x)$
- 99%CI = $x \pm 2,56 \cdot SE(x)$

$$P(Z \leq z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du$$

<http://www.fourmilab.ch/rpkp/experiments/analysis/zCalc.html> (pazi: p/2!)

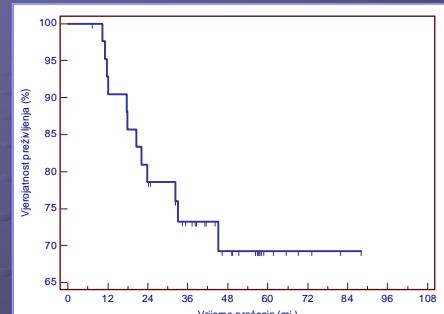


Granice pouzdanosti

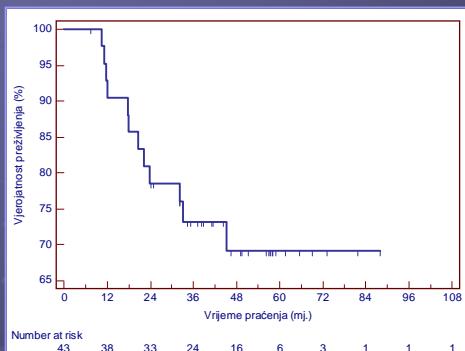
ney for at least 12 months.

Results from life table analysis are usually presented in a survival curve rather than in tables. The solid line in Fig 11–3 is a survival curve for the kidney transplant data. The dashed lines on either side of the survival curve represent 95% **confidence bands** for the curve. Although confidence bands are often not presented in journal articles, they should be included because they help readers interpret the amount of variability in the results. Typically, as the time interval from entry into the study becomes longer, the number of patients who have been in the study that long becomes increasingly smaller. The confidence

