

SESSION 2: WEIGHTED LOG RANK TESTS

Module 16: Survival Analysis for Clinical Trials
Summer Institute in Statistics for Clinical Research
University of Washington
June, 2017

Susanne May, Ph.D.
Professor
Department of Biostatistics
University of Washington

OVERVIEW

- Session 1
 - Review basics
 - Cox model for adjustment and interaction
 - Estimating baseline hazards and survival
- **Session 2**
 - **Weighted logrank tests**
- Session 3
 - Other two-sample tests
- Session 4
 - Choice of outcome variable
 - Power and sample size
 - Information accrual under sequential monitoring
 - Time-dependent covariates

KEY IN CLINICAL TRIALS

- Group comparisons
 - Two groups
 - k groups
 - Test for (linear) trend

- Assume, H_0 : no differences between groups

SISCR 2017: SA in Clinical Trials - SMay

2 - 3

EXAMPLE

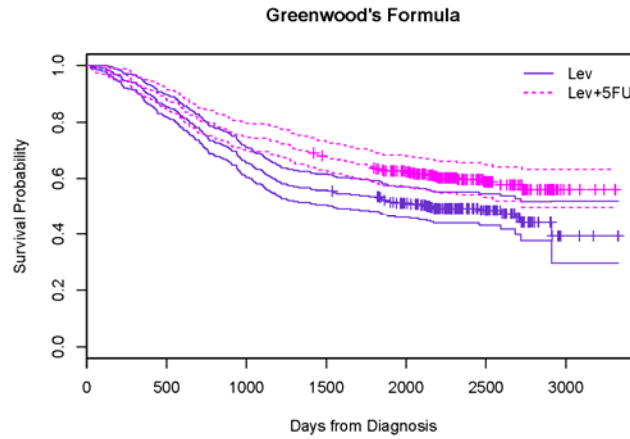
- Levamisole and Fluorouracil for adjuvant therapy of resected colon carcinoma
Moertel et al, 1990, 1995
- 1296 patients
- Stage B₂ or C
- 3 unblinded treatment groups
 - Observation only
 - Levamisole (oral, 1yr)
 - Levamisole (oral, 1yr) + fluorouracil (intravenous 1yr)

SISCR 2017: SA in Clinical Trials - SMay

2 - 4

COLON DATA EXAMPLE

- Kaplan-Meier plots and pointwise CIs



SISCR 2017: SA in Clinical Trials - SMay

2 - 5

THE P-VALUE QUESTION

- Statistical significance?

SISCR 2017: SA in Clinical Trials - SMay

2 - 6

TWO-GROUP COMPARISONS

- A number of statistical tests available
- The calculation of each test is based on a contingency table of group by status at each observed survival (event) time $t_j, j=1, \dots, m$, as shown in the Table below.

Event/Group	1	2	Total
Die	$d_{1(j)}$	$d_{2(j)}$	$D_{(j)}$
Do Not Die	$n_{1(j)} - d_{1(j)} = s_{1(j)}$	$n_{2(j)} - d_{2(j)} = s_{2(j)}$	$N_{(j)} - D_{(j)} = S_{(j)}$
At Risk	$n_{1(j)}$	$n_{2(j)}$	$N_{(j)}$

SISCR 2017: SA in Clinical Trials - SMay

2 - 7

TWO-GROUP COMPARISONS

- The contribution to the test statistic at each event time is obtained by calculating the expected number of deaths in group 1 (or 0), **assuming that the survival function is the same in each of the two groups.**
- This yields the usual **“row total times column total divided by grand total”** estimator. For example, using group 1, the estimator is

$$\hat{E}_{1(j)} = \frac{n_{1(j)} D_{(j)}}{N_{(j)}}$$

- Most software packages base their estimator of the variance on the hypergeometric distribution, defined as follows:

$$\hat{V}_{(j)} = \frac{n_{1(j)} n_{2(j)} D_{(j)} (N_{(j)} - D_{(j)})}{N_{(j)}^2 (N_{(j)} - 1)}$$

