

Analiza cenzuriranih podataka i krivulje preživljjenja

Mladen Petrovečki



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

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



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Obrada podataka o preživljjenju bolesnika

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



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Analiza preživljjenja

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

Komet, 1986.
(bjedoci: 2061.)



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



www.aktuari.hr

- aktuar – stručnjak koji se bavi problemima finansijske neizvjesnosti i rizika koristeći matematičke metode teorije vjerojatnosti, statistike i finansijske 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 finansijska kretanja, vještina komunikacije



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Aktuarske tablice preživljjenja (tablice smrtnosti)

8 National Vital Statistics Reports, Vol. 54, No. 14, April 19, 2006
Table 1. Life table for the total population: United States, 2003

Age	Probability of dying between ages x to $x+1$ d_x	Number surviving to age x \bar{n}_x	Number dying between ages x to $x+1$ d_{x+1}	Person-years lived between ages x to $x+1$ t_{x+1}	Total number of person-years lived above age x T_x	Expectation of life at age x e_x
					d_{x+1}	\bar{n}_{x+1}
0-1	0.000865	100,000	687	99,934	7,748,956	77.5
1-2	0.000868	99,933	46	99,929	7,740,471	77.0
2-3	0.000831	99,925	33	99,921	7,550,181	76.1
3-4	0.000859	99,924	26	99,922	7,450,930	75.1
4-5	0.000898	99,920	20	99,919	7,351,709	74.1
5-6	0.000930	99,918	17	99,917	7,251,510	73.1
6-7	0.000951	99,917	15	99,915	7,153,329	72.1
7-8	0.000942	99,915	14	99,913	7,054,164	71.1
8-9	0.000943	99,913	14	99,912	6,954,019	70.2
9-10	0.000934	99,910	13	99,903	6,855,877	69.2
10-11	0.000916	99,906	16	99,908	6,756,754	68.2
11-12	0.000917	99,900	15	99,903	6,657,646	67.2
12-13	0.000976	99,095	17	99,077	6,559,539	66.2
13-14						66.2



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

http://en.wikipedia.org/wiki/Actuarial_table

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

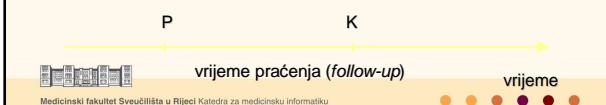


Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



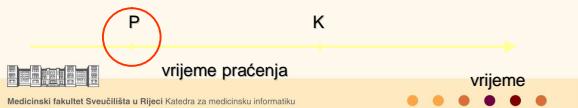
Analiza preživljjenja

- analiza podataka vezanih uz vremensko praćenje događaja
- dvije 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)

uskladeno
preživljenje
engl. *adjusted survival rate*

ukupno preživljenje
engl. *observed survival rate*

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

Manual for Staging of Cancer
TNM



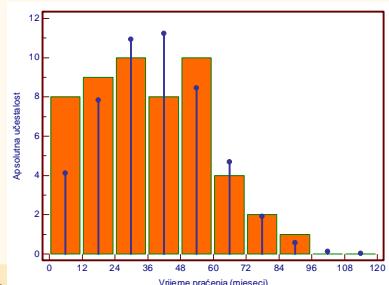
Vrijeme praćenja

- 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 probani uzorak
 - dio populacije tog razdoblja
 - prva dg. karcinoma, bez regionalnih metastaza, itd.



Vrijeme praćenja (1)

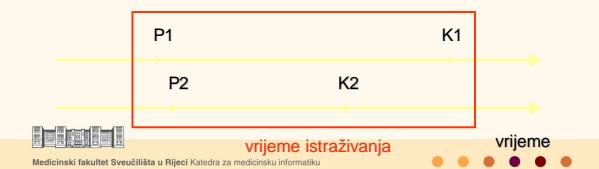
- raspodjela u pravilu nije simetrična



卷之三十一

Vrijeme praćenja (2)

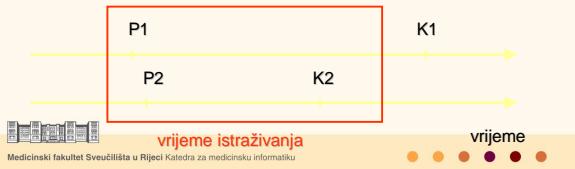
- potpuni podaci (potpuno praćenje)



vrijeme istraživanja

Vrijeme praćenja (2)

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



Vrijeme istraživanja

Vrijeme praćenja (3)

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



Vrijeme istraživanja

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

A decorative horizontal bar at the bottom of the page. It features a dark grey silhouette of a city skyline with various building heights, positioned above a thin yellow horizontal line.