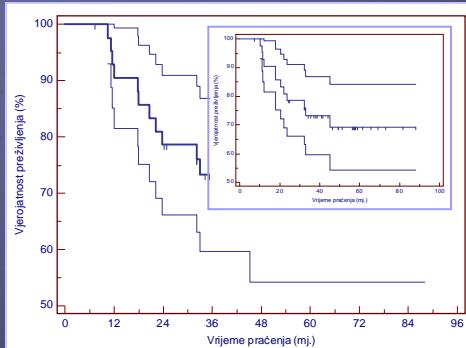
Cenzurirani podaci



Broj bolesnika pod rizikom



Standardni prikaz podataka



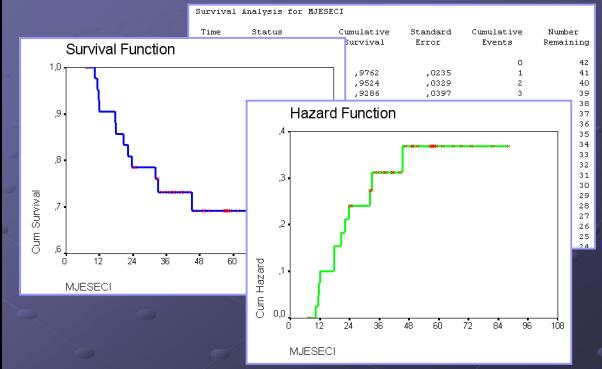
Tablice preživljenja, SPSS®

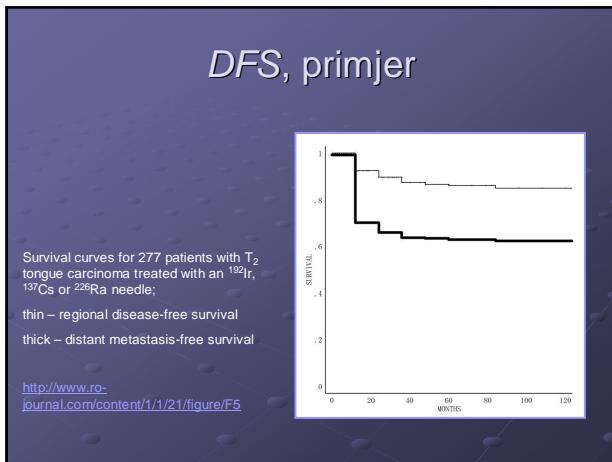
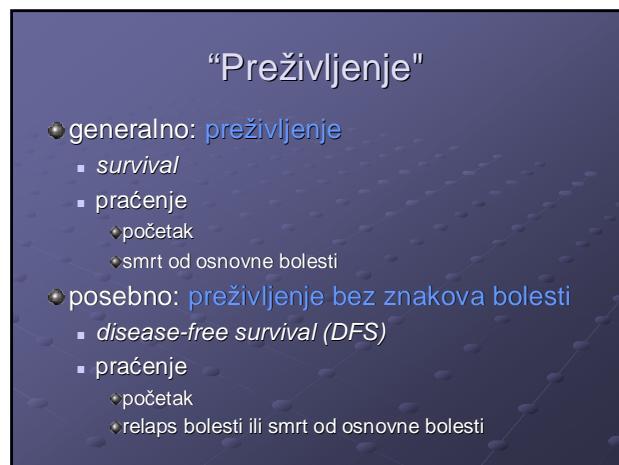
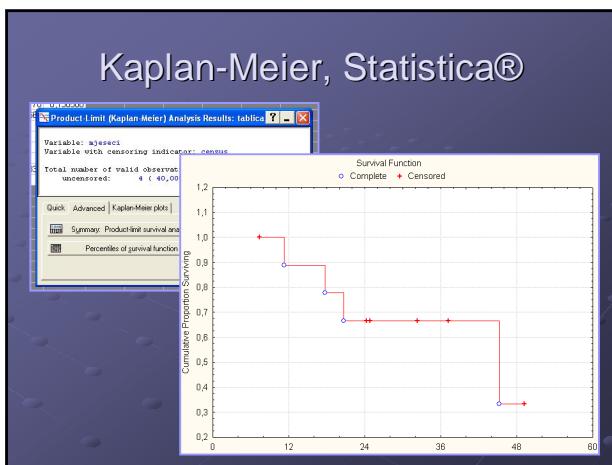
Life Table Survival Variable MJESECI										
Start Time	Intvl	Number Entering	Number Withdrawn	Number Exposed	Number of Events	Terminal Events	Events entering	Surviving at End	Hazard Rate	Number at Risk
0,0	0,0	43,0	0,0	43,0	0,0	,0000	,0000	43,0	,0000	43,0
6,0	6,0	43,0	1,0	42,5	3,0	,070	,0000	40,0	,0000	38,0
12,0	12,0	39,0	0,0	39,0	3,0	,078	,0000	36,0	,0000	36,0
18,0	18,0	36,0	0,0	36,0	3,0	,083	,0000	33,0	,0000	33,0
24,0	24,0	33,0	2,0	32,0	0,0	,0000	,0000	31,0	,0000	31,0
30,0	30,0	31,0	5,0	28,5	2,0	,070	,0000	26,0	,0000	26,0
36,0	36,0	26,0	0,0	26,0	1,0	,038	,0000	25,0	,0000	25,0
42,0	42,0	19,0	2,0	18,0	1,0	,055	,0000	17,0	,0000	17,0
48,0	48,0	16,0	3,0	14,5	,0	,0000	,0000	13,0	,0000	13,0
54,0	54,0	13,0	7,0	9,5	,0	,0000	,0000	6,0	,0000	6,0
60,0	60,0	6,0	2,0	4,0	,0	,0000	,0000	4,0	,0000	4,0
66,0	66,0	4,0	1,0	3,5	,0	,0000	,0000	3,0	,0000	3,0
72,0	72,0	3,0	1,0	2,5	,0	,0000	,0000	2,0	,0000	2,0
78,0	78,0	2,0	1,0	1,5	,0	,0000	,0000	1,0	,0000	1,0
84,0	84,0	1,0	1,0	,5	,0	,0000	,0000	0,0	,0000	0,0

The median survival time for these data is 18.0 months.

