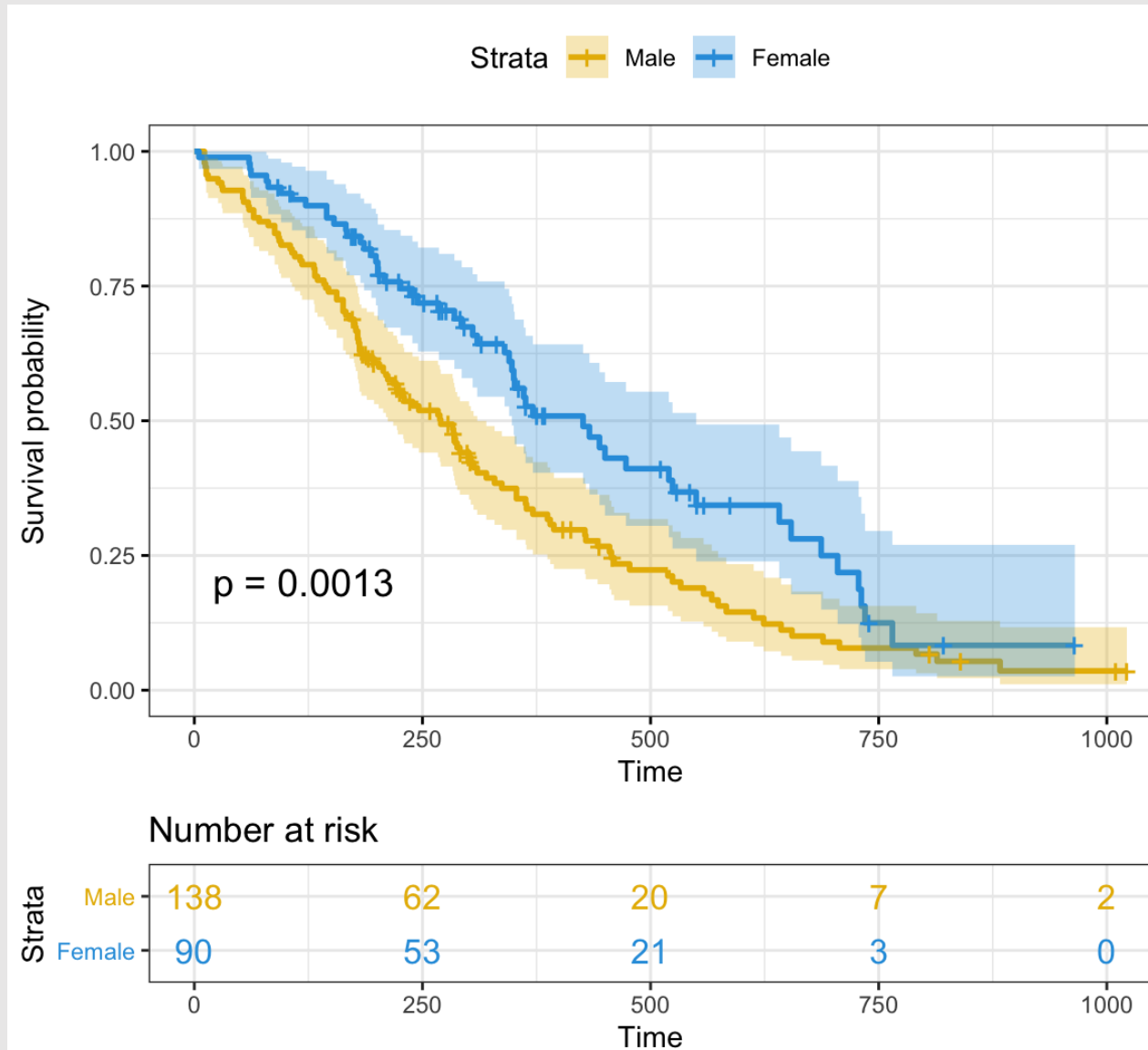


Survival Analysis II part

- Comparing survival curves
- Cox Regression



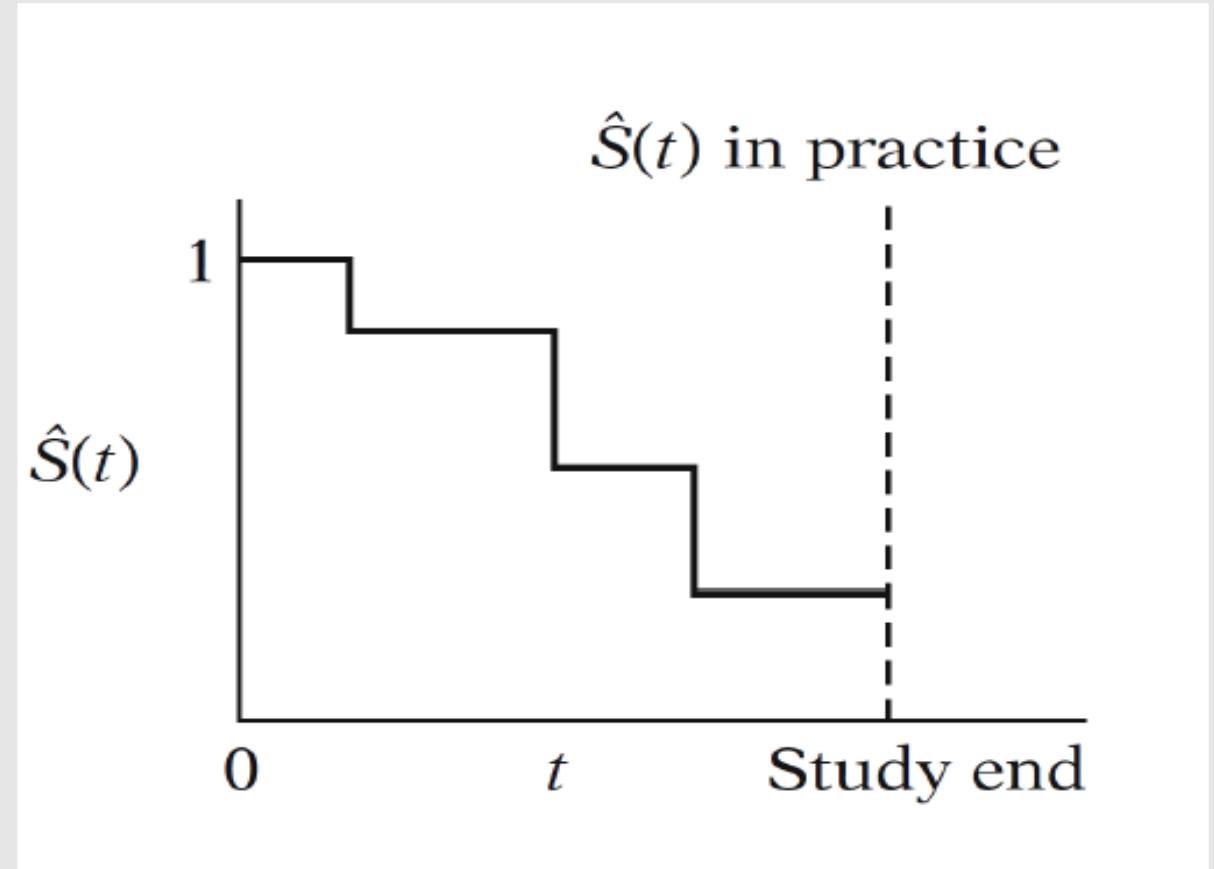
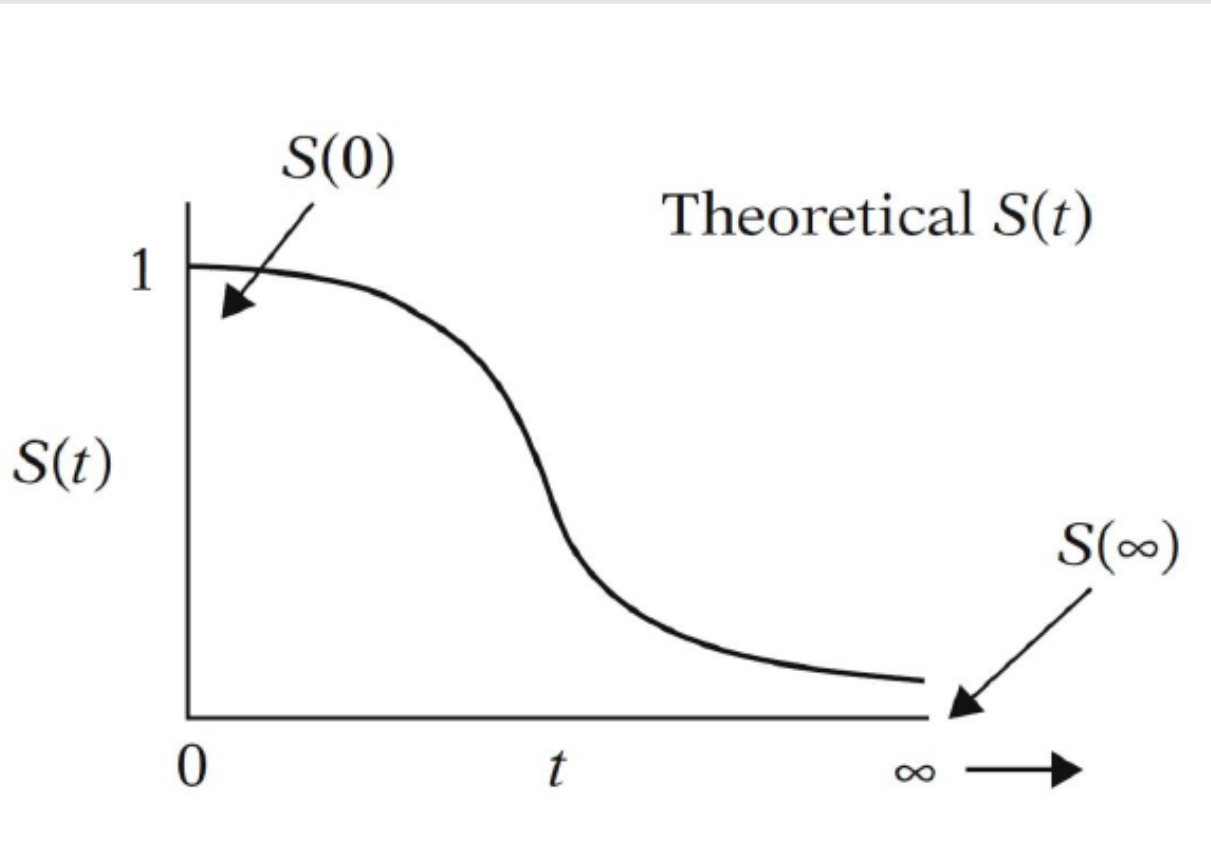


Aims of Survival Analysis

- Estimate time-to-event for a group of individuals, such as time until hospitalization or death for a group of patients.
- **To compare time-to-event between two or more groups**, such as treated vs. placebo patients in a randomized controlled trial.
- To assess the relationship of co-variables to time-to-event, such as: does weight, insulin resistance, or cholesterol influence survival time of CV patients?

Survival function:

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



Comparison of two groups of survival data

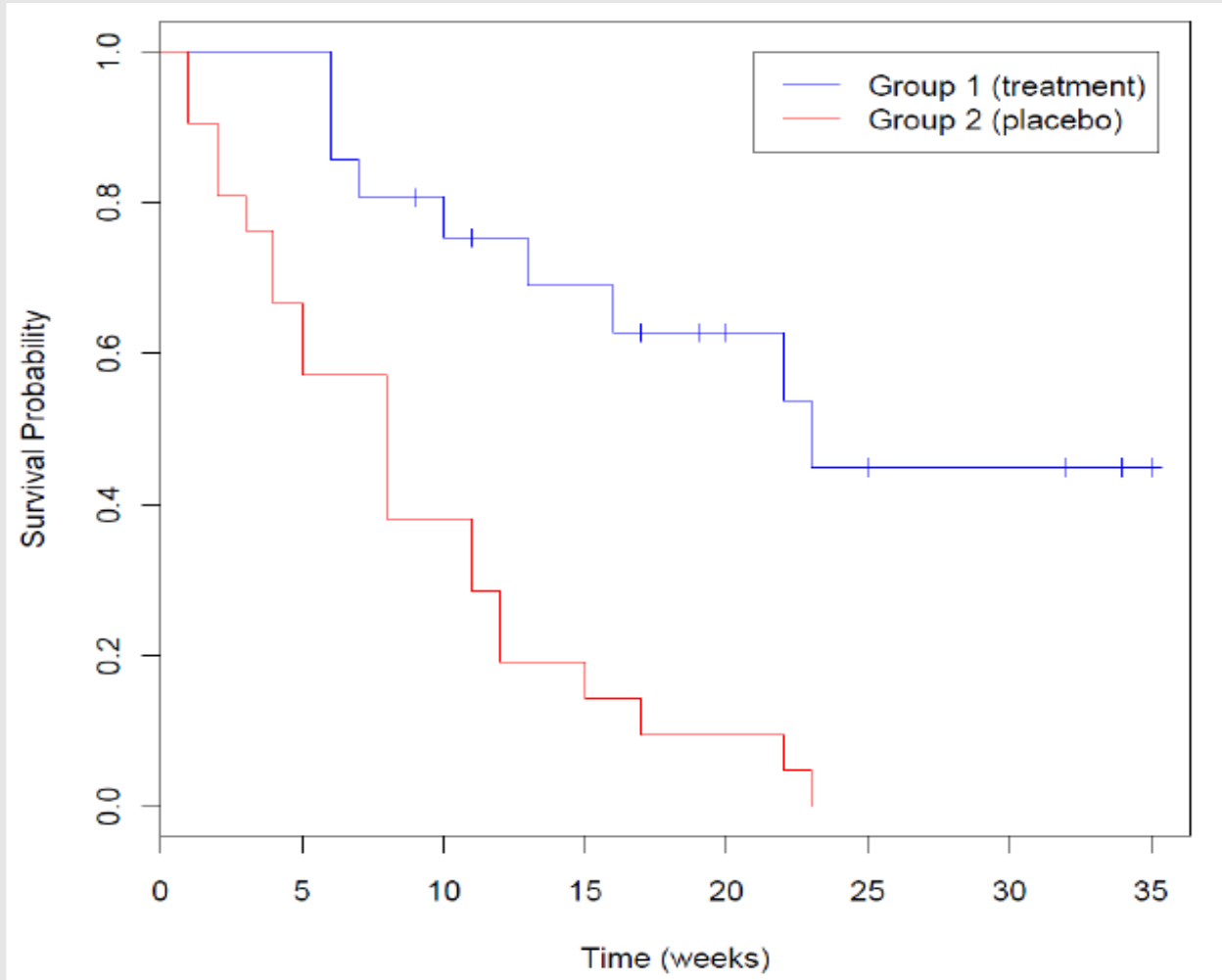
The aim is to compare survival times of two (or more) groups of patients: one **exposed** to a certain *treatment/risk factor* another **not exposed**.