“Izgubljen iz praćenja”

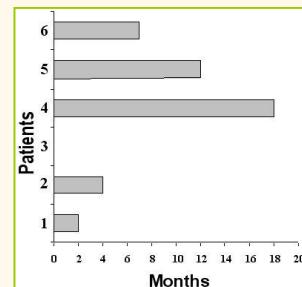
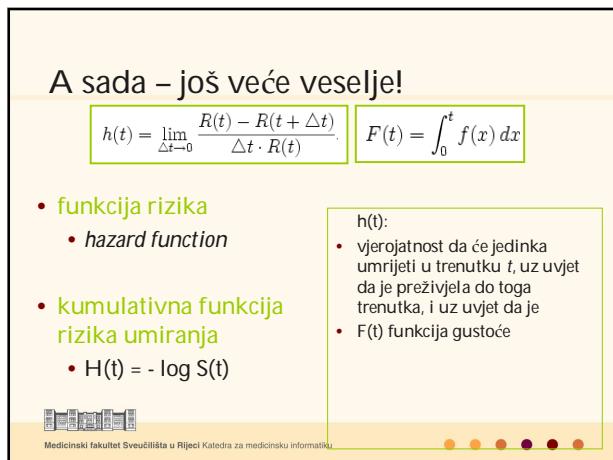
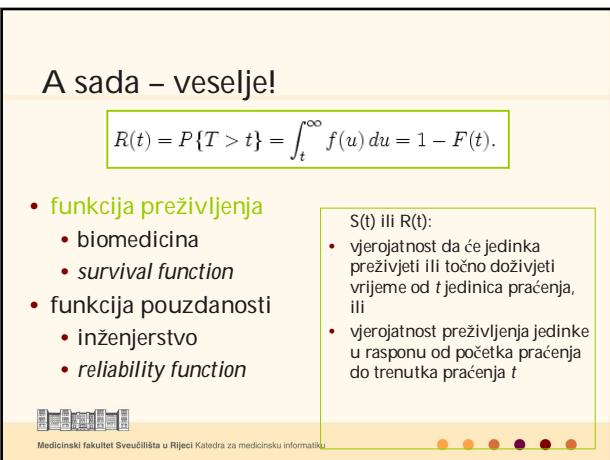
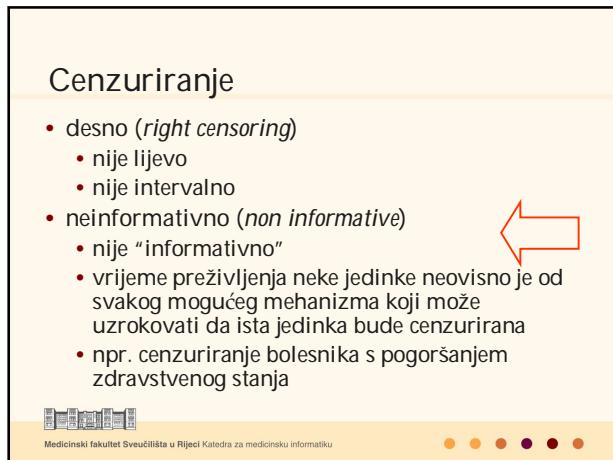
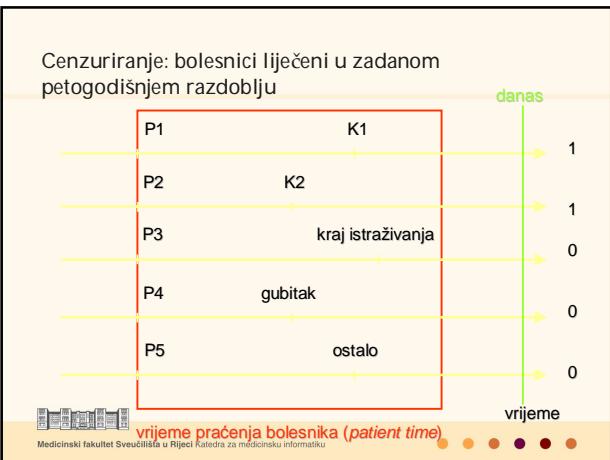
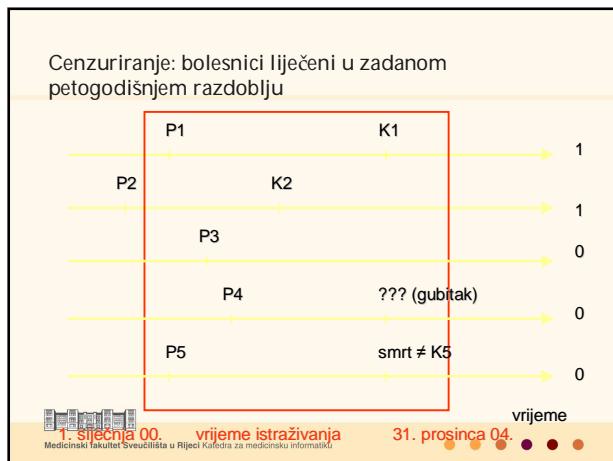