SISCR 2017: SA in Clinical Trials - SMay

2 - 8

TWO-GROUP COMPARISONS

- Each test may be expressed in the form of a ratio of weighted sums over the observed survival times as follows

$$Q = \frac{\left[\sum_{j=1}^m W_{(j)} (d_{(j)} - \hat{E}_{(j)}) \right]^2}{\sum_{j=1}^m W_{(j)}^2 \hat{V}_{(j)}}$$

- Where $j = 1, \dots, m$ are the ordered unique event times
- Under the null hypothesis and assuming that the censoring experience is independent of group, and that the total number of observed events and the sum of the expected number of events is large, then the p -value for Q may be obtained using the chi-square distribution with one degree-of-freedom,

$$p = \Pr(\chi^2(1) \geq Q)$$

WEIGHTING

- Weights used by different tests

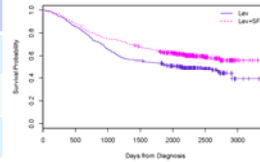
- Log Rank: $W_j = 1$
- Wilcoxon: $W_j = N_j$
- Tarone-Ware: $W_j = \sqrt{N_j}$
- Peto-Prentice: $W_j = \tilde{S}(t_{(j)})$ where $\tilde{S}(t) = \prod_{t_{(j)} \leq t} \left(\frac{N_j + 1 - D_j}{N_j + 1} \right)$
- Fleming-Harrington: $W_j = [\hat{S}(t_{(j-1)})]^p \times [1 - \hat{S}(t_{(j-1)})]^q$
 $p = q = 0 \Rightarrow W_j = 1$
 $p = 1, q = 0 \Rightarrow W_j =$ Kaplan-Meier estimate at previous survival time
- and $\hat{S}(t_{(j-1)})$ is the Kaplan-Meier estimator at time t_{j-1}

Most frequently used test weights later times relatively more heavily, while Wilcoxon weights early times more heavily

COLON CANCER EXAMPLE

- Comparing Lev vs Lev+5FU

Group	N	Obs	Exp
Lev	310	161	136.9
Lev+5FU	304	123	147.1
Total	614	284	284.0



- Log-rank test: $\chi^2(1) = 8.2$, p-value = 0.0042
- Peto-Prentice: $\chi^2(1) = 7.6$, p-value = 0.0058
- Wilcoxon: $\chi^2(1) = 7.3$, p-value = 0.0069
- Tarone-Ware: $\chi^2(1) = 7.7$, p-value = 0.0055
- Flem-Harr(1,.0): $\chi^2(1) = 7.6$, p-value = 0.0056
- Flem-Harr(1,.3): $\chi^2(1) = 9.5$, p-value = 0.0020

SISCR 2017: SA in Clinical Trials - SMay

2 - 11

- Example where choice of weights makes a difference

SISCR 2017: SA in Clinical Trials - SMay

2 - 12

EXAMPLE: LOW BIRTH WEIGHT INFANTS

- Data from UMass
- Goal: determine factors that predict the length of time low birth weight infants (<1500 grams) with bronchopulmonary dysplasia (BPD) were treated with oxygen
- Note: observational study, not clinical trial
- 78 infants total, 35 (43 not) receiving surfactant replacement therapy
- Outcome variable: total number of days the baby required supplemental oxygen therapy

SUMMARY STATISTICS - LBWI

- The estimated median number of days of therapy
 - for those babies **who did not have** surfactant replacement therapy
 - 107 {95% CI: (71, 217)},
 - for those **who had** the therapy is
 - 71 {95% CI: (56, 110)}
- The median number of days of therapy for the babies not on surfactant is about 1.5 times longer than those using the therapy.