We have to perform an **hypothesis test**

H_0 : There is no difference in survival among groups.

The **logrank test** is the most widely used method of comparing two or more survival curves

Comparing survival curves



Do we have any reason to claim that group 1 (treatment) has a **significant** better survival prognosis than group 2 (placebo)?

Log-rank test

We look at **2** groups [→ extensions to **several** groups are possible]

When are two KM curves *statistically* equivalent?

→ we need a **testing procedure** to compare the two curves

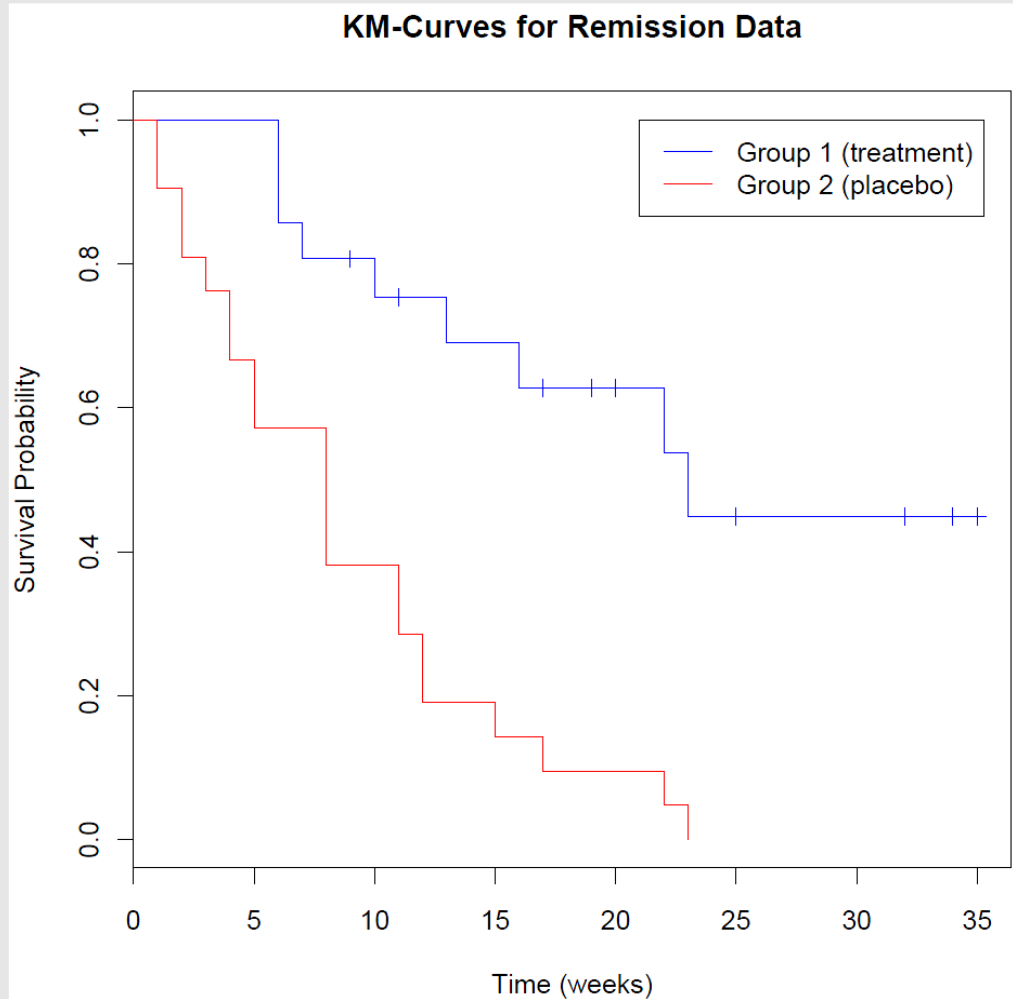
→ when we have evidence that the **true** survival curves are different?

Null hypothesis (H_0): no difference between (*true*) survival curves

Goal: To find an expression (depending on the data) from which we know the distribution (or at least approximately) **under the null hypothesis**

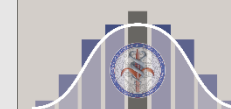
Assumption : **Proportional Hazards** over time (see later) !

Example: remission times (weeks) for two groups of leukemia patients



Remission data: n=42

$t_{(j)}$	# failures		# in risk set	
	m_{1j}	m_{2j}	n_{1j}	n_{2j}
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	12	12	6
13	1	0	12	4
15	0	1	11	4
16	1	0	11	3
17	0	1	10	3
22	1	1	7	2
23	1	1	6	1



UNITÀ DI BIOSTATISTICA

Dipartimento Universitario Clinico di Scienze Mediche Chirurgiche e della Salute

Remission data: n=42

$t_{(j)}$	# failures		# in risk set	
	m_{1j}	m_{2j}	n_{1j}	n_{2j}
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	12	12	6
13	1	0	12	4
15	0	1	11	4
16	1	0	11	3
17	0	1	10	3
22	1	1	7	2
23	1	1	6	1

Expected cell counts:

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

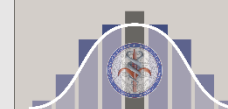
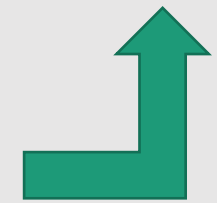
Proportion
in risk set

of failures
over both
groups

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



We expect no differences between groups under H_0



EXAMPLE

Expanded Table (Remission Data)

j	t _(j)	# failures		# in risk set		# expected		Observed-expected	
		m _{1j}	m _{2j}	n _{1j}	n _{2j}	e _{1j}	e _{2j}	m _{1j} -e _{1j}	m _{2j} -e _{2j}
1	1	0	2	21	21	(21/42) × 2	(21/42) × 2	-1.00	1.00
2	2	0	2	21	19	(21/40) × 2	(19/40) × 2	-1.05	1.05
3	3	0	1	21	17	(21/38) × 1	(17/38) × 1	-0.55	0.55
4	4	0	2	21	16	(21/37) × 2	(16/37) × 2	-1.14	1.14
5	5	0	2	21	14	(21/35) × 2	(14/35) × 2	-1.20	1.20
6	6	3	0	21	12	(21/33) × 3	(12/33) × 3	1.09	-1.09
7	7	1	0	17	12	(17/29) × 1	(12/29) × 1	0.41	-0.41
8	8	0	4	16	12	(16/28) × 4	(12/28) × 4	-2.29	2.29
9	10	1	0	15	8	(15/23) × 1	(8/23) × 1	0.35	-0.35
10	11	0	2	13	8	(13/21) × 2	(8/21) × 2	-1.24	1.24
11	12	0	2	12	6	(12/18) × 2	(6/18) × 2	-1.33	1.33
12	13	1	0	12	4	(12/16) × 1	(4/16) × 1	0.25	-0.25
13	15	0	1	11	4	(11/15) × 1	(4/15) × 1	-0.73	0.73
14	16	1	0	11	3	(11/14) × 1	(3/14) × 1	0.21	-0.21
15	17	0	1	10	3	(10/13) × 1	(3/13) × 1	-0.77	0.77
16	22	1	1	7	2	(7/9) × 2	(2/9) × 2	-0.56	0.56
17	23	1	1	6	1	(6/7) × 2	(1/7) × 2	-0.71	0.71
Totals		9	21			19.26	10.74	-10.26	+10.26

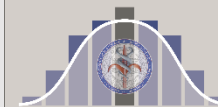
$$O_i - E_i = \sum_{j=1}^{\# \text{ failure times}} (m_{ij} - e_{ij})$$

$$\text{Log-rank statistic} = \frac{(O_2 - E_2)^2}{\text{Var}(O_2 - E_2)}$$

Remark: We could also work with Group 1 and we would get the same statistic

$$O_1 - E_1 = -10.26$$

$$O_2 - E_2 = 10.26$$



$$\text{Log-rank statistic for two groups} = \frac{(O_2 - E_2)^2}{\text{Var}(O_2 - E_2)} \sim \chi_1^2$$

Call:

```
survdiff(formula = Surv(time, status) ~ treatment)
```

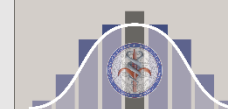
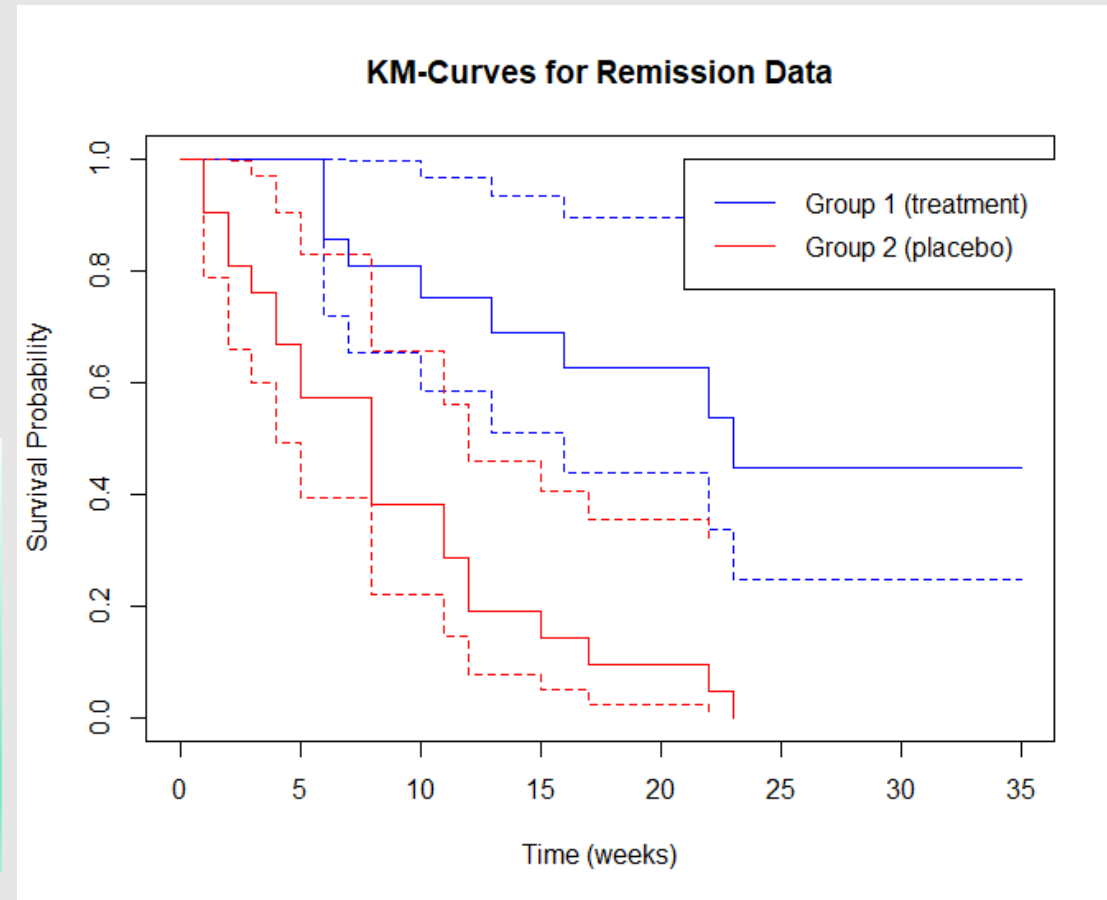
	N	Observed	Expected	(O-E)^2/E	(O-E)^2/V
treatment=1	21	9	19.3	5.46	16.8
treatment=2	21	21	10.7	9.77	16.8

Chisq= 16.8 on 1 degrees of freedom, p= 4e-05



What does this tell us?

probability of obtaining a test statistic **at least as extreme** as the one that was actually observed, under H0.



LG test for Several Groups

H_0 : **All** survival curves are the same

- ▶ Suppose we have $K > 2$ groups and we wish to simultaneously compare them with respect to survival time distributions (or equivalently, hazards)

$$H_0: \lambda_1(t) = \lambda_2(t) = \dots = \lambda_K(t), \text{ for all } t > 0$$

(i.e. the survival curves for the all groups are equal everywhere)

- ▶ We are particularly concerned with the alternatives

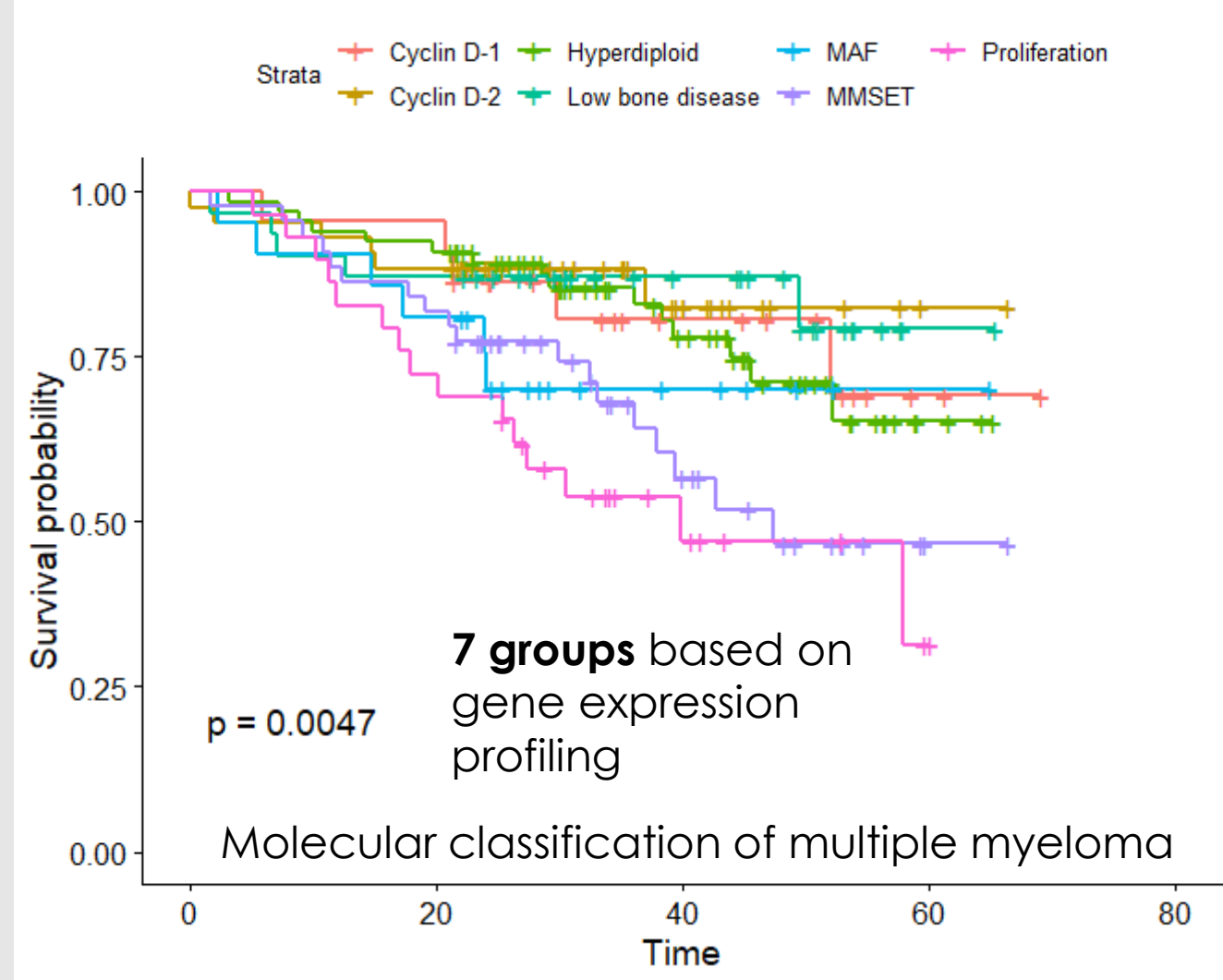
$$H_A: \lambda_k(t) > \lambda_{k'}(t), \text{ for some } t > 0$$

or

$$\lambda_k(t) < \lambda_{k'}(t), \text{ for some } t > 0$$

for at least some $k \neq k'$

- Log-rank statistic for > 2 groups involves computing variances and covariances of $O_i - E_i$
- $G (\geq 2)$ groups: log-rank statistic $\sim \chi^2$ with $G-1$ df