Figure 3. Outcome
The mean length of survival for our patients was 7 months.
Patient 3 was lost to follow up.



Zaključak: podaci o preživljjenju

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

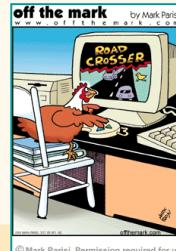


Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Izračunavanje preživljjenja

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



© Mark Parisi. Permission required for use.

I. Tablice preživljjenja

- osiguravateljske tablice
- tablice smrtnosti

vrijeme praćenja	vjerojatnost smrti $q=d/(n-wZ)$	vjerojatnost preživljjenja $p=1-q$	kumulativna vjerojatnost preživljjenja $S(t)=1-p^t$
0-12 mј.	0,11	0,89	0,89
13-24 mј.	0,27	0,73	0,65
25-36 mј.	0	1	0,65
37-48 mј.	0,4	0,6	0,39
49-60 mј.	0	1	0,39



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Kako do preživljjenja?

- upis podataka
- preuređenje podataka
- izračun podataka



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



1. Upis podataka, Excel®



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

2. Preuređenje podataka

A	B	C	D	E
1 pacijent	datumop	datumkraj	cenzus	mjeseci
2 23456	23.6.2000	15.4.2007	0	81,8
3 24485	15.10.2003	8.11.2005	0	24,8
4 23080	25.7.2000	29.8.2004	0	49,2
5 23511	28.12.2001	15.2.2007	0	61,6
6 24188	20.2.2002	29.10.2004	0	32,3
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
11 23309	9.7.2004	18.7.2006	0	24,3

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



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

3. Izračun podataka

vrijeme praćenja	vjerojatnost smrti $q=d/(n-w/2)$	vjerojatnost preživljjenja $p=1-q$	kumulativna vjerojatnost preživljjenja $S(t)=\prod p$
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 n smrtni ishod u intervalu d cenzurirani u intervalu w

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

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

3. Izračun podataka

vrijeme praćenja	vjerojatnost smrti $q=d/(n-w/2)$	vjerojatnost preživljjenja $p=1-q$	kumulativna vjerojatnost preživljjenja $S(t)=\prod p$
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

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

3. Izračun podataka

vrijeme praćenja	vjerojatnost smrti $q=d/(n-w/2)$	vjerojatnost preživljjenja $p=1-q$	kumulativna vjerojatnost preživljjenja $S(t)=\prod p$
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 m			0,39
49-60 m			0,39