Survival Function

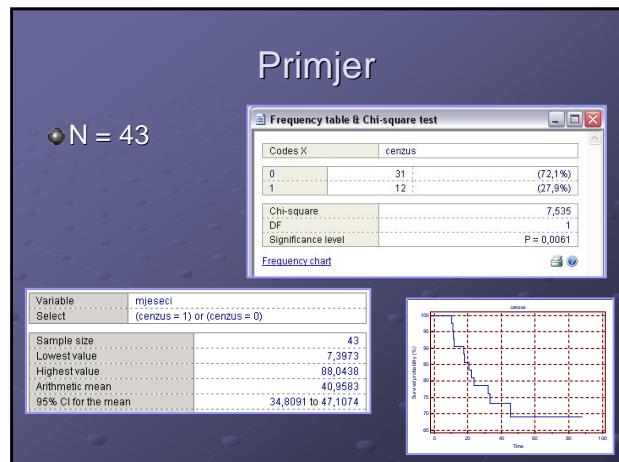
Kaplan-Meier, SPSS®

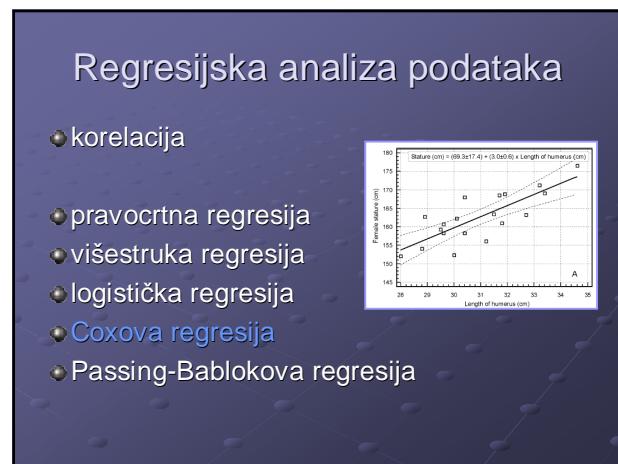
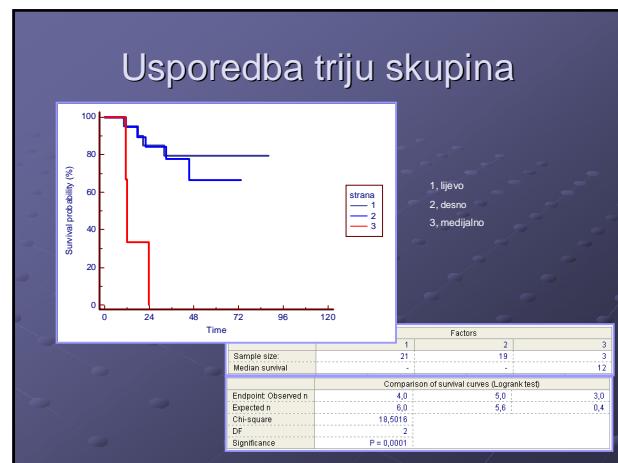
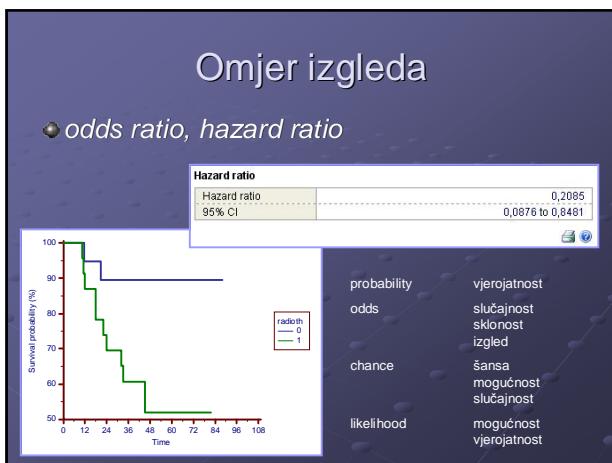
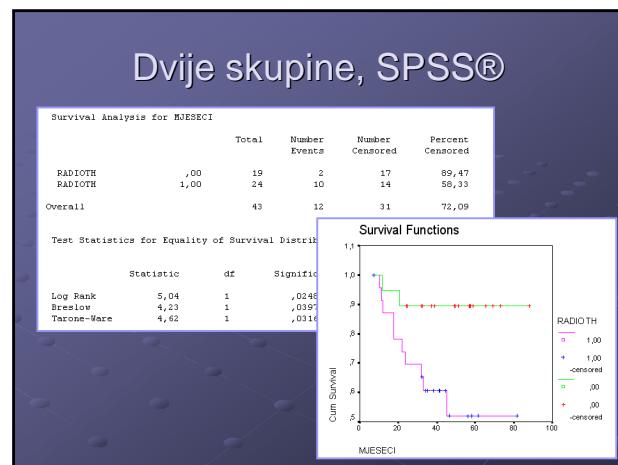
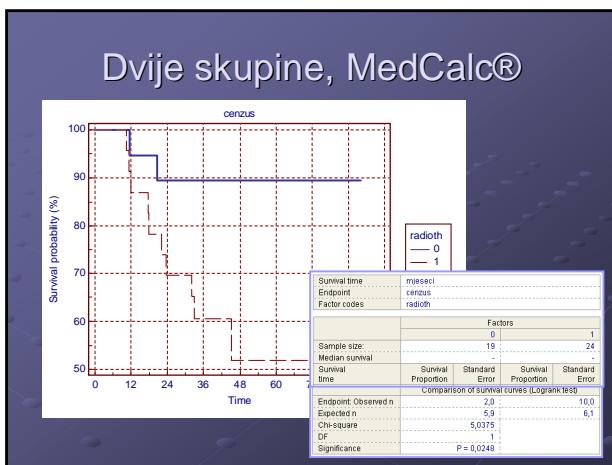




Usporediti dvije skupine...