TWO-GROUP COMPARISONS LBWI

- Different weighting approaches

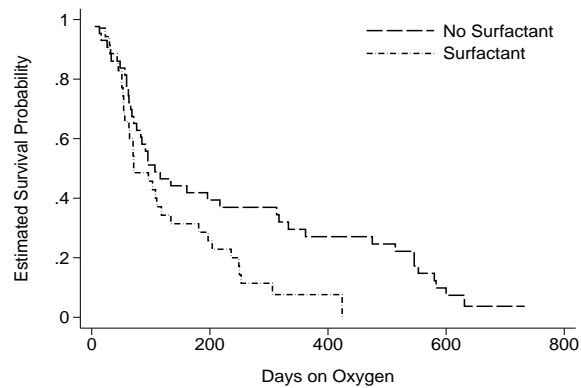
Test	Statistic	p – value
Log-rank	5.62	0.018
Wilcoxon	2.49	0.115
Tarone-Ware	3.70	0.055
Peto-Prentice	2.53	0.111
Flem-Harr(1,0)	2.66	0.103
Flem-Harr(0,1)	9.07	0.0026

SISCR 2017: SA in Clinical Trials - SMay

2 - 15

EXAMPLE: LBWI

- Kaplan-Meier plot



SISCR 2017: SA in Clinical Trials - SMay

2 - 16

WEIGHTS

- Determine weights up front
- Clinical considerations
- Ordinarily: No weights = log rank test

TRIALS WHERE WEIGHTS ARE IMPORTANT ?

- Question: Examples of settings where log rank and Cox model
 - Might be inappropriate?
 - Have low power?

- K – groups

K-GROUPS

- K-Group Comparisons

Group	1	2	...	k	...	K	Total
Die	$d_{1(j)}$	$d_{2(j)}$...	$d_{k(j)}$...	$d_{K(j)}$	$D_{(j)}$
Not Die	$s_{1(j)}$	$s_{2(j)}$...	$s_{k(j)}$...	$s_{K(j)}$	$S_{(j)}$
At Risk	$n_{1(j)}$	$n_{2(j)}$...	$n_{k(j)}$...	$n_{K(j)}$	$N_{(j)}$

- In a manner similar to the two-group case, we estimate the expected number of events for each group under an assumption of equal survival functions as

$$\hat{E}_{k(j)} = \frac{D_{(j)} n_{k(j)}}{N_{(j)}}, \quad k = 1, 2, \dots, K$$

K-GROUP COMPARISON

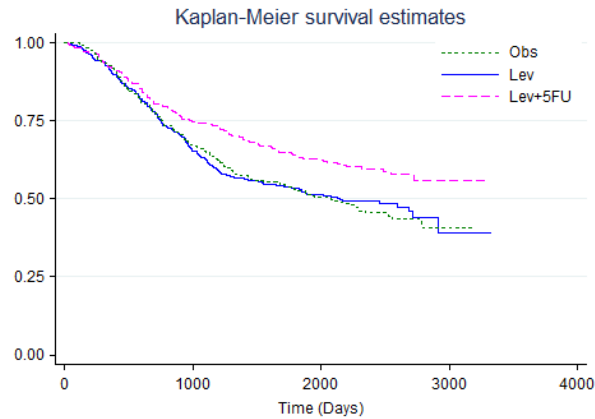
- Again, compare observed vs expected
- Quadratic form Q
- Under the null hypothesis and if the summed estimated expected number of events is large
- Test statistic $p = \Pr(\chi^2(K-1) \geq Q)$

COLON CANCER EXAMPLE

- **Obs vs Lev vs Lev+5FU**
- Log-rank test: $\chi^2(2) = 11.7$, p-value = 0.0029
- Wilcoxon: $\chi^2(2) = 9.7$, p-value = 0.0078
- Peto-Prentice: $\chi^2(2) = 10.3$, p-value = 0.0059
- Tarone-Ware: $\chi^2(2) = 10.6$, p-value = 0.0049
- Flem-Harr(1,0): $\chi^2(2) = 10.4$, p-value = 0.0056
- Flem-Harr(1,.3): $\chi^2(2) = 13.7$, p-value = 0.0011

COLON CANCER EXAMPLE

- Obs vs Lev vs Lev+5FU



SISCR 2017: SA in Clinical Trials - SMay

2 - 23

TREND TEST – EXAMPLE 1 (COLON)

- Obs vs Lev vs Lev+5FU
- Coding ?
- Pretend you did not see any results yet ...