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

II. Kaplan-Meierov postupak

- nema zadane intervale praćenja
- vjerojatnost preživljjenja se izračunava za svakog bolesnika koji umre
- cenzurirani ispitanici nisu dio izračunavanja vjerojatnosti
- dugotrajno izračunavanje kod velikih skupina ispitanika

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

KM podaci o preživljienju

Survival time	mjeseci	Endpoint	cenzus
Sample size	10	Median survival	45,2
Survival time	Survival Proportion	Standard Error	
7,4	-	-	
11,2	0,889	0,105	
17,7	0,778	0,139	
20,6	0,667	0,157	
24,3	-	-	
24,8	-	-	
32,3	-	-	
37,2	-	-	
45,2	0,333	0,248	
49,2	-	-	

MedCalc® Version 9.3.0.0 Windows 98/NT/ME/2000/XP/2003
Copyright © 1993-2007 Frank Schoonjans
This product is licensed to:
Mladen Petrović

KM krivulja preživljjenja

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

Kaplan, Meier

- Kaplan EL, Meier P. Nonparametric estimation from incomplete observations. J Am Stat Assoc 1958;53:457-81.



Medicina i informatica u Rijeci Katedra za medicinsku informatiku

Kaplan, Meier

- među 5 najcitanijih radova u znanosti od trenutka objavljivanja (M. Zhou, Kentucky University; <http://www.ms.uky.edu/~mai/>)

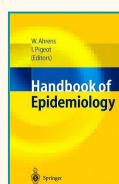
- prikaz krivulje u zavisnosti od N
<http://www.ms.uky.edu/~mai/java/stat/KapMei.html>

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Rizik umiranja

- $H(t) = -\log S(t)$
 - snaga mortaliteta (epidemiologija)
 - ...



 http://www.bips.uni-bremen.de/handbook_of_epidemiology.php



Primier

Sample 1	
Variable	mijescl
Select	census = 0
Sample 2	
Variable	mijescl
Select	census = 1
Sample size	100
Arithmetic mean	50.00
95% CI for the mean	44.75 - 55.25
Variance	100.00
Standard deviation	10.00
Standard error of the mean	1.00

The figure displays two box plots side-by-side. The left box plot represents Sample 1 (census = 0) with a median of 50, an interquartile range (IQR) from 44 to 55, and whiskers extending from 10 to 88. The right box plot represents Sample 2 (census = 1) with a median of 15, an IQR from 10 to 20, and whiskers extending from 5 to 40. Both plots show individual data points as small circles.

Preživljenje, MedCalc®

Survival time	mjeseči	cenzus	
Sample size			43
Median survival			
Survival time	Survival Proportion	Standard Error	
7,3973	0,876	0,024	
10,3562	0,876	0,024	
11,2438	0,952	0,033	
11,6384	0,829	0,040	
12	0,805	0,045	
17,7205	0,681	0,050	
17,8849	0,857	0,054	
20,6466	0,833	0,058	
22,1918	0,810	0,061	
23,7041	0,786	0,063	
24,2959	-	-	

npr., prosječno petogodišnje preživljenje: 69,2+7,7%

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku in

Krivulja preživljienja, MedCalc®

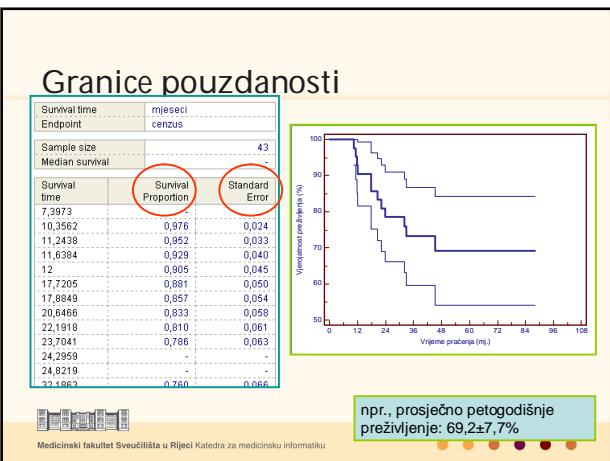
Survival time mjeseci
 Endpoint cenzus

Sample size 43
 Median survival

Survival time	Survival Proportion	Standard Error
7,3973	0,976	0,024
10,3562	0,976	0,024
11,2438	0,952	0,033
11,6394	0,929	0,040
12	0,905	0,045
17,7205	0,881	0,060
17,8849	0,857	0,054
20,6466	0,833	0,058
22,1918	0,810	0,061
23,7041	0,786	0,063
24,2959	-	-
24,8219	-	-
32,1163	0,760	0,066