Pairwise comparisons between group levels with corrections for **multiple testing issue [alpha inflation...]**

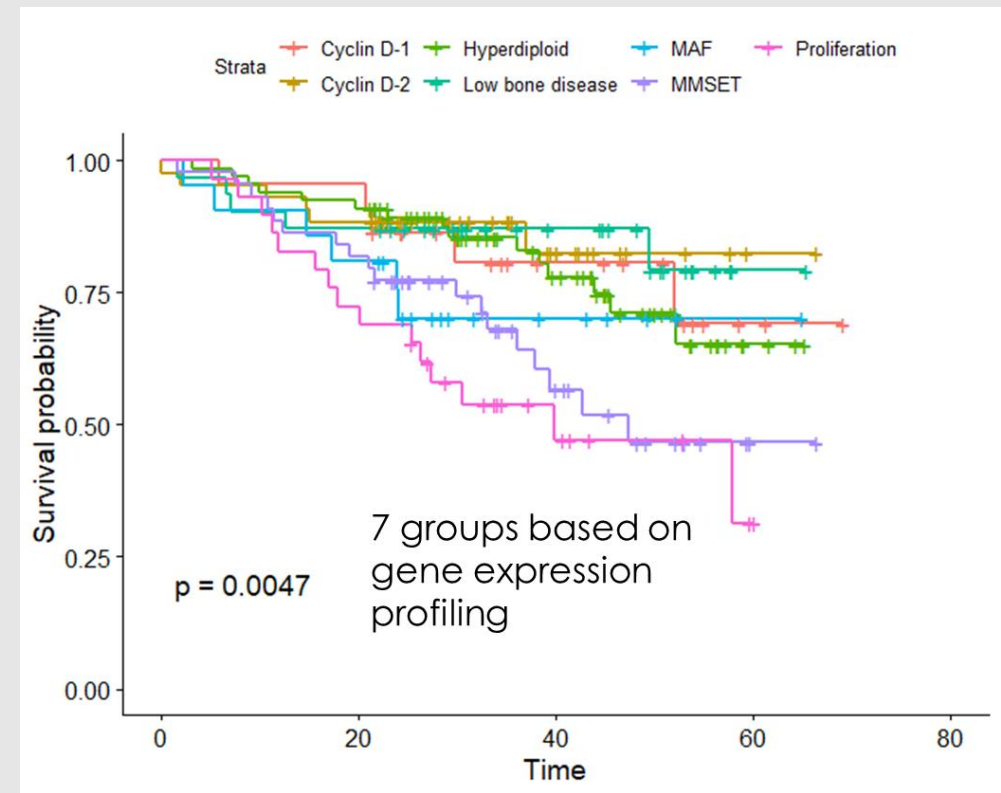


R function: `pairwise_survdiff {survminer}`

	Cyclin D-1	Cyclin D-2	Hyperdiploid	Low bone disease	MAF	MMSET
Cyclin D-2	0.723	-	-	-	-	-
Hyperdiploid	0.943	0.723	-	-	-	-
Low bone disease	0.723	0.988	0.644	-	-	-
MAF	0.644	0.447	0.523	0.485	-	-
MMSET	0.328	0.103	0.103	0.103	0.723	-
Proliferation	0.103	0.038	0.038	0.062	0.485	0.527

p value adjustment method: BH

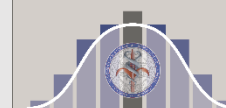
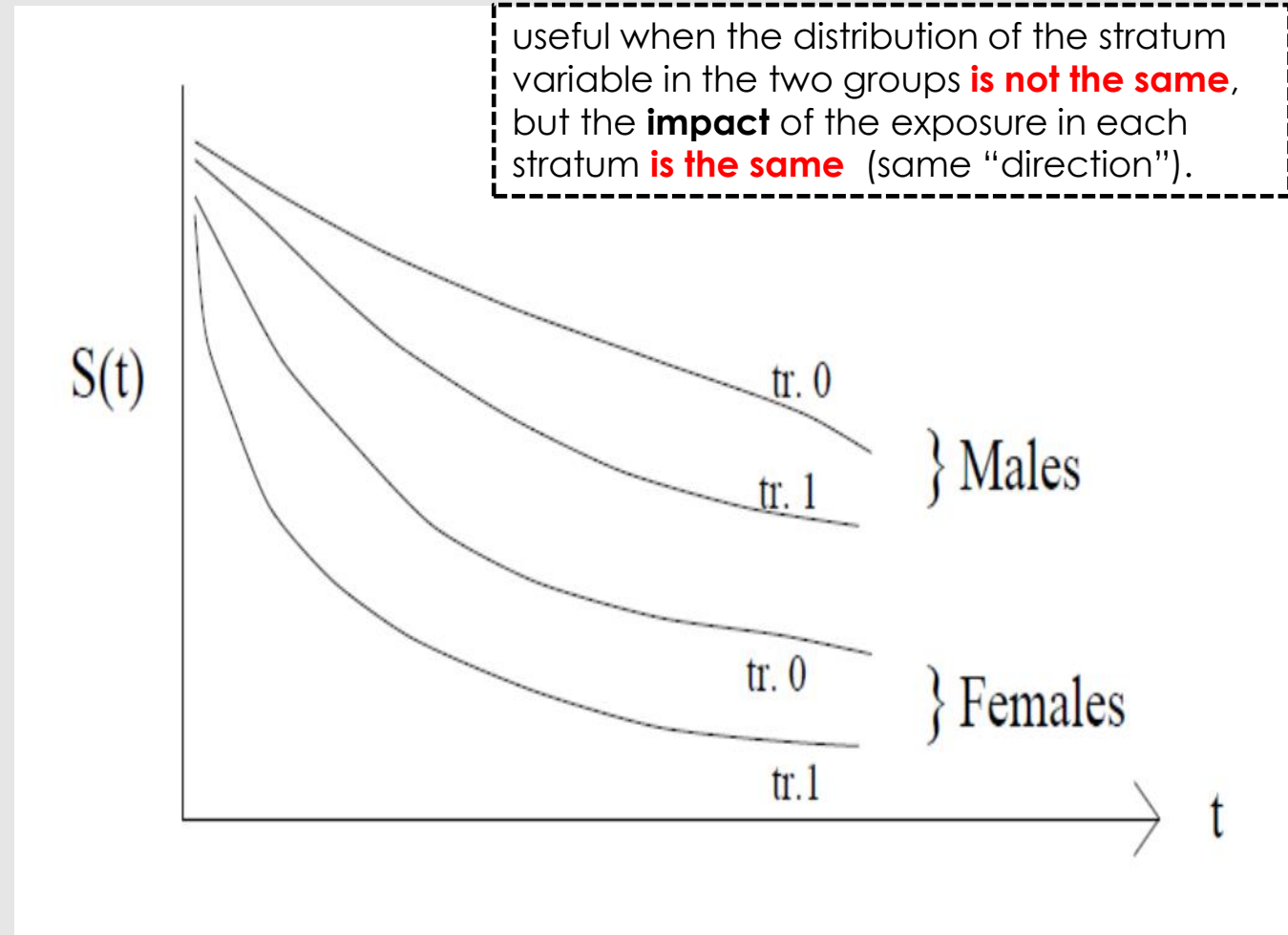
→ Various choices for the adjustment method



Stratified Log-rank test

Variation of log rank test:

- Allows controlling for *additional* (stratified:categorical) variable [**confounder**]
- Split data into strata, based on values of confounder
- Calculate $O-E$ **within** strata
- **Sum** $O-E$ across strata



Stratified log rank test – Example:

- Remission data
- Stratified variable: 3-level variable (LWBC3) indicating low, medium, or high log white blood cell count (coded 1, 2, and 3, respectively)

Treated Group: rx=0 Placebo Group: rx=1

```
->lwbc3 = 1
```

rx	Events observed	Events expected
0	0	2.91
1	4	1.09
Total	4	4.00

```
->lwbc3 = 2
```

rx	Events observed	Events expected
0	5	7.36
1	5	2.64
Total	10	10.00

```
->lwbc3 = 3
```

rx	Events observed	Events expected
0	4	6.11
1	12	9.89
Total	16	16.00

Recap: Non-stratified test : χ^2 -value of 16.79
and corresponding p-value rounded to 0.0000

Call:

```
survdif(formula = Surv(time, status) ~ treatment)
```

	N	Observed	Expected	(O-E)^2/E	(O-E)^2/V
treatment=1	21	9	19.3	5.46	16.8
treatment=2	21	21	10.7	9.77	16.8

Chisq= 16.8 on 1 degrees of freedom, p= 4e-05

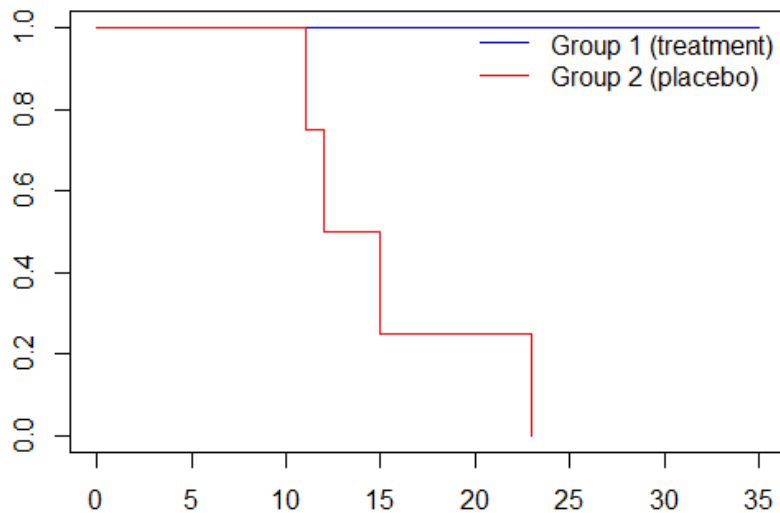
```
fit <- survdiff(Surv(data$V1, data$V2) ~ data$V5 + strata(lwbc3))
fit
Call:
survdiff(formula = Surv(data$V1, data$V2) ~ data$V5 + strata(lwbc3))
```

	N	Observed	Expected	(O-E) ^2/E	(O-E) ^2/V
data\$V5=0	21	9	16.4	3.33	10.1
data\$V5=1	21	21	13.6	4.00	10.1

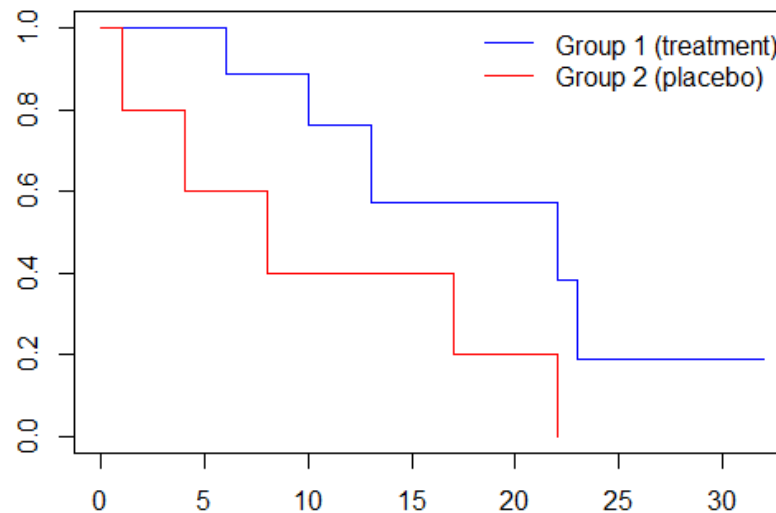
Chisq= 10.1 on 1 degrees of freedom, p= 0.001

Always significant, same direction of the effect, but **magnitude** of the effect varies across strata (varying sample size..)

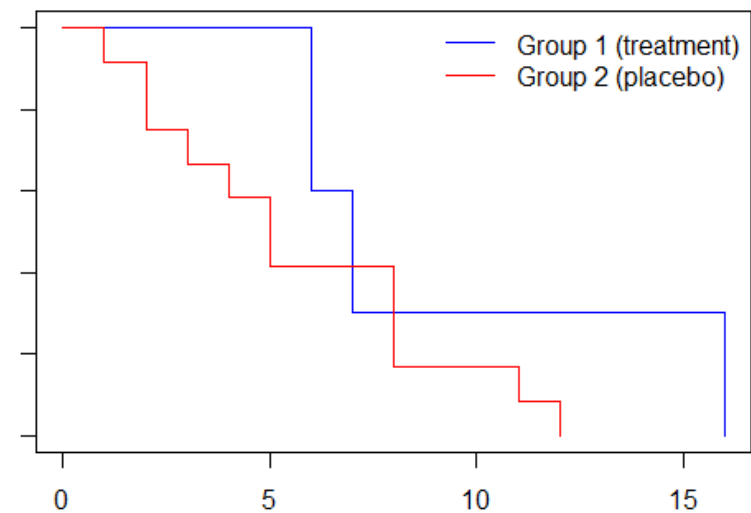
lwbc=1



lwbc=2



lwbc=3



Stratified vs. unstratified approach

Log rank unstratified*

$$O_i - E_i = \sum_j (m_{ij} - e_{ij})$$

i = group #, j = jth failure time

Log rank stratified*

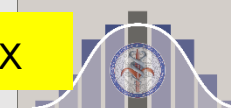
$$O_i - E_i = \sum_s \sum_j (m_{ijs} - e_{ijs})$$

i = group #, j = jth failure time,
s = stratum #

Limitations:

- Sample size may be **small** within strata
- **Categorical** stratifying variable and exposure
- **Interactions** ?