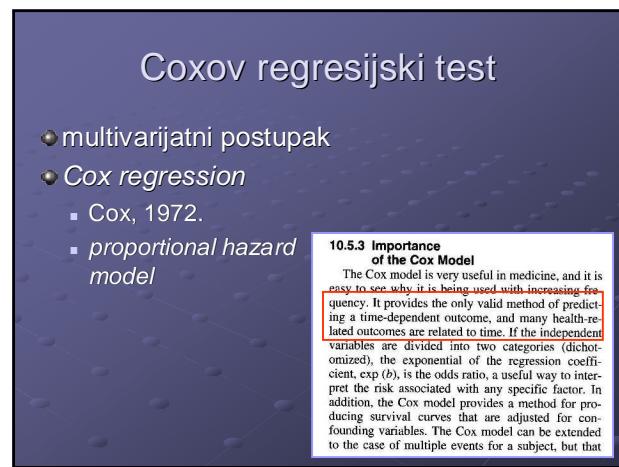
ceiving a kidney in 1984 is above the curve for patients receiving a kidney in 1978, indicating a higher proportion of patients retaining a functioning graft at any one point in time. However, variation in samples may be expected to occur simply by chance, and a reasonable question is whether the differences between the two patient cohorts is greater than expected by chance. To test this hypothesis, we need methods to compare survival distributions. If there are no censored observations, the Wilcoxon rank-sum test in-







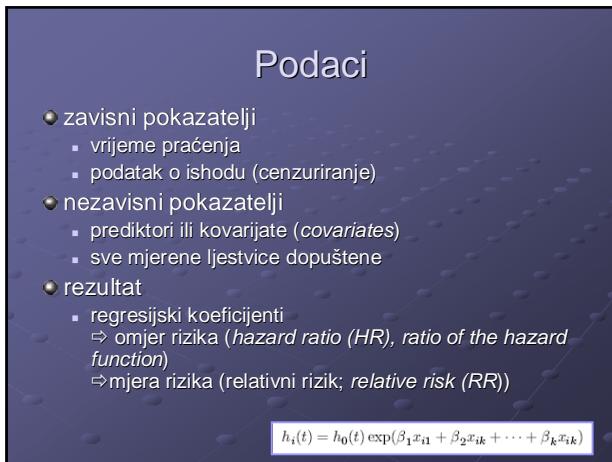
Linearni sustavi i kaos



Coxov regresijski test

- multivarijatni postupak
- Cox regression
 - Cox, 1972.
 - proportional hazard model

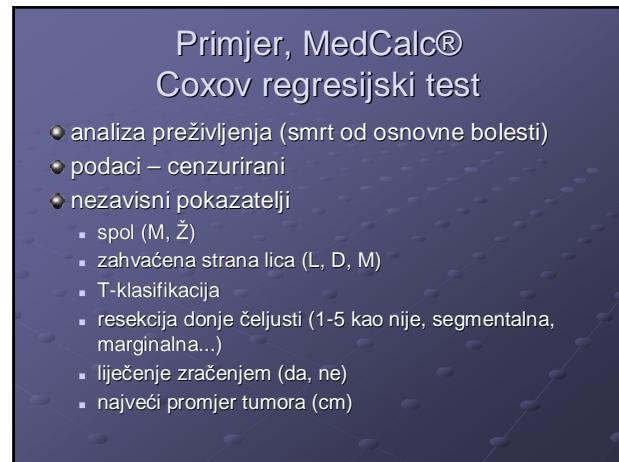
10.5.3 Importance of the Cox Model
The Cox model is very useful in medicine, and it is easy to see why it is being used with increasing frequency. It provides the only valid method of predicting a time-dependent outcome, and many health-related outcomes are related to time. If the independent variables are divided into two categories (dichotomized), the exponential of the regression coefficient, $\exp(b)$, is the odds ratio, a useful way to interpret the risk associated with any specific factor. In addition, the Cox model provides a method for producing survival curves that are adjusted for confounding variables. The Cox model can be extended to the case of multiple events for a subject, but that