npr., prosječno petogodišnje preživljjenje: $69,2 \pm 7,7\%$

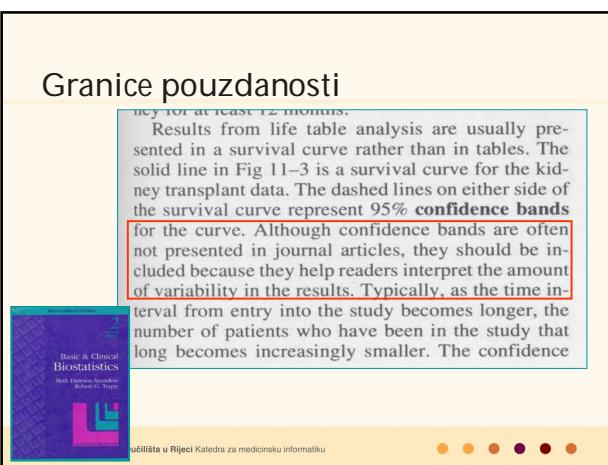
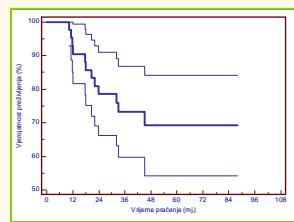
Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



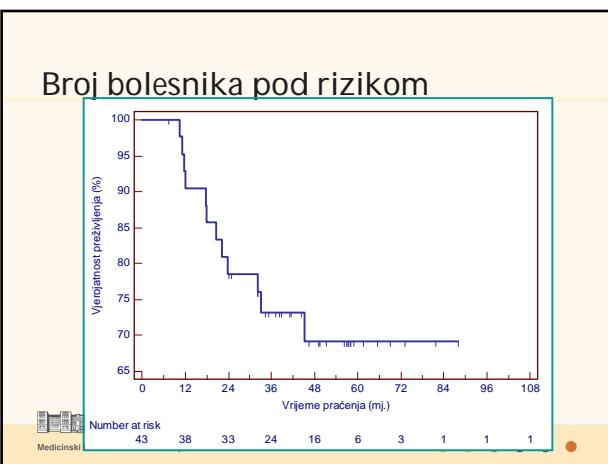
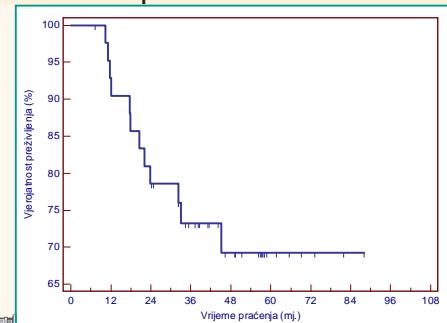
Izračun granica pouzdanosti

- granice pouzdanosti (*Confidence Intervals*)

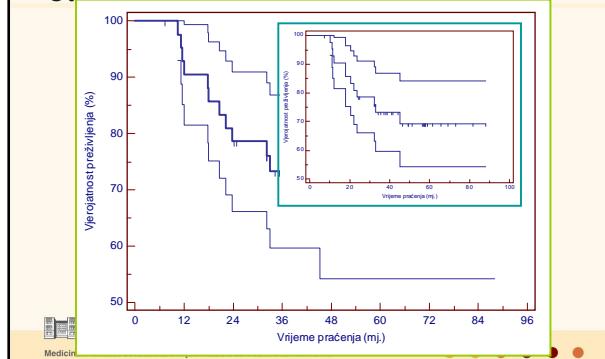
- $CI = x \pm z SE(x)$
 - $95\% CI = x \pm 1,96 SE(x)$
 - $99\% CI = x \pm 2,56 SE(x)$