*At the denominator there is always an estimate of the variance-covariance matrix

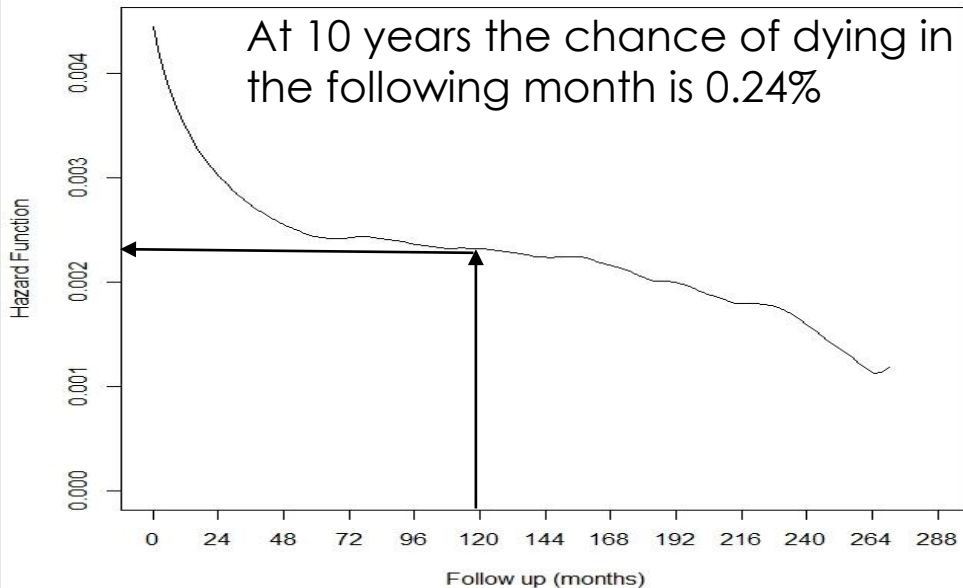


Again: Hazard Function & Cumulative Hazard Rate

The probability that **if you survive to t** , you will succumb to the event in the next instant.

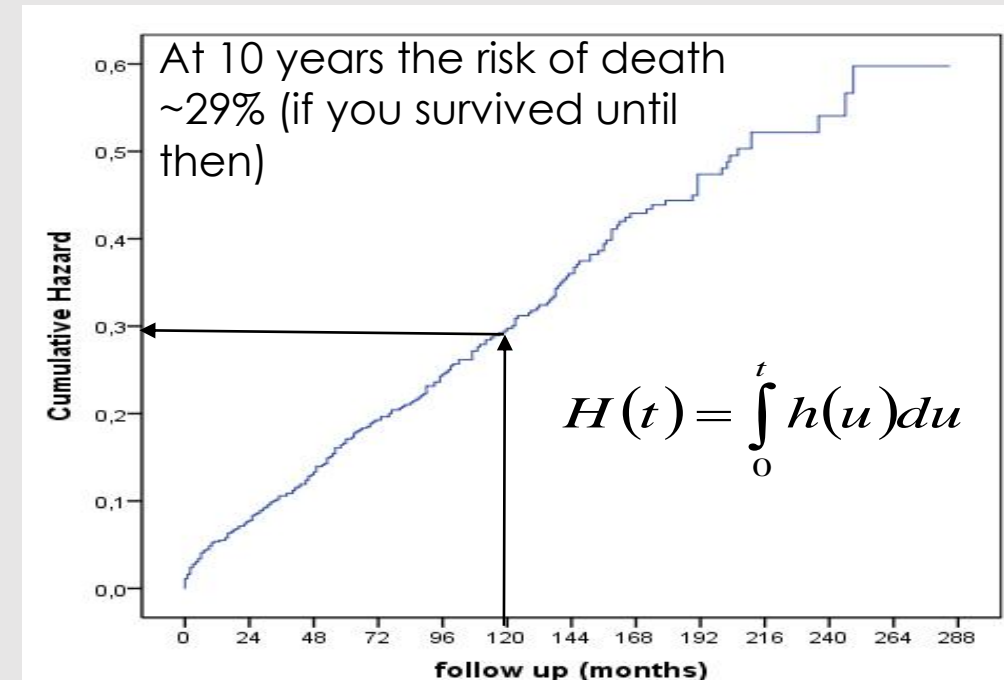
$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t | T \geq t)}{\Delta t}$$

instantaneous event rate



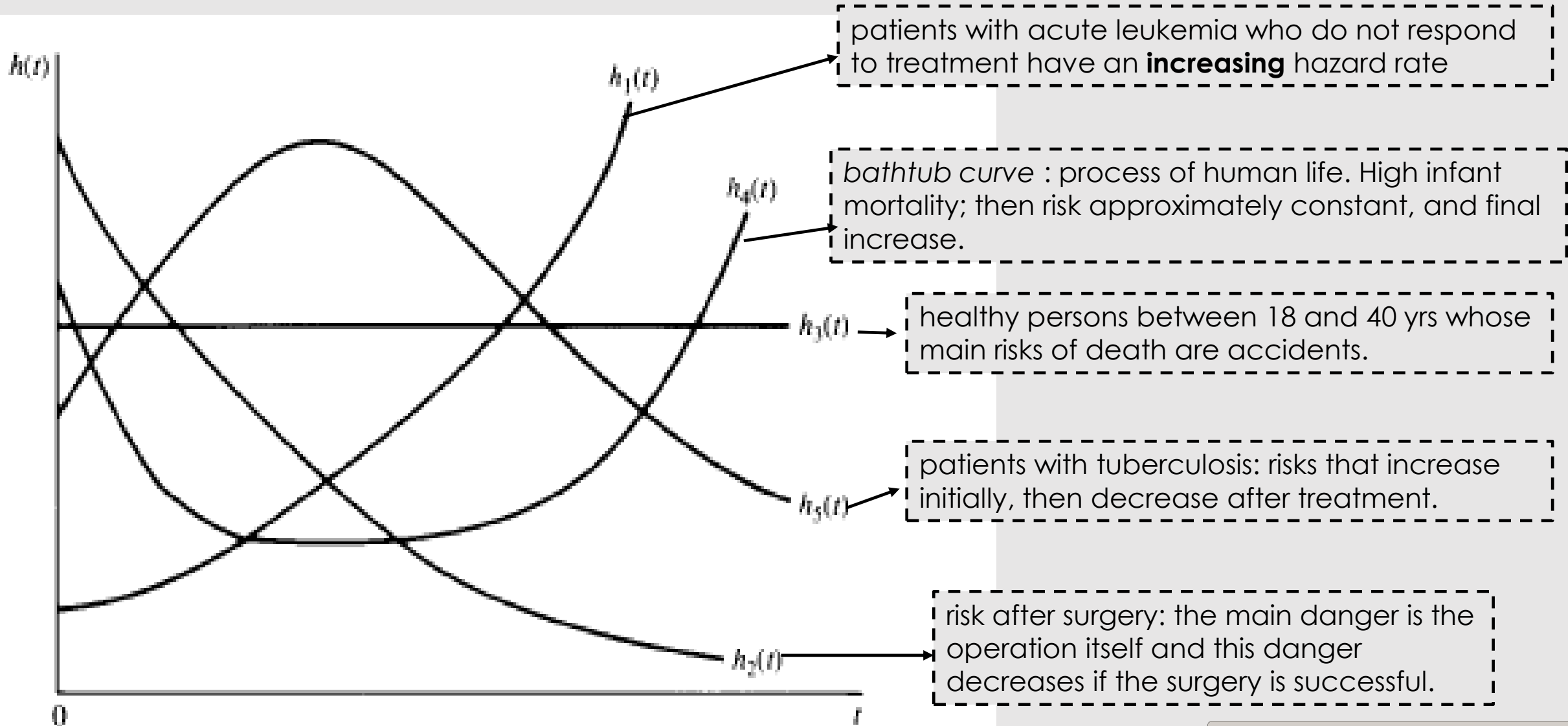
Risk of event **up to time t** given that the event has not occurred before t

Cumulative Hazard Function
(Cumulative Hazard Rate)

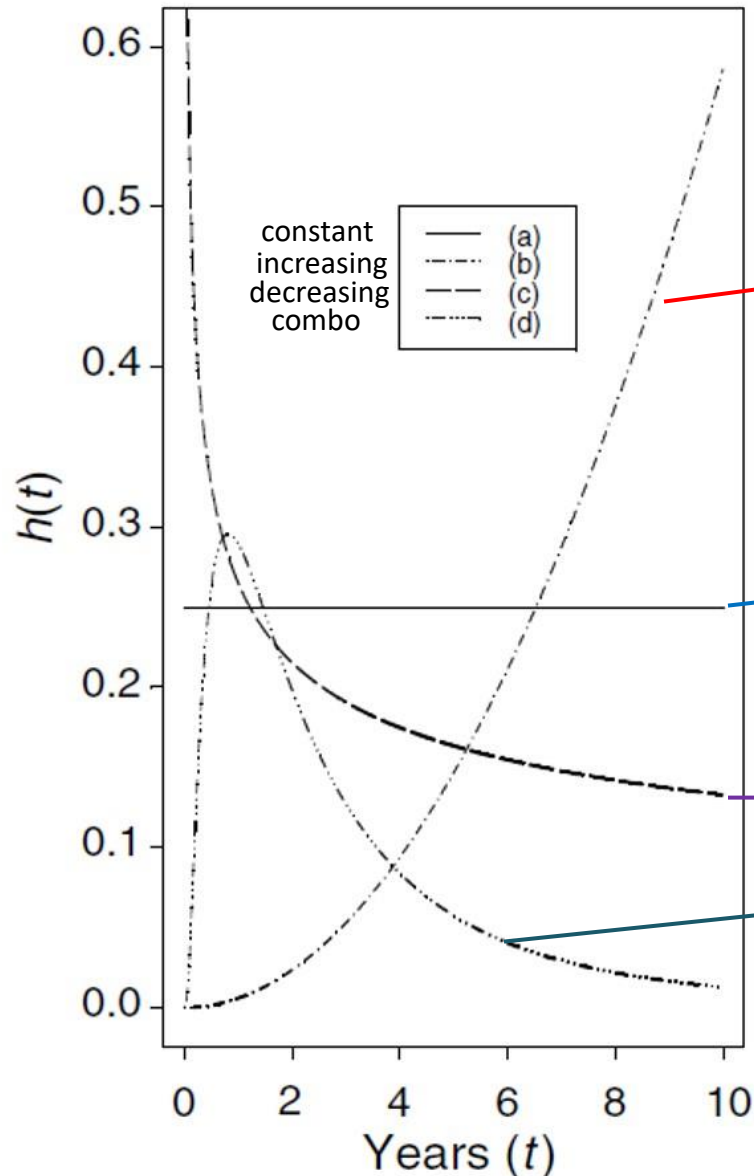


Hazard function: “force of mortality over time”