Podaci

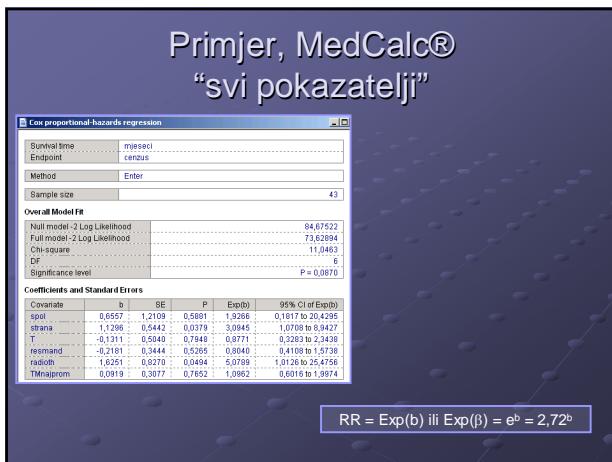
- zavisni pokazatelji
 - vrijeme praćenja
 - podatak o ishodu (cenzuriranje)
- nezavisni pokazatelji
 - prediktori ili kovarijate (covariates)
 - sve mjerene ljestvice dopuštene
- rezultat
 - regresijski koeficijenti
 - ⇒ omjer rizika (hazard ratio (HR), ratio of the hazard function)
 - ⇒ mjera rizika (relativni rizik; relative risk (RR))

$$h_i(t) = h_0(t) \exp(\beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik})$$



Primjer, MedCalc® Coxov regresijski test

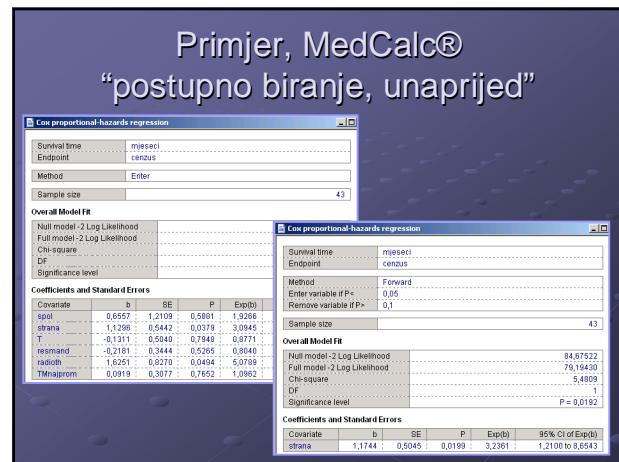
- analiza preživljivanja (smrt od osnovne bolesti)
- podaci – cenzurirani
- nezavisni pokazatelji
 - spol (M, Ž)
 - zahvaćena strana lica (L, D, M)
 - T-klasifikacija
 - resekcija donje čeljusti (1-5 kao nije, segmentalna, marginalna...)
 - liječenje zračenjem (da, ne)
 - najveći promjer tumora (cm)



Primjer, MedCalc® “svi pokazatelji”

Cox proportional-hazards regression					
Overall Model Fit					
Survival time	mjeseci				
Endpoint	cenzus				
Method	Enter				
Sample size	43				
Null model -2 Log Likelihood	84,67522				
Full model -2 Log Likelihood	73,62694				
Chi-square	11,0463				
Df	6				
Significance level	P = 0,0070				
Coefficients and Standard Errors					
Covariate	b	SE	P	Exp(b)	95% CI of Exp(b)
spol	0,6557	1,2109	0,5801	0,1936	0,1917 to 0,4395
strana	-1,1296	0,5442	0,0378	3,0945	1,0709 to 8,9427
T	-0,1311	0,5840	0,7848	0,8771	0,3283 to 2,3438
rezmand	-0,2181	0,3444	0,5265	0,8040	0,4105 to 1,5738
radioth	1,6251	0,8270	0,0494	6,0769	1,0126 to 25,4756
TMajajrom	0,0919	0,3077	0,7652	1,0862	0,6016 to 1,9874

$$RR = \text{Exp}(b) \text{ ili } \text{Exp}(\beta) = e^b = 2,72^b$$



Primjer, MedCalc® “postupno biranje, unaprijed”

Cox proportional-hazards regression					
Overall Model Fit					
Survival time	mjeseci				
Endpoint	cenzus				
Method	Forward				
Enter variable if P <	0,05				
Remove variable if P >	0,1				
Sample size	43				
Null model -2 Log Likelihood	84,67522				
Full model -2 Log Likelihood	73,19430				
Chi-square	5,4809				
Df	1				
Significance level	P = 0,0192				
Coefficients and Standard Errors					
Covariate	b	SE	P	Exp(b)	95% CI of Exp(b)
strana	1,1744	0,5045	0,0199	3,2361	1,2100 to 8,6543