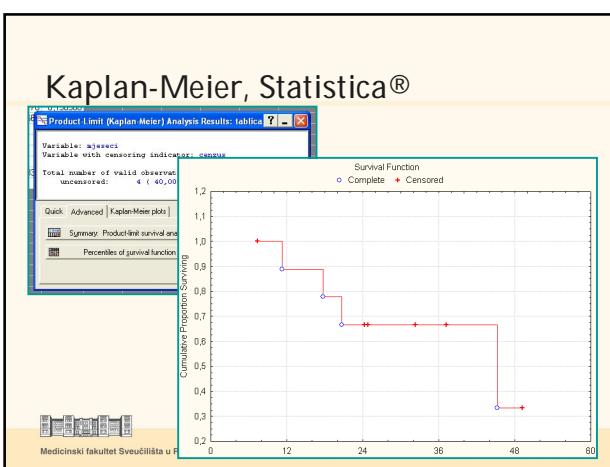
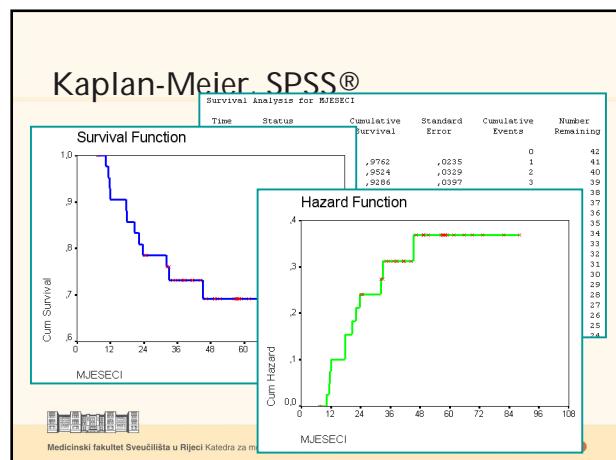
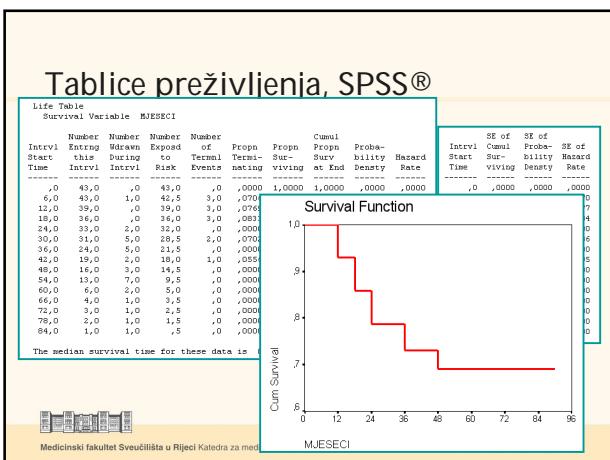


Cenzurirani podaci



Standardni prikaz podataka

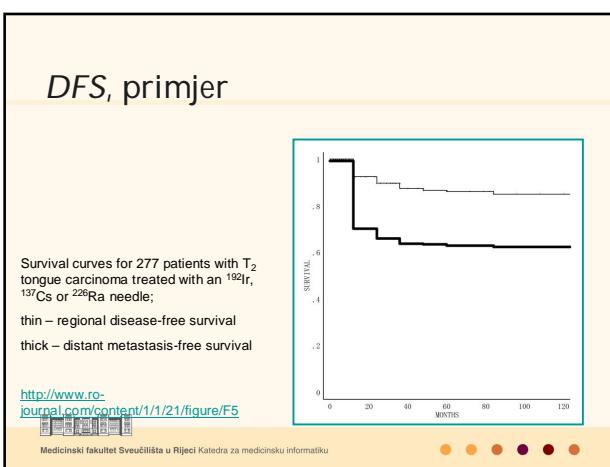




"Preživljenje"