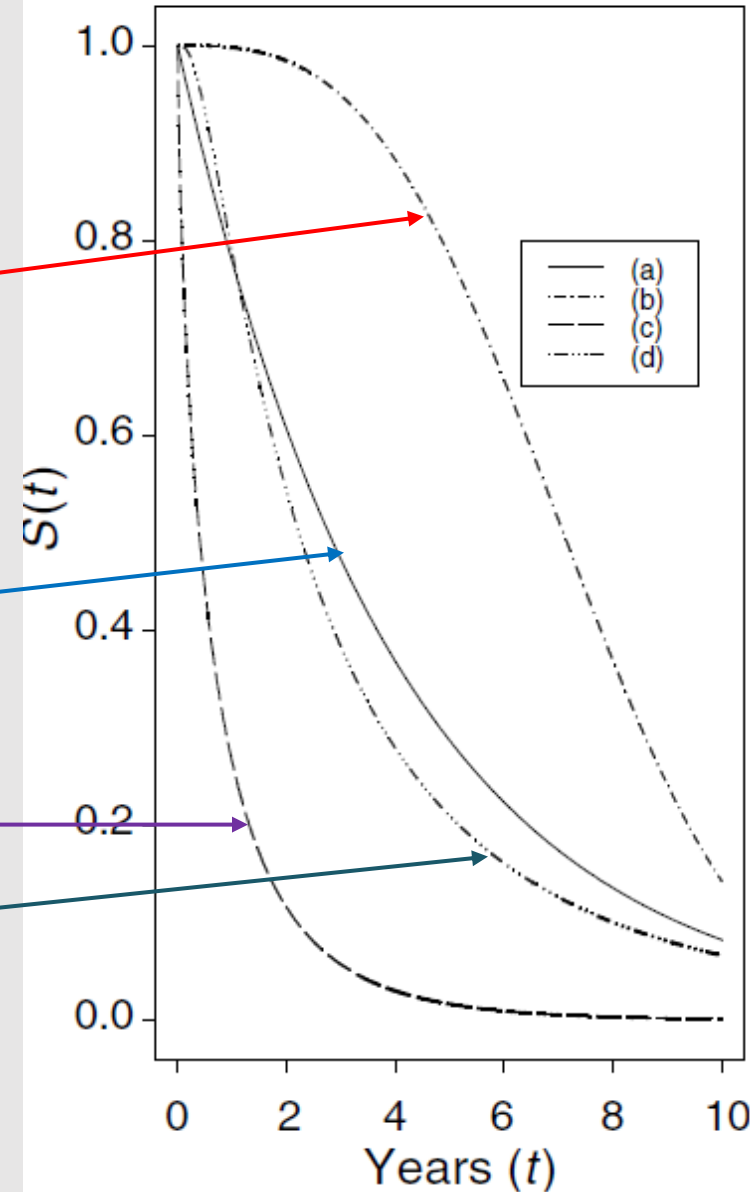
Hazard function



Hazard curves



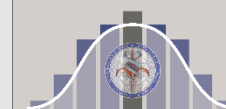
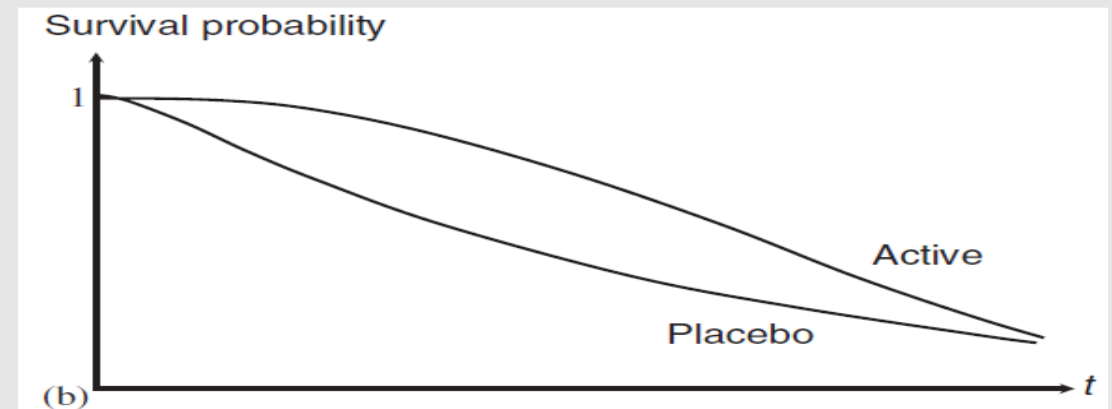
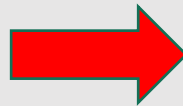
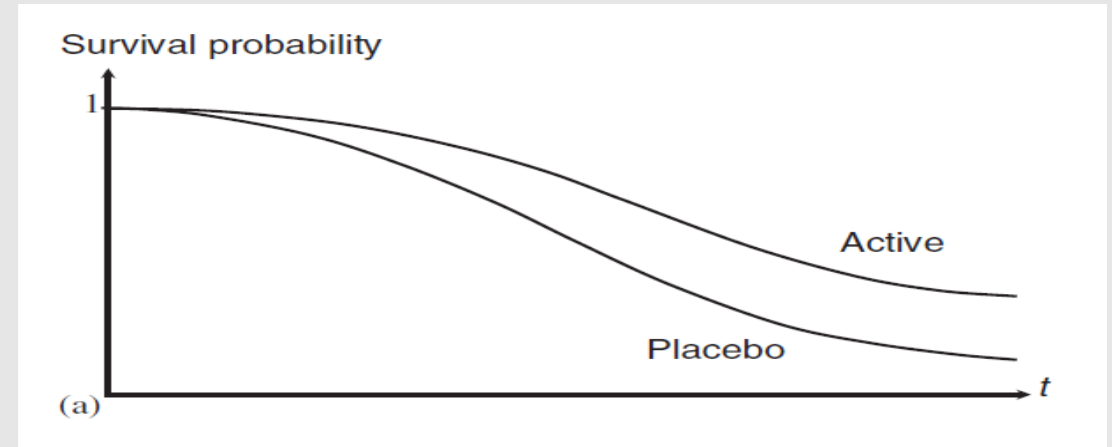
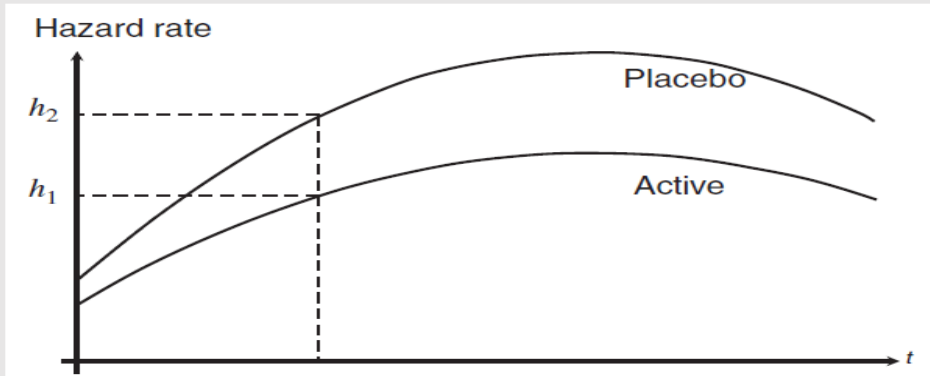
Survival curves



differences between survival curves **are much less evident**, compared to differences in the corresponding hazard functions

Proportional hazards (PH) assumption

The hazard **at any given time** for an individual in one group **is proportional to** the hazard **at any given time** for an individual in the other group. If the hazard functions are proportional \rightarrow survival functions **do not cross** one another. **[log-rank test assumes PH !!]**. There is a test that could be used to verify PH.



A common problem in medical applications is how to check for the **overall homogeneity** of survival curves **when the PH assumption does not hold**.

A treatment may offer a **short-term** benefit but does not provide **long-term** advantages.

Or two survival curves **cross** each other.

Using the log-rank test under conditions of **non-proportional hazards** may lead to misleading results.

1. **Weighted** log-rank tests:

- Wilcoxon
- Tarone-Ware
- Peto
- Fleming-Harrington

Different **weights** are applied to differences between expected and observed deaths to **emphasize** certain times more than others...

2. The **supremum (Renyi)** tests designed to detect differences in survival curves which **cross**. That is, an early difference in survival in favor of one group is balanced by a later reversal.

So far, we have discussed methods to compare **two** or **more** survival curves, both under the PH assumption.

These methods are useful when the exposure/risk factor is **categorical** and there is at maximum another confounder that is also categorical (with the same «*direction of effect*» of the exposure across **strata...**).

Next step will be to introduce a **more general regression model** (on the scale of the hazard function) that will allow us to estimate the **joint effect** of one or more risk factors (in whatever scale of measure) on the time-to-event or to evaluate the **specific effect** of one exposure of interest adjusting for multiple confounders.



Aims of Survival Analysis

- Estimate time-to-event for a group of individuals, such as time until hospitalization or death for a group of patients.
- To compare time-to-event between two or more groups, such as treated vs. placebo patients in a randomized controlled trial.
- **To assess the relationship of co-variables to time-to-event**, such as: does weight, insulin resistance, or cholesterol influence survival time of CV patients?

The Most-Cited Statistical Papers

Journal of Applied Statistics
Vol. 32, No. 5, 461–474, July 2005

(1) With 25,869 citations (currently cited 1,984 times per year),

Kaplan, E. L. & Meier, P. (1958) Nonparametric estimation from incomplete observations, *Journal of the American Statistical Association*, 53, pp. 457–481.

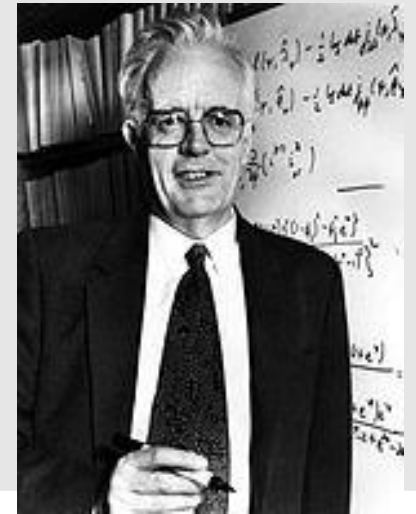
(among the **top five most cited** papers for the entire field of science)

(2) With 18,193 citations (1,342 per year),

Cox, D. R. (1972) Regression models and life tables, *Journal of the Royal Statistical Society, Series B*, 34, pp. 187–220.

Semi-parametric regression approach that estimates the effect of **covariates** on the **hazard function**

David Cox



Regression Models and Life-Tables

BY D. R. COX

Imperial College, London

[Read before the ROYAL STATISTICAL SOCIETY, at a meeting organized by the Research Section, on Wednesday, March 8th, 1972, Mr M. J. R. HEALY in the Chair]

SUMMARY

The analysis of censored failure times is considered. It is assumed that on each individual are available values of one or more explanatory variables. The hazard function (age-specific failure rate) is taken to be a function of the explanatory variables and unknown regression coefficients multiplied by an arbitrary and unknown function of time. A conditional likelihood is obtained, leading to inferences about the unknown regression coefficients. Some generalizations are outlined.



Why don't we use **others** regression methods ?

Logistic regression [binary outcome]:

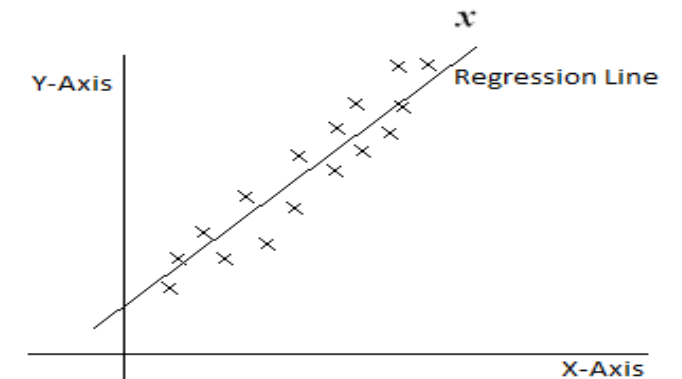
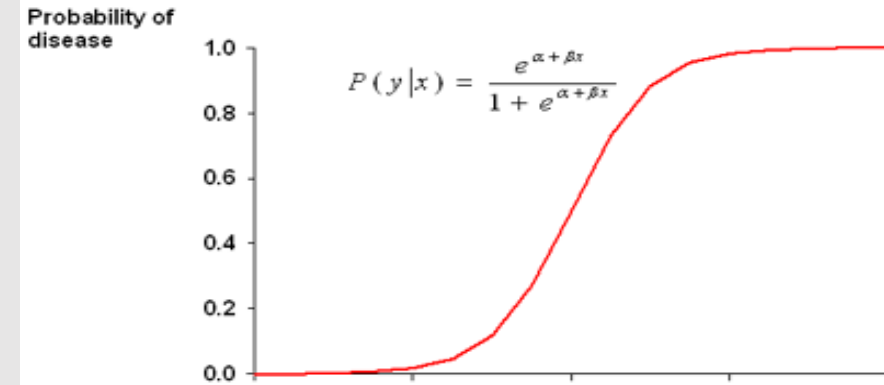
- ignores information about the **time** to the event

Linear regression [continuous outcome]:

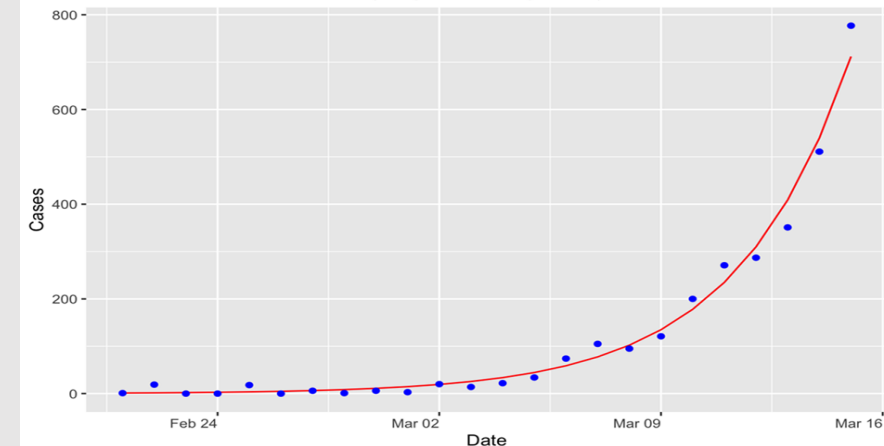
- not suitable for **non-symmetric** [>0] distributions [like follow up times]
- does not take into account **censoring**

Poisson regression [event counts/rates]:

- #events/RR **in a given interval** (**\neq time to the event**)



Predicted vs. Actual Number of COVID-19 Cases
(using Poisson Regression)

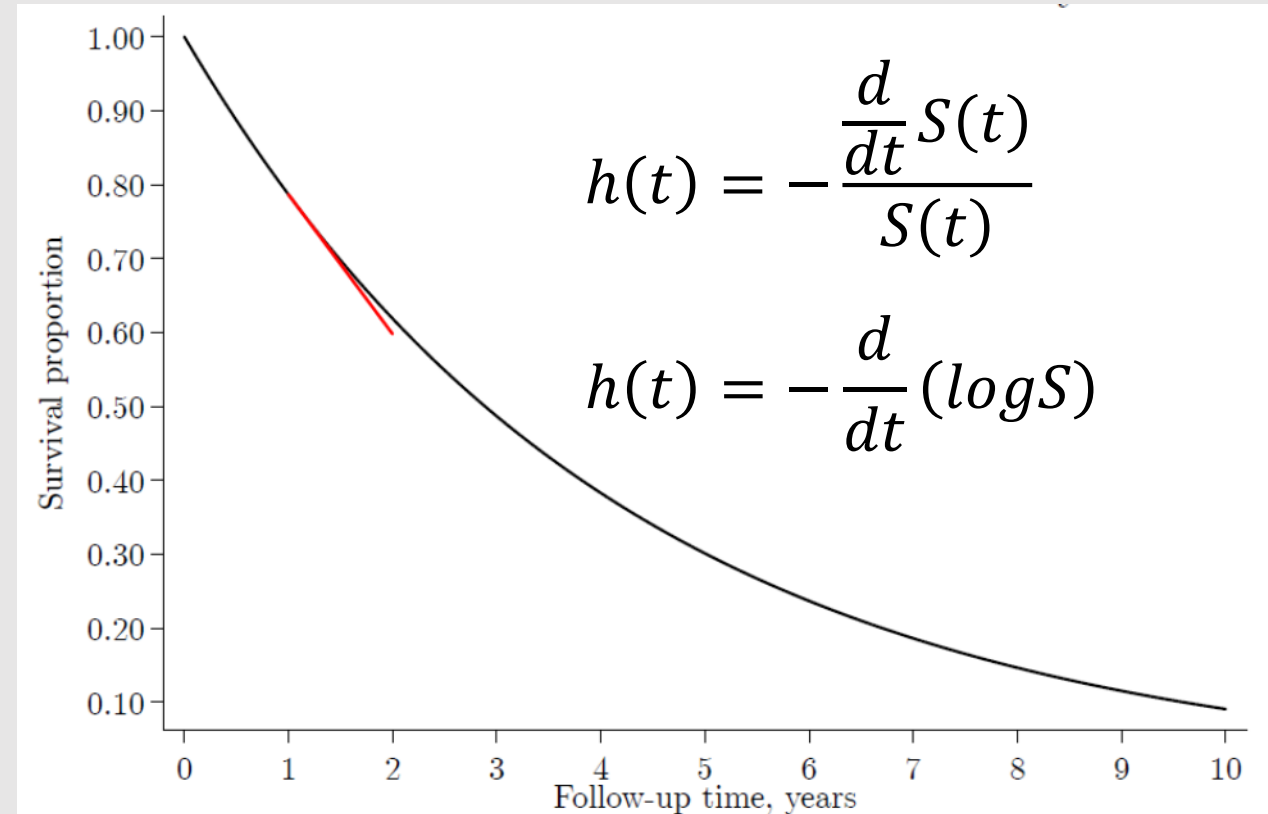
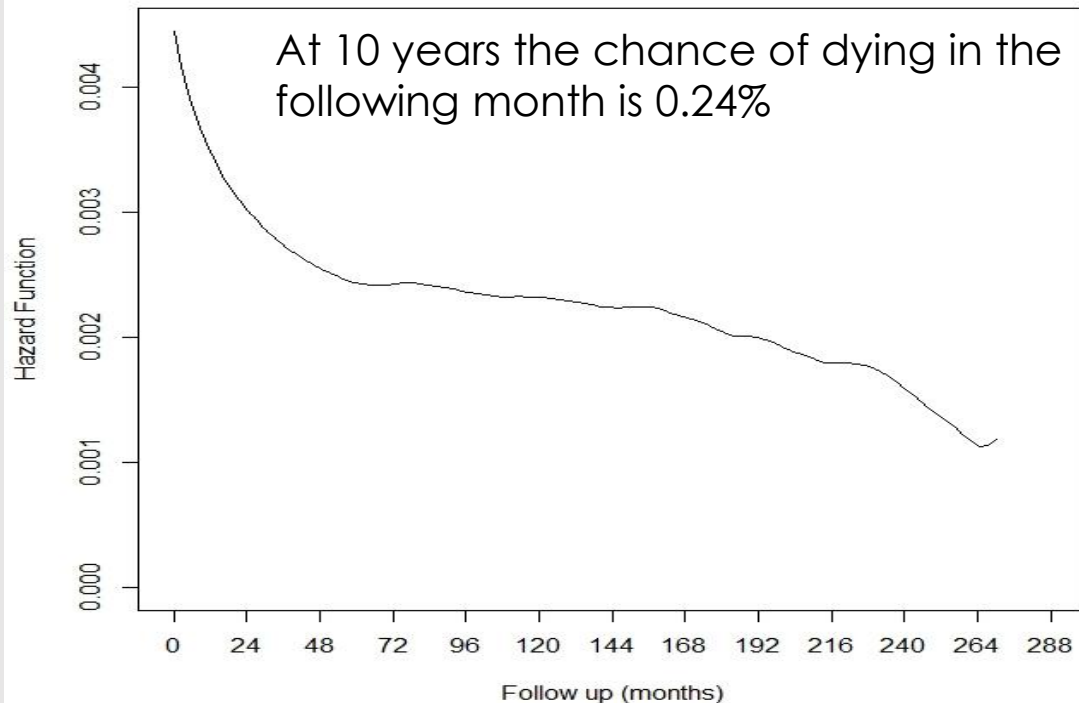