SISCR 2017: SA in Clinical Trials - SMay

2 - 24

TREND TEST

- H_0 : survival functions are equal
- H_A : survival functions are rank-ordered and follow the trend specified by a vector of coefficients

- Examples
 - Drug dosing
 - Age

SISCR 2017: SA in Clinical Trials - SMay

2 - 25

TREND ANALYSIS

- Trend test

Groups				
Obs	0			
Lev	1			
Lev+5FU	2			
	p - value			
Log-rank				
Wilcoxon				
Tarone-Ware				
Peto-Prentice				

SISCR 2017: SA in Clinical Trials - SMay

2 - 26

TREND ANALYSIS

- Trend test

Groups				
Obs	0			
Lev	1			
Lev+5FU	2			
	<i>p</i> – value			
Log-rank	0.002			
Wilcoxon	0.007			
Tarone-Ware	0.004			
Peto-Prentice	0.005			

SISCR 2017: SA in Clinical Trials - SMay

2 - 27

TREND ANALYSIS

- Trend test

Groups				
Obs	0	0		
Lev	1	0.25		
Lev+5FU	2	1		
	<i>p</i> – value			
Log-rank	0.002	0.0007		
Wilcoxon	0.007	0.002		
Tarone-Ware	0.004	0.001		
Peto-Prentice	0.005	0.002		

SISCR 2017: SA in Clinical Trials - SMay

2 - 28

TREND ANALYSIS

- Trend test

Groups				
Obs	0	0	0	
Lev	1	0.25	0.75	
Lev+5FU	2	1	1	
	<i>p</i> – value			
Log-rank	0.002	0.0007	0.01	
Wilcoxon	0.007	0.002	0.008	
Tarone-Ware	0.004	0.001	0.02	
Peto-Prentice	0.005	0.002	0.02	

SISCR 2017: SA in Clinical Trials - SMay

2 - 29

TREND ANALYSIS

- Trend test

Groups				
Obs	0	0	0	0
Lev	1	0.25	0.75	?
Lev+5FU	2	1	1	1
	<i>p</i> – value			
Log-rank	0.002	0.0007	0.01	0.79
Wilcoxon	0.007	0.002	0.008	0.96
Tarone-Ware	0.004	0.001	0.02	0.87
Peto-Prentice	0.005	0.002	0.02	0.93
Flem-Harr(1,.3)	0.0007	0.0002	0.004	0.69

SISCR 2017: SA in Clinical Trials - SMay

2 - 30

- Another example regarding trend

TREND – EXAMPLE 2

- Thomas et al. (1977)
- Also Marubini and Valsecchi (1995, p 126)
- 29 Animals
- 3 level of carcinogenic agent (0, 1.5, 2.0)
- Outcome: time to tumor formation

Group	Dose	N	Times to event (<i>t</i>) or censoring (<i>t</i> +)
0	0	9	73+,74+,75+,76,76,76+,99,166,246+
1	1.5	10	43+,44+,45+,67,68+,136,136,150,150,150
2	2.0	10	41+,41+,47,47+,47+,58,58,58,100+,117

TREND TEST

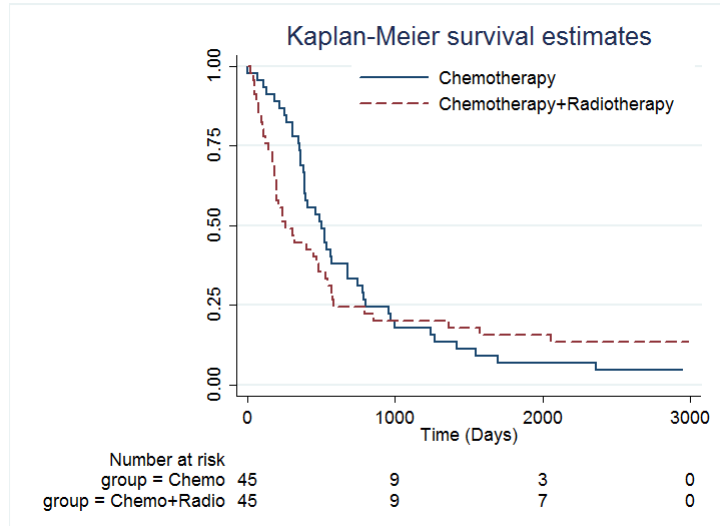
- Dose example, 29 animals

Test (Group differences)	df	Chi2	P-value
Log-rank	2	8.05	0.018
Wilcoxon	2	9.04	0.011
Trend test			
Log-rank (1,2,3)	1	5.87	0.015
Wilcoxon (1,2,3)	1	6.26	0.012
Log-rank (0,1.5,2)	1	3.66	0.056
Wilcoxon (0,1.5,2)	1	3.81	0.051