- generalno: **preživljenje**
 - survival
 - praćenje
 - početak
 - smrt od osnovne bolesti
- posebno: **preživljenje bez znakova bolesti**
 - disease-free survival (DFS)
 - praćenje
 - početak
 - relaps bolesti ili smrt od osnovne bolesti

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Usporedba podataka o preživljenju

- usporedba dvije skupine podataka
 - log-rank (logrank) test
 - Mentelov ili Mantel-Coxov test
 - Wilcoxonov test
 - generalizirani Wilcoxonov test
 - Gehanov test
 - Gehan-Breslow-Ljhev test
 - opći Kruskal-Wallisov test za cenzurirane podatke
 - Mantel-Haenszelov test
 - Tarone-Wareov test
- usporedba triju i više skupina

Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Usporediti dvije skupine...

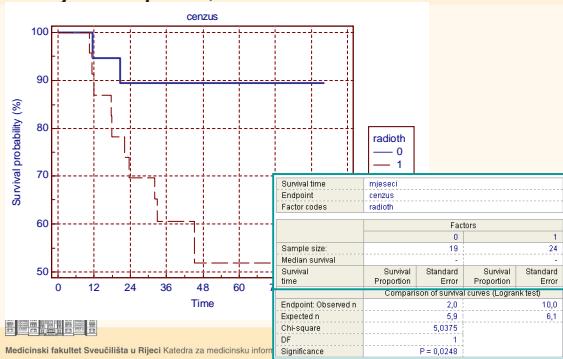
ceiving a kidney in 1984 is above the curve for patients receiving a kidney in 1978, indicating a higher proportion of patients retaining a functioning graft at any one point in time. However, variation in samples may be expected to occur simply by chance, and a reasonable question is whether the differences 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-



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



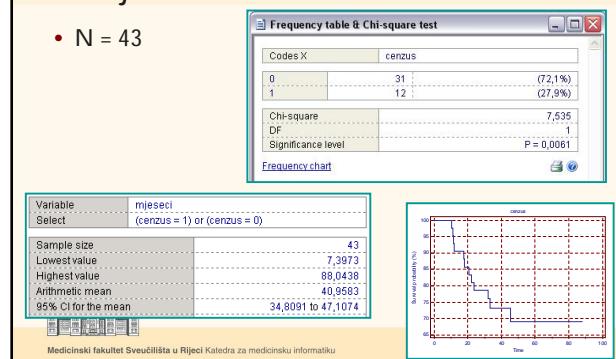
Dvije skupine, MedCalc®



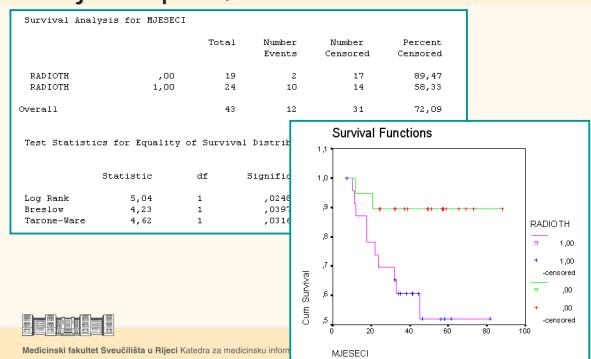
Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