The *dependent* variable of the Cox model

The probability that **if you survive to t** , you will succumb to the event in the next instant.

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t | T \geq t)}{\Delta t}$$

instantaneous event rate



Cox Regression Model

The scale on which *linearity* is assumed is the **log-hazard** scale:

$$h(t|X) = h_0(t) \exp(X_1\beta_1 + X_2\beta_2 + X_3\beta_3 + \dots + X_p\beta_p)$$

$$\log \left(\frac{h(t|X)}{h_0(t)} \right) = X_1\beta_1 + X_2\beta_2 + X_3\beta_3 + \dots + X_p\beta_p$$

- $h_0(t)$ is the **baseline hazard function**
- the exponential function represents the effect of the linear combination of the covariates X on the hazard

The aim is to determine the **joint** effect of the covariates on the hazard or to focus on a **specific** effect.

The dependent variable of the Cox model is the hazard function.

The model assumes that the risk at time t for subject i is: $h_i(t|X_i) = h_0(t)exp(X_i\beta)$

Remind of **censored** data:

someone who is followed for 18 months is a part of the computations *until the interval that contains the censoring time* (risk set) and not thereafter (**partial likelihood**).

Why $exp(\text{linear predictor})$? To avoid **negative** hazard rates.

- Implies that factors are multiplicative, e.g., treatment reduces the hazard by X %.
- Two covariates multiply in effect
- For *biological phenomena* it seems to fit well

The baseline hazard in Cox is estimated “**non-parametrically**”:

- estimated on the specific dataset
- does not extrapolate....

- To estimate β Cox proposed a **partial** likelihood (PL) procedure based on conditional probability:

$$L(\beta) = \prod_{j=1}^n \frac{\exp(X(t_{(j)})\beta)}{\sum_{i \in R_j} \exp(X_i(t_{(j)})\beta)}$$

- Maximizing the PL function we obtain:
 1. Estimates of β
 2. Standard errors for β
 3. p values for β

$t_{(1)}, \dots, t_{(n)}$ ordered event times

R_j **Risk set** at time $t_{(j)}$

$X(t_{(j)})$ covariates for the individual who fails at time $t_{(j)}$

(non-parametric estimate of cumulative baseline hazard could be obtained after β estimation)

In the Cox model the statistical independence between censoring and survival time is obtained conditional to the covariates

Interpretation of parameter estimates

Let us consider two subjects i e i' :

$$\begin{aligned}\eta_i &= \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} \\ \eta_{i'} &= \beta_1 x_{i'1} + \beta_2 x_{i'2} + \dots + \beta_p x_{i'p}\end{aligned}$$

linear part of the Cox model

The **hazard ratio** between them is:

$$\frac{h_i(t|X)}{h_{i'}(t|X)} = \frac{h_0(t)\exp(\eta_i)}{h_0(t)\exp(\eta_{i'})} = \frac{\exp(\eta_i)}{\exp(\eta_{i'})}$$

Suppose to have a single continuous variable X : $h_i(t) = h(t)\exp(\beta x_i)$

The ratio of the hazard for a subject with value $x+1$ with respect to one with value x is:

$$\frac{\exp\{\beta(x+1)\}}{\exp(\beta x)} = \exp(\beta)$$

'HAZARD RATIO'

$$\text{Exp}(\beta) = \text{HAZARD RATIO (HR)}$$

If **HR** ~ 1 (95% CI contains 1) : there is **not a significant** impact of the covariate X on the hazard of event

If **HR** > 1 (95% CI > 1) : presence or increasing values of X **increase** the hazard of event (=decrease survival)

If **HR** < 1 (95% CI < 1) : presence or increasing values of X **decrease** the hazard of event (=increase survival)

Impact of gender (M=0,F=1) and level of education (school yrs) with respect to time to the first marriage:

Cox model results	β	se(β)	exp(β) HR	lower 95% CI	upper 95% CI
Gender (F vs M)	0,48	0,20	1,61	1,09	2,40
School years	-0,07	0,02	0,93	0,51	0,98

At a given instant in time, the hazard of marriage for women is **1.61 times higher** than men (at the same level of education)

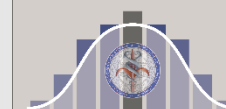
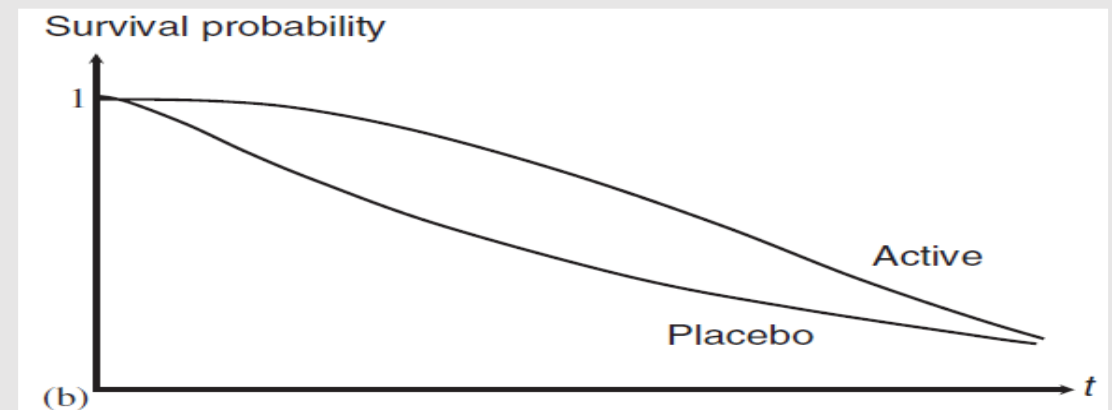
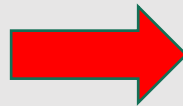
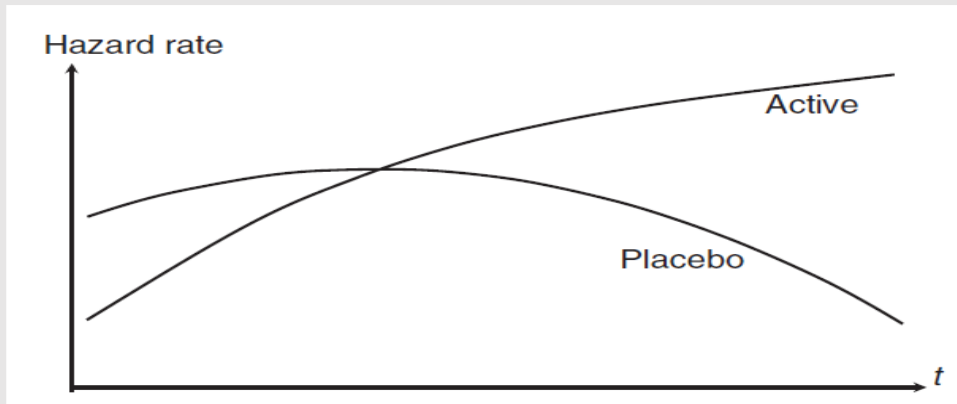
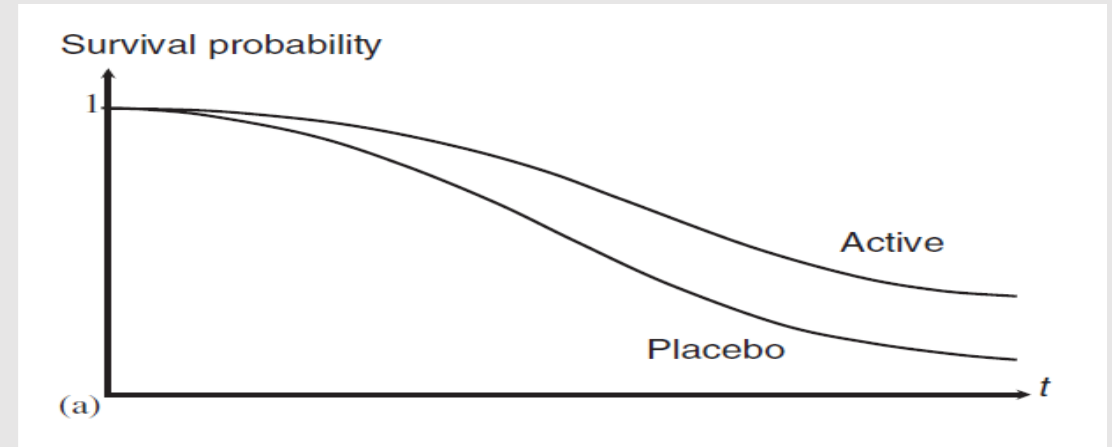
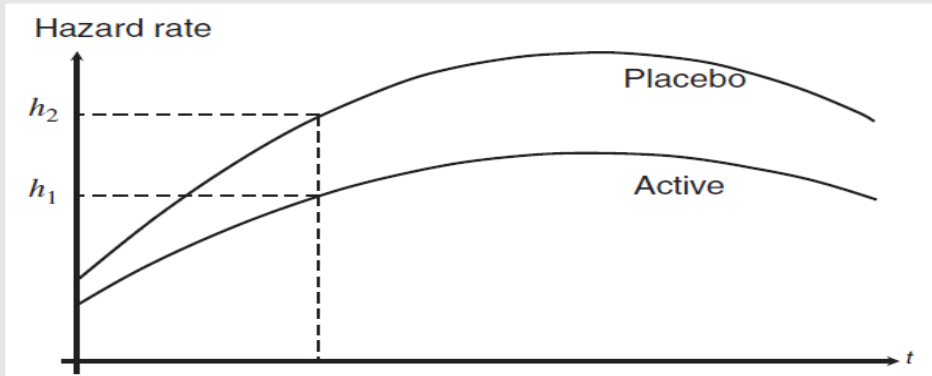
At a given instant in time, the hazard of marriage for women is **61% higher** than for men (at the same level of education)

For people (men or women) with an additional +1yr school the hazard of marriage, *at a given instant in time*, is **0.93 times** than for those without...

For each extra yr of school the hazard of marriage (men or women) *at a given instant in time* is **7% less**.

Again : proportional hazards (PH) !!

The hazard **at any given time** for an individual in one group is proportional to the hazard **at any given time** for an individual in the other group. If the hazard functions are proportional \rightarrow survival functions **do not cross** one another...



Cox model assumes **proportional hazards** (PH). Covariates X have always *the same relative effect* along time:

$$h(t|X) = h_0(t) \exp(X_1\beta_1 + X_2\beta_2 + \dots + X_p\beta_p) = h_0(t) \exp(\mathbf{X}\boldsymbol{\beta})$$

The function $\exp(\mathbf{X}\boldsymbol{\beta})$ does not depend on t

Hazard Ratio between two subjects, with covariates X and X^* does not depend on t :

$$\frac{h_0(t) \exp(\mathbf{X}\boldsymbol{\beta})}{h_0(t) \exp(\mathbf{X}^*\boldsymbol{\beta})} = \exp((\mathbf{X} - \mathbf{X}^*)\boldsymbol{\beta})$$

If PH assumption does not hold, «standard» Cox model could be no longer valid
[we could check for this] [there are extensions]

$$h_i(t|X_i) = h_0(t)\exp(X_i\beta)$$

- β_k is the **difference** in the log-hazard function comparing two subpopulations differing in x_k by “1-unit” and that are similar with respect to all other covariates in the model
- the **effect** expressed by β_k is **adjusted for** all other covariates in the model, so it has the interpretation of a log-relative hazard associated with a change in x_k , holding other covariates constant **at some fixed value**
- is it possible to compare **hypothetical patients** with different covariates values and check how their **estimated survival curves** appear; [remind: the baseline hazard depends on the study cohort...]
- the Cox PH model is indeed a model for the hazard more than a model for survival time, **although they are related one-to-one if no competing risks exists**

Survival function derived from the Cox regression model (no competing risks !!!)