EXAMPLE 3

- Stablein and Koutrouvelis (1985)
- Gastrointestinal Tumor Study Group (1982)
- Chemotherapy vs.
Chemotherapy and Radiotherapy
- 90 patients (45 per group)

KAPLAN-MEIER SURVIVAL CURVES



SISCR 2017: SA in Clinical Trials - SMay

2 - 35

TEST STATISTICS – EXAMPLE 3

Test	Statistic	p – value
Log-rank		?
Wilcoxon		?
Peto-Prentice		?
Tarone-Ware		?
FI-Ha(1,0)		?
FI-Ha(0,1)		?

SISCR 2017: SA in Clinical Trials - SMay

2 - 36

TEST STATISTICS – EXAMPLE 3

Test	Statistic	p – value
Log-rank	0.23	0.64
Wilcoxon		
Peto-Prentice		
Tarone-Ware		
FI-Ha(1,0)		
FI-Ha(0,1)		

SISCR 2017: SA in Clinical Trials - SMay

2 - 37

TEST STATISTICS – EXAMPLE 3

Test	Statistic	p – value
Log-rank	0.23	0.64
Wilcoxon	3.96	0.047
Peto-Prentice		
Tarone-Ware		
FI-Ha(1,0)		
FI-Ha(0,1)		

SISCR 2017: SA in Clinical Trials - SMay

2 - 38

TEST STATISTICS – EXAMPLE 3

Test	Statistic	p – value
Log-rank	0.23	0.64
Wilcoxon	3.96	0.047
Peto-Prentice	4.00	0.046
Tarone-Ware	1.90	0.17
FI-Ha(1,0)		
FI-Ha(0,1)		

SISCR 2017: SA in Clinical Trials - SMay

2 - 39

TEST STATISTICS – EXAMPLE 3

Test	Statistic	p – value
Log-rank	0.23	0.64
Wilcoxon	3.96	0.047
Peto-Prentice	4.00	0.046
Tarone-Ware	1.90	0.17
FI-Ha(1,0)	2.59	0.11
FI-Ha(0,1)	4.72	0.03

SISCR 2017: SA in Clinical Trials - SMay

2 - 40

TEST STATISTICS – EXAMPLE 3

Test	Statistic	p – value
Log-rank	0.23	0.64
Wilcoxon	3.96	0.047
Peto-Prentice	4.00	0.046
Tarone-Ware	1.90	0.17
FI-Ha(1,0)	2.59	0.11
FI-Ha(0,1)	4.72	0.03

- Why the difference?

SISCR 2017: SA in Clinical Trials - SMay

2 - 41

GROUP COMPARISONS

- $H_0: S_1(t) = S_2(t) \quad \lambda_1(t) = \lambda_2(t)$
- Possible alternative
 - Survival function: $S_2(t) = S_1(t)^C, C \neq 1$
 - Hazard function: $\lambda_2(t) = C\lambda_1(t), C \neq 1$
 $\ln(\lambda_2(t)) = \ln(\lambda_1(t)) + C, C \neq 1$
- Log-rank test most powerful
if hazards are proportional

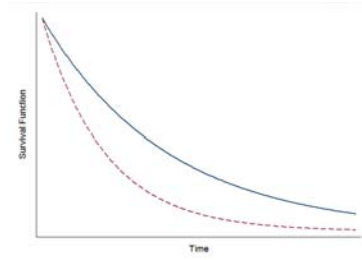
SISCR 2017: SA in Clinical Trials - SMay