Primjer

- N = 43



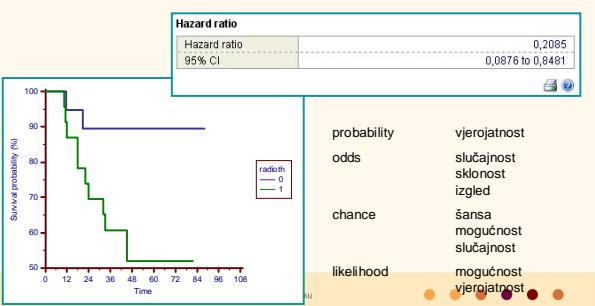
Dvije skupine, SPSS®



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

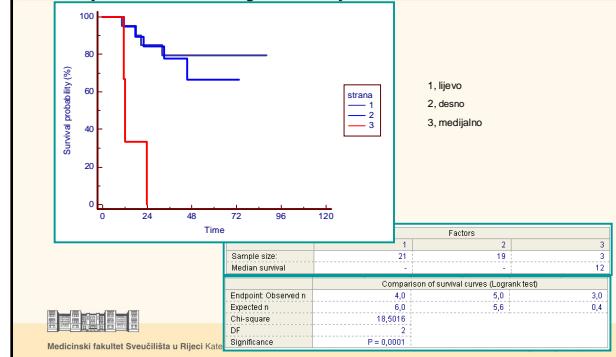
Omjer izgleda

- odds ratio, hazard ratio*



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

Usporedba triju skupina



Zaključivanje

- granice pouzdanosti
- p-vrijednosti



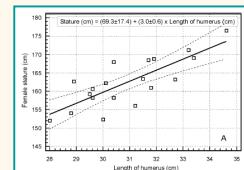
Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Regresijska analiza podataka

- korelacija

- pravocrtna regresija
- višestruka regresija
- logistička regresija
- Coxova regresija
- Passing-Bablokova regresija



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



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



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

Podaci &

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



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

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

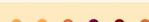
Primjer, MedCalc®

Coxov regresijski test

- analiza preživljavanja (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)



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Primjer, MedCalc®

"svi pokazatelji"

Cox proportional-hazards regression					
	Survival time	Endpoint	Method	Sample size	Overall Model Fit
	mjesec	cenzus	Enter	43	Null model -2 Log Likelihood 84,67522 Full model -2 Log Likelihood 73,62894 Chi-square 11,0463 DF 6 Significance level P = 0,0870
Coefficients and Standard Errors					
Covariate	B	SE	P	Exp(B)	95% CI of Exp(B)
spol	0,6557	1,2109	0,5801	1,9286	0,0117 to 30,4295
strana	1,1296	0,5442	0,0379	3,0945	1,0705 to 8,9427
T	-0,1311	0,5840	0,7948	0,8771	0,3263 to 2,3438
rezmand	-0,2181	0,3444	0,5265	0,8840	0,4108 to 1,5738
radioth	1,6251	0,8270	0,0494	5,0799	1,0126 to 25,4756
Thmajstrom	0,0919	0,3077	0,7652	1,0962	0,6015 to 1,9874

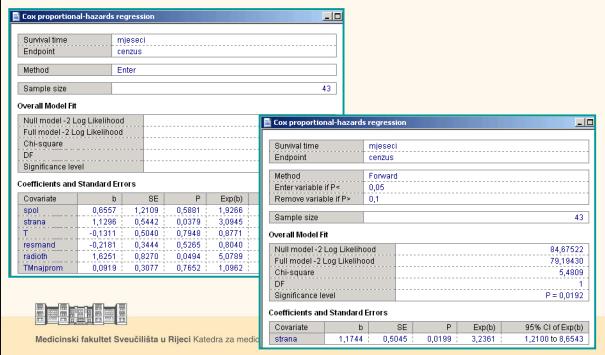


Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

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

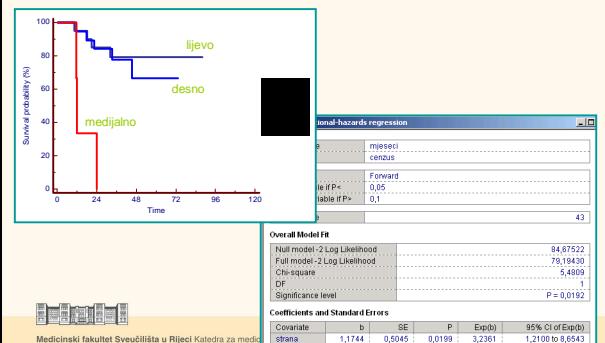


Primjer, MedCalc® "postupno biranje, unaprijed"



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

Primjer, MedCalc® "postupno biranje, unaprijed"



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku



mladenp@kbd.hr



Medicinski fakultet Sveučilišta u Rijeci Katedra za medicinsku informatiku