Once the β are estimated, we can obtain the corresponding survival function:

$$S(t|x) = S_0(t)\exp(\beta x)$$

$S_0(t)$ is derived from an estimate of the **cumulative baseline hazard** (a complex derivation in the non-parametric form, similar to the Nelson-Aalen formulation)

The estimate of $S_0(t)$ and a fixed set of values for the explanatory variables produce an estimate of the survival function **for a specific person or group**.

The expression for $S(t|x)$ shows that proportional hazard functions dictate that the estimated survival functions **do not intersect**.

Summary: basic assumptions (all standard methods, KM, log rank & basic Cox):

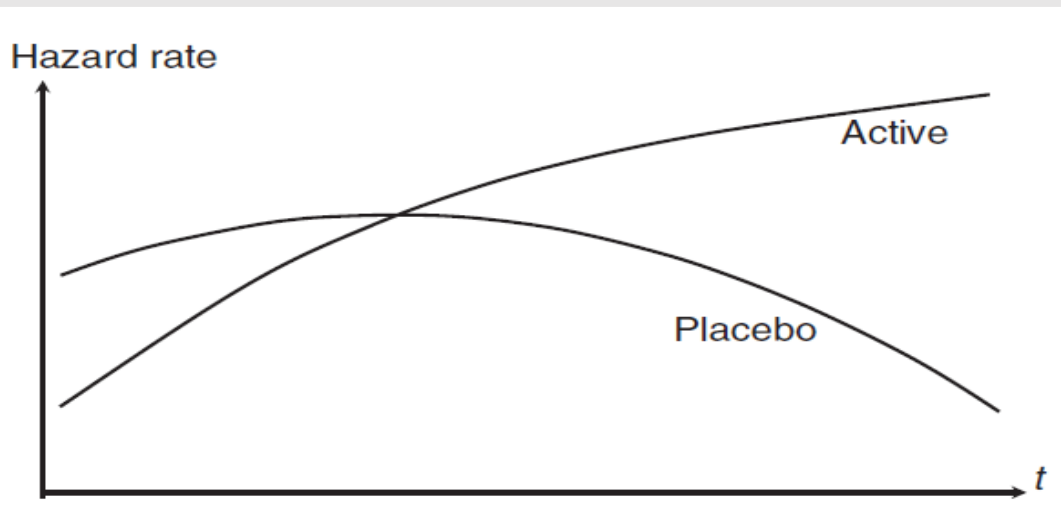
1. Events of the individuals occur **independently** of one another

Acceptable in «**time to the first event**» analyses

2. Hazard of event *at any given time* for an individual in one group is **proportional to** the hazard *at that time* for an individual in the other group...

hazard functions **do not cross** one another

3. Hazard ratios **are independent of time**



what if the '**treatment**' effect changes with time* ?

*...or we have repeated measures of a covariate ???

Last (but not least!):

4. Censoring mechanism is «**independent**» of the event [**conditional** on covariates in Cox]:

Those still at risk at time t are **a random sample** of the population at risk at time t , for all t ...

...is that always true???

Primary and Secondary end point:

Table 2. Efficacy Outcomes.*

Outcome	Apixaban Group (N= 9120)		Warfarin Group (N= 9081)		Hazard Ratio (95% CI)	P Value
	Patients with Event	Event Rate	Patients with Event	Event Rate		
	no.	%/yr	no.	%/yr		
Primary outcome: stroke or systemic embolism	212	1.27	265	1.60	0.79 (0.66–0.95)	0.01
Stroke	199	1.19	250	1.51	0.79 (0.65–0.95)	0.01
Ischemic or uncertain type of stroke	** 162	0.97	175	1.05	0.92 (0.74–1.13)	0.42
Hemorrhagic stroke	40	0.24	78	0.47	0.51 (0.35–0.75)	<0.001
Systemic embolism	15	0.09	17	0.10	0.87 (0.44–1.75)	0.70
Key secondary efficacy outcome: death from any cause	603	3.52	669	3.94	0.89 (0.80–0.998)	0.047

Are patients that die **before** experiencing the primary outcome *similar* to the others?

Regression Models and Life-Tables

BY D. R. COX

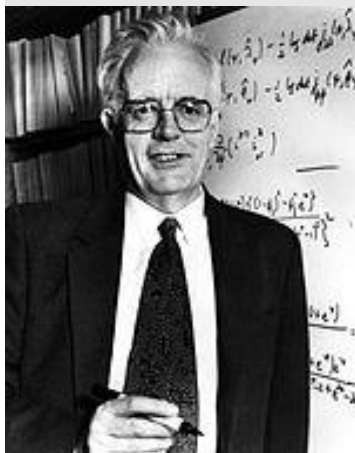
Imperial College, London

[Read before the ROYAL STATISTICAL SOCIETY, at a meeting organized by the Research Section, on Wednesday, March 8th, 1972, Mr M. J. R. HEALY in the Chair]

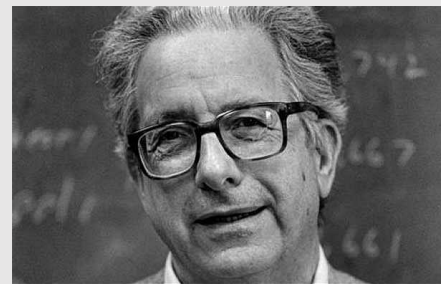
SUMMARY

The analysis of censored failure times is considered. It is assumed that on each individual are available values of one or more explanatory variables. The hazard function (age-specific failure rate) is taken to be a function of the explanatory variables and unknown regression coefficients multiplied by an arbitrary and unknown function of time. A conditional likelihood is obtained, leading to inferences about the unknown regression coefficients. Some generalizations are outlined.

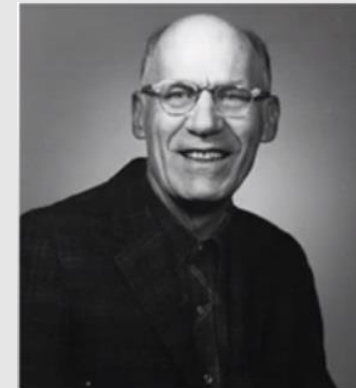
David Cox (1924-2022)



Edward L. Kaplan (1920-2006)



Paul Meier (1924-2011)



NONPARAMETRIC ESTIMATION FROM INCOMPLETE OBSERVATIONS*

E. L. KAPLAN

University of California Radiation Laboratory

AND

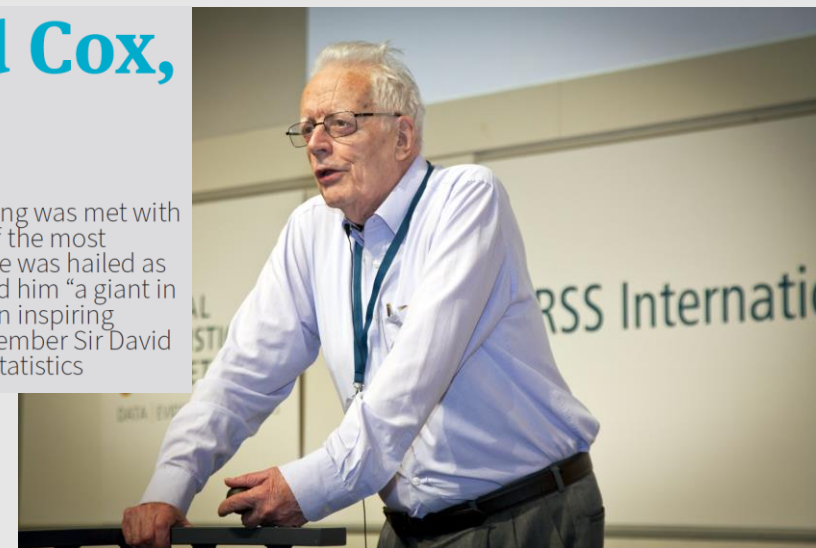
PAUL MEIER

University of Chicago

Sir David Cox and me
(London, sept. 2016)

Remembering Sir David Cox, 1924-2022

Sir David Cox died on 18 January 2022 at the age of 97. News of his passing was met with an outpouring of tributes. To the Royal Statistical Society, he was “one of the most important statisticians of the past century”. At Nuffield College, Oxford, he was hailed as “a pioneering statistician”. The MRC Biostatistics Unit at Cambridge called him “a giant in the field”, while at St John’s College, Cambridge, he was celebrated as “an inspiring scholar”. In this special collection of articles, friends and colleagues remember Sir David in their own way, while also reflecting on his immense contributions to statistics

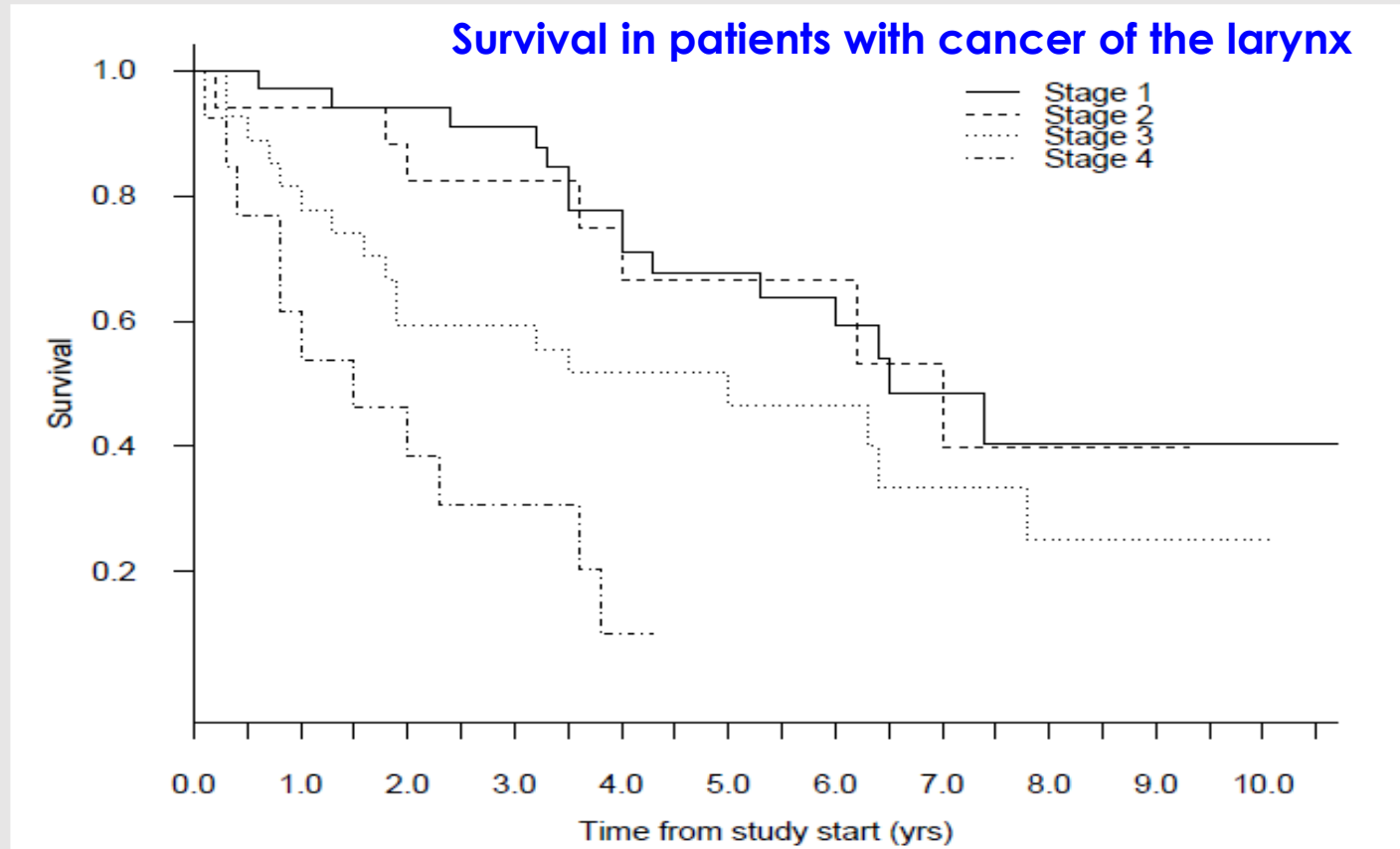


Supplementary materials

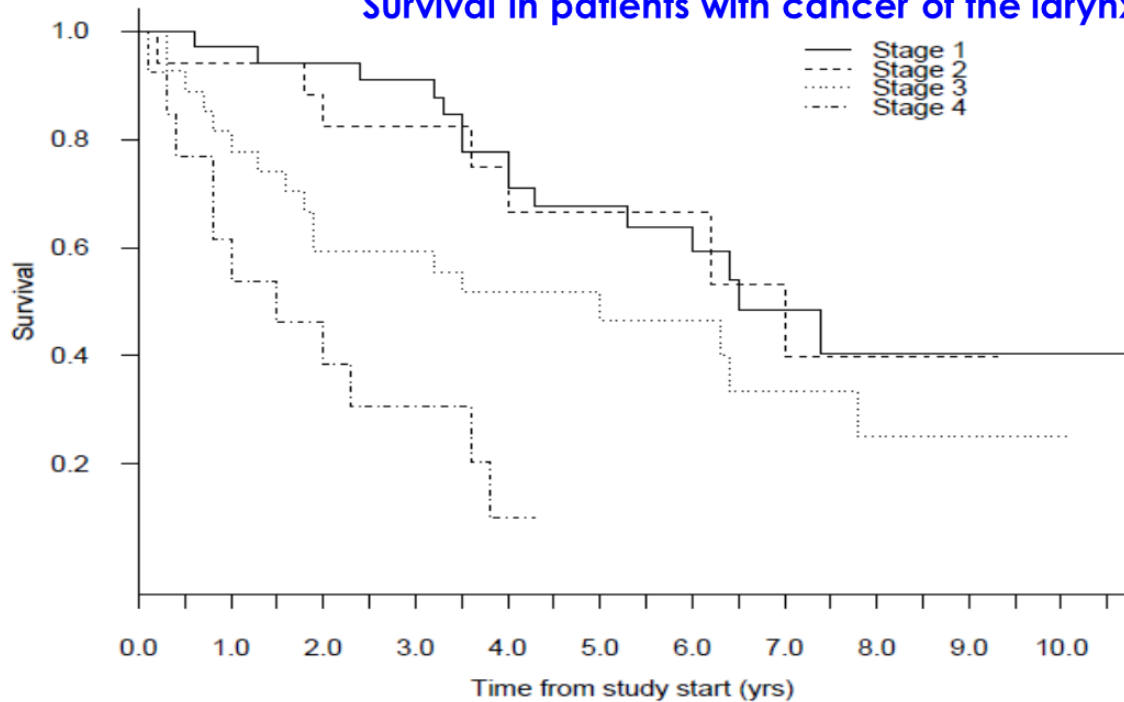
Linear Trend test between survival curves

When there is a **natural ordering** of the groups that we want to compare (i.e: by stage of disease) we can make use of a test **for linear trend**.