2 - 42

SURVIVAL FUNCTIONS

- We can detect

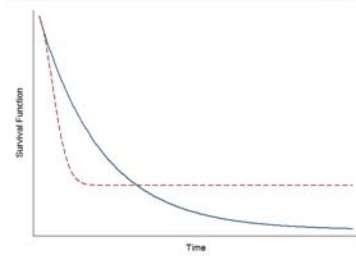
this



proportional

(generated as 2 exponential distributions)

but ordinarily not this



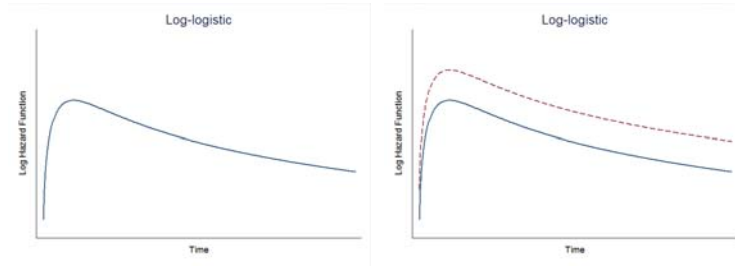
not proportional

PROPORTIONAL HAZARDS

- Easier to visualize on log hazard scale

GROUP COMPARISONS

- Proportional hazards – use log hazards scale
- Example: log-logistic survival times
- Hazards plotted on log scale



SISCR 2017: SA in Clinical Trials - SMay

2 - 45

SO FAR

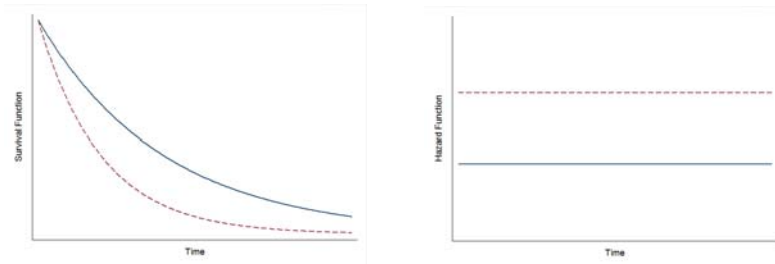
- Two and K – group comparisons
- Trend tests
- Non-parametric
- Did not make use of actual values of time

SISCR 2017: SA in Clinical Trials - SMay

2 - 46

PARAMETRIC MODELS

- Control group: Exponential(0.5)
- Example
- Survival functions Hazard functions

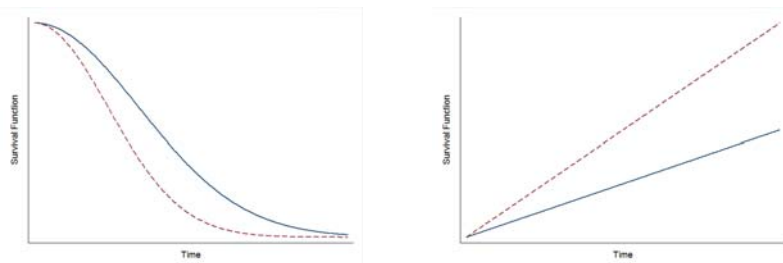


SISCR 2017: SA in Clinical Trials - SMay

2 - 47

PARAMETRIC MODELS

- Control group: Weibull(0.5,2)
- Example
- Survival Functions Hazard Functions

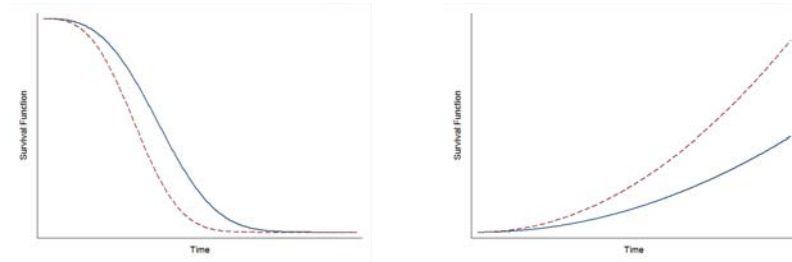


SISCR 2017: SA in Clinical Trials - SMay

2 - 48

PARAMETRIC MODELS

- Control group: Weibull(0.5,3)
- Example
- Survival Functions Hazard Functions



SISCR 2017: SA in Clinical Trials - SMay

2 - 49

PARAMETRIC APPROACHES

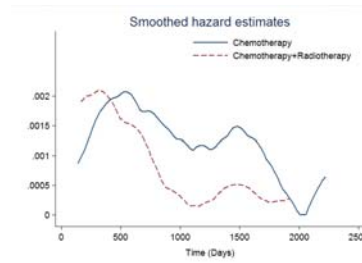
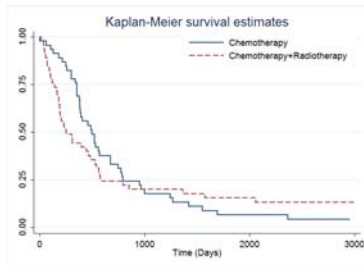
- Weibull and exponential
 - Proportional hazards assumption
 - Distributional assumptions

SISCR 2017: SA in Clinical Trials - SMay

2 - 50

BACK TO EXAMPLE 3

- Gastrointestinal Tumor Study
- Survival Functions
- Hazard Functions



SISCR 2017: SA in Clinical Trials - SMay

2 - 51

- Other covariates

SISCR 2017: SA in Clinical Trials - SMay

2 - 52

EXAMPLE 1: COLON CANCER – REVISITED

- Tumor differentiation and survival

Group	Observed Events	Expected Events
Well	42	47.5
Moderate	311	334.9
Poor	88	58.6
	441	441

- $\chi(2) = 17.2$,
- p – value = 0.0002

SISCR 2017: SA in Clinical Trials - SMay

2 - 53

EXAMPLE 1 REVISITED

- Tumor differentiation by treatment group

Groups	Obs	Lev	Lev+5FU	Total
Well	27	37	29	93
Moderate	229	219	215	663
Poor	52	44	54	150
Total	308	300	298	906

SISCR 2017: SA in Clinical Trials - SMay

2 - 54

STRATIFIED LOG-RANK TEST

- Assume R strata ($r = 1, \dots, R$)
- Recall (non-stratified) log-rank test statistic

$$Q = \frac{\left[\sum_{j=1}^m (d_{1(j)} - \hat{E}_{1(j)}) \right]^2}{\sum_{j=1}^m \hat{V}_{(j)}}$$

- Stratified log-rank test

$$Q = \frac{\left[\sum_{j=1}^{m_1} (d_{11(j)} - \hat{E}_{11(j)}) + \dots + \sum_{j=1}^{m_r} (d_{1r(j)} - \hat{E}_{1r(j)}) + \dots + \sum_{j=1}^{m_R} (d_{1R(j)} - \hat{E}_{1R(j)}) \right]^2}{\sum_{j=1}^{m_1} \hat{V}_{1(j)} + \dots + \sum_{j=1}^{m_r} \hat{V}_{r(j)} + \dots + \sum_{j=1}^{m_R} \hat{V}_{R(j)}}$$

SISCR 2017: SA in Clinical Trials - SMay

2 - 55

STRATIFIED LOG-RANK TEST

- H_0 : $\lambda_{1r}(t) = \lambda_{2r}(t)$ for all $r = 1, \dots, R$
- H_A : $\lambda_{1r}(t) = c\lambda_{2r}(t), c \neq 1$ for all $r = 1, \dots, R$
- Under H_0 test statistic $\sim \chi^2(K-1)$
- The $d_{1r(j)}, \hat{E}_{1r(j)}$ and $\hat{V}_{r(j)}$ are solely based on subjects from the r -th strata

SISCR 2017: SA in Clinical Trials - SMay

2 - 56

STRATIFIED LOG-RANK TEST

Well differentiated	Observed Events	Expected Events
Obs	18	16.7
Lev	16	10.6
Lev+5FU	8	14.7
	42	42

Moderately differentiated	Observed Events	Expected Events
Obs	109	98.7
Lev	115	105.4
Lev+5FU	87	106.9
	311	311.0

SISCR 2017: SA in Clinical Trials - SMay

2 - 57

STRATIFIED LOG-RANK TEST

Poorly differentiated	Observed Events	Expected Events
Obs	27	24.8
Lev	34	30.5
Lev+5FU	27	32.7
	88	88.0