In this case, for example, the research question is whether survival **deteriorates** with increasing severity stage, rather than the more general question whether there are **any differences** in survival between stages.



Survival in patients with cancer of the larynx



```
> survdiff( Surv(t2death, death) ~ stage, data=larynx )  
Call:  
survdiff(formula = Surv(t2death, death) ~ stage, data = larynx)
```

	N	Observed	Expected	(O-E)^2/E	(O-E)^2/V
stage=1	33	15	22.57	2.537	4.741
stage=2	17	7	10.01	0.906	1.152
stage=3	27	17	14.08	0.603	0.856
stage=4	13	11	3.34	17.590	19.827

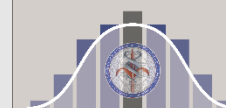
Chisq= 22.8 on 3 degrees of freedom, p= 4.53e-05

Extended Log-Rank test for G groups (no ordering)



- ▶ Time origin: diagnosis with cancer
- ▶ Failure event: death
- ▶ Question of interest: How does survival time from diagnosis to death vary by *stage of disease* at presentation?

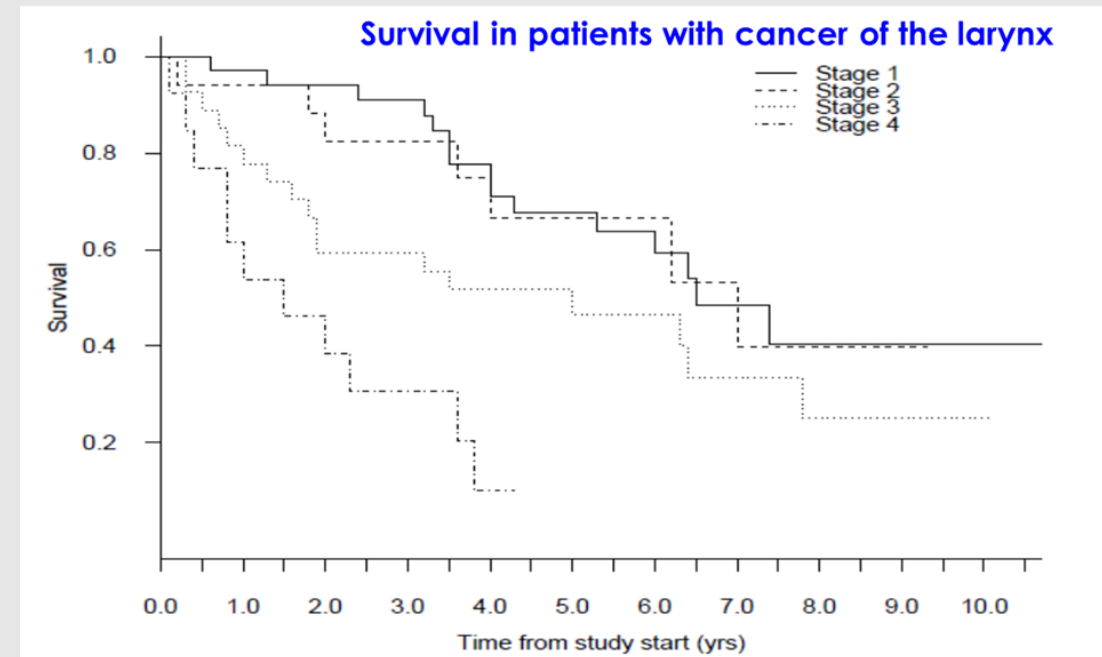
Conclusion: The hypothesis that all four survival curves are equal is clearly rejected. We conclude that **at least one group is different** with respect to survival



This test says nothing about **how** the groups differ; which one is the worst, the best, etc...

That can be further explored with a linear trend test.

- 4 stages of disease recorded at the baseline (the origin)
- 4 stages of disease groups can be ordered in a meaningful way
- H_1 : survival by stage of disease is either **progressively worse** or **progressively better**
- That is, we wish to take advantage of the ordinal nature of the stage of the grouping variable



$$D_j = a_j(O_j - E_j)$$

a_j = "score" of the stage

$$F_j = a_j E_j$$

$$G_j = a_j^2 E_j$$

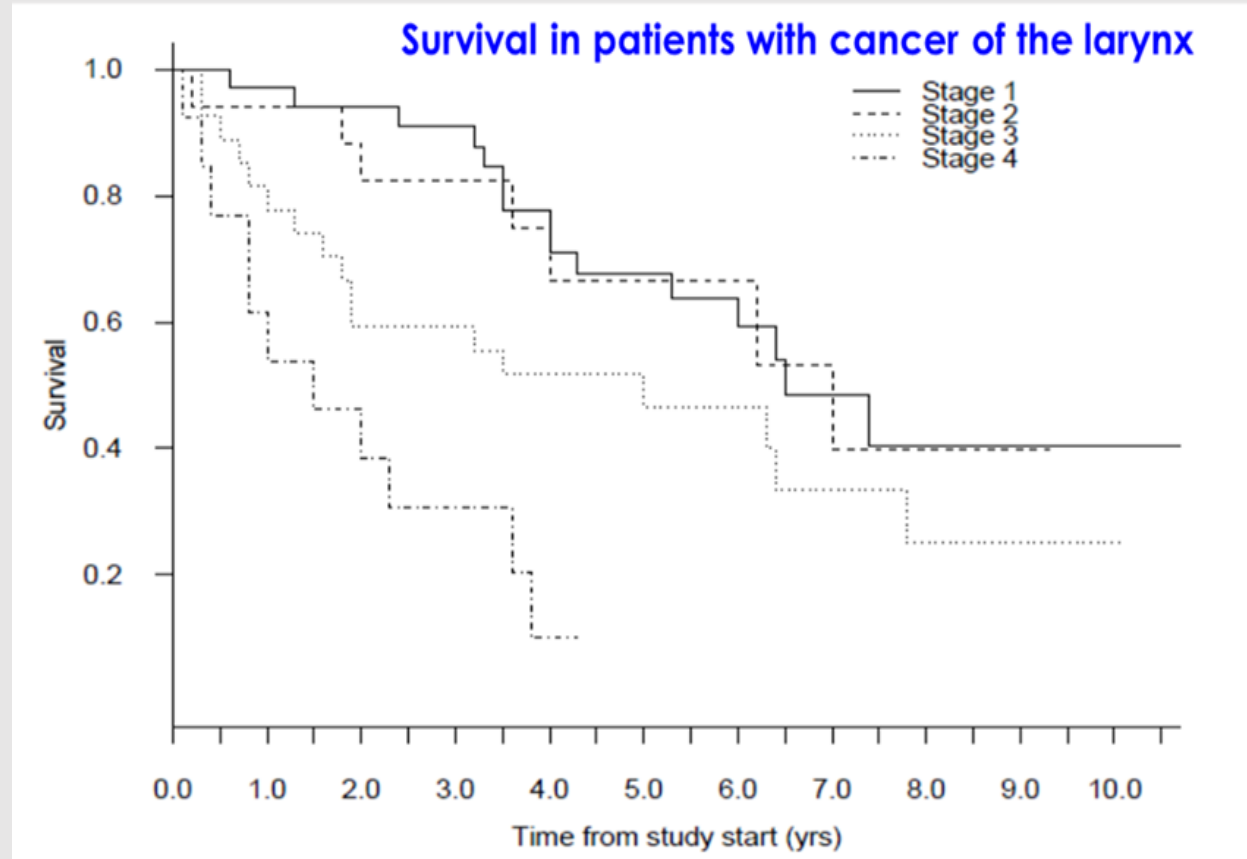
Sum over j

$$\chi^2_{trend} = \frac{D^2}{V_{trend}}$$

$$V_{trend} = G - \left(\frac{F^2}{E} \right)$$

Variance-covariance matrix

Conclusion: Reject the hypothesis that all four survival curves are equal and conclude that stage is **positively associated** with the hazard for death

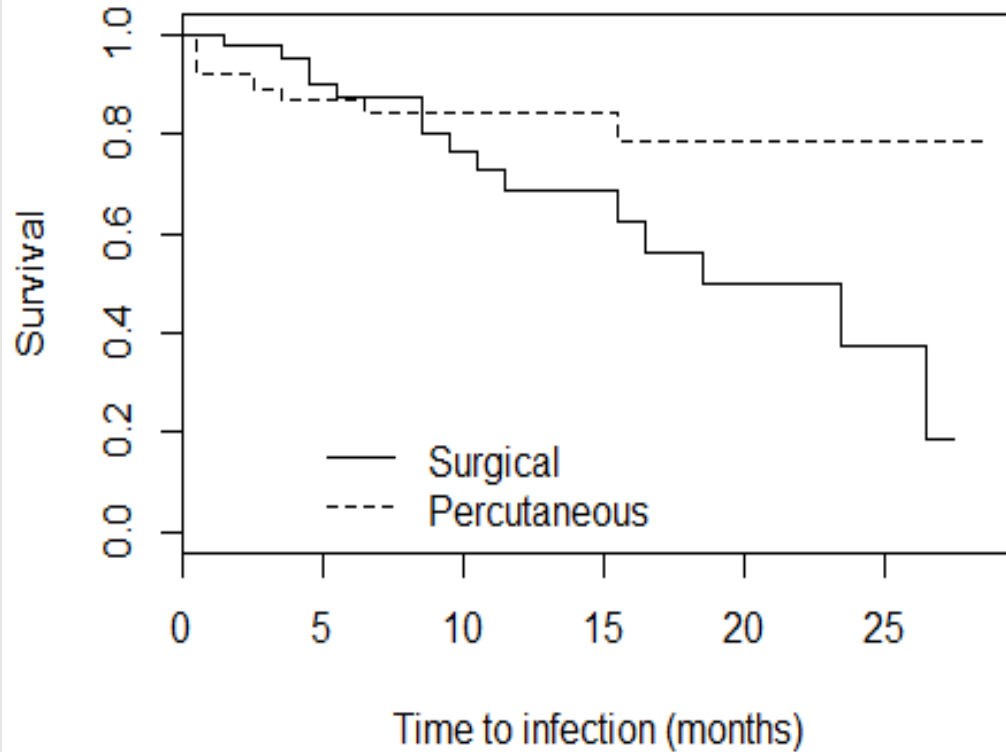


```
Q          Var          Z          pNorm
-25.80    48.15    -3.7190  0.00020005 ***
$scores
[1] 1 2 3 4
```

Kidney dialysis data

The kidney-dialysis trial was designed to assess the time of the infection in patients with renal insufficiency. In 43 patients, the catheter was surgically implanted, whereas in 76 patients it was percutaneously placed.

Kidney dialysis data



Positive differences in favor of the surgical method are cancelled out by the **negative** differences later.

##	Q	Var	Z	pNorm
## 1	-3.96355	6.23675	-1.587104	0.1124892
## n	9.00000	38919.18761	0.045621	0.9636127
## sqrtN	-13.20293	433.84450	-0.633875	0.5261626
## S1	-2.46920	4.37254	-1.180837	0.2376674
## S2	-2.31343	4.20869	-1.127672	0.2594583
## FH_p=1_q=1	-1.02064	0.10661	-3.125846	0.0017729 **

The first 5 tests are unable to detect the overall differences.

Instead, the last, F-H test yield significant results.

Why?

Weighted Logrank tests

Consider weighting (Observed - Expected) *differently* over failure times ($k=1, \dots, D$).

This will enable us to inflate early or late differences

[increased power under non-proportional hazards]

$$T_w = \frac{\sum_{k=1}^D w_k (O_k - E_k)^2}{\sum_{k=1}^D w_k^2 \text{Var}_k}$$

S1 and S2 are based on *pooled* KM estimates, therefore weights depend on the survival experience in the pooled sample

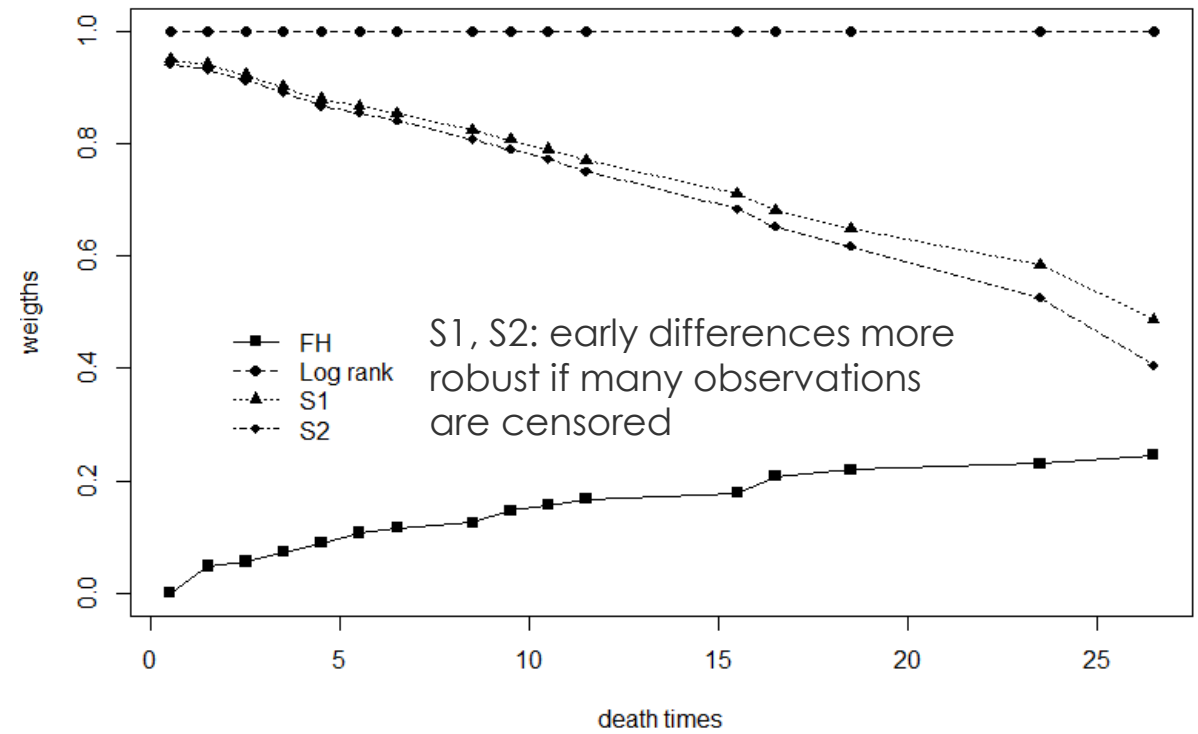