Combined over differentiation strata	Observed Events	Expected Events
Obs	154	140.1
Lev	165	146.5
Lev+5FU	122	154.4
	441	441.0

- $\chi(2) = 10.5$
- P-value: 0.005

SISCR 2017: SA in Clinical Trials - SMay

2 - 58

COMPARISON STRATA VS NO STRATA

- $\chi(2) = 10.5$
- P-value: 0.005

Combined over differentiation strata	Observed Events	Expected Events
Obs	154	140.1
Lev	165	146.5
Lev+5FU	122	154.4
	441	441.0

- $\chi(2) = 11.7$
- P-value: 0.003

Without strata	Observed Events	Expected Events
Obs	161	146.1
Lev	168	148.4
Lev+5FU	123	157.5
	452	452

SISCR 2017: SA in Clinical Trials - SMay

2 - 59

COMPARISON STRATA VS NO STRATA

- Why are the observed and expected different?

SISCR 2017: SA in Clinical Trials - SMay

2 - 60

COMPARISON STRATA VS NO STRATA

- Why are the observed and expected different?
- Answer: There are 23 individuals with missing differentiation level

SISCR 2017: SA in Clinical Trials - SMay

2 - 61

(FAIR) COMPARISON STRATA VS NO STRATA

- $\chi(2) = 10.5$
- P-value: 0.005

Combined over differentiation strata	Observed Events	Expected Events
Obs	154	140.1
Lev	165	146.5
Lev+5FU	122	154.4
	441	441.0

- $\chi(2) = 10.6$
- P-value: 0.005

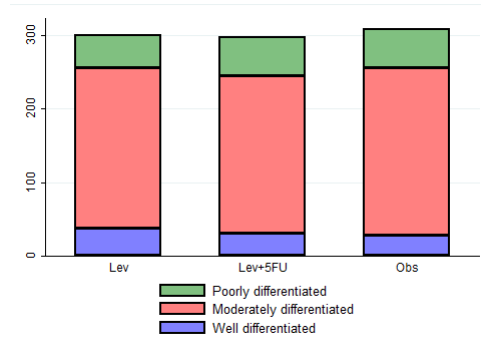
Without strata	Observed Events	Expected Events
Obs	154	141.4
Lev	165	145.3
Lev+5FU	122	154.3
	441	441.0

SISCR 2017: SA in Clinical Trials - SMay

2 - 62

DIFFERENTIATION BY TREATMENT GROUP

- Randomization worked



SISCR 2017: SA in Clinical Trials - SMay

2 - 63

- Example with more strata

SISCR 2017: SA in Clinical Trials - SMay

2 - 64

MORE STRATA - EXAMPLE 5

- Van Belle et al (Biostatistics, 2nd Edition)
- Based on Passamani et al (1982)
- Patients with chest pain
- Studied for possible coronary artery disease
 - Definitely angina
 - Probably angina
 - Probably not angina
 - Definitely not angina
- Physician diagnosis
- Outcome: Survival

SISCR 2017: SA in Clinical Trials - SMay

2 - 65

30 STRATA

# vessels	# of prox. vessels				Left Ventricular Score
	0	1	2	3	
0	5-11				
0	12-16				
0	17-30				
1	5-11	5-11			
1	12-16	12-16			
1	17-30	17-30			
2	5-11	5-11	5-11		
2	12-16	12-16	12-16		
2	17-30	17-30	17-30		
3	5-11	5-11	5-11	5-11	
3	12-16	12-16	12-16	12-16	
3	17-30	17-30	17-30	17-30	

SISCR 2017: SA in Clinical Trials - SMay

2 - 66

30 STRATA

- $\text{Chi}^2(3) = 1.47$
- P – value = 0.69

- Comparing 4 groups across 30 strata

SUMMARY

- Two sample tests
- Different flavors (weighted) two sample tests
- K – sample test
- Trend test
- Stratified test

TO WATCH OUT FOR:

- Only ranks are used for “standard” tests
- Observations with time = 0
- Crossing survival functions
- Independent censoring
- Clinical relevance
 - Log rank test and Cox
 - A difference between 3 and 6 days is judged the same as a difference between 3 years and 6 years

- Questions ?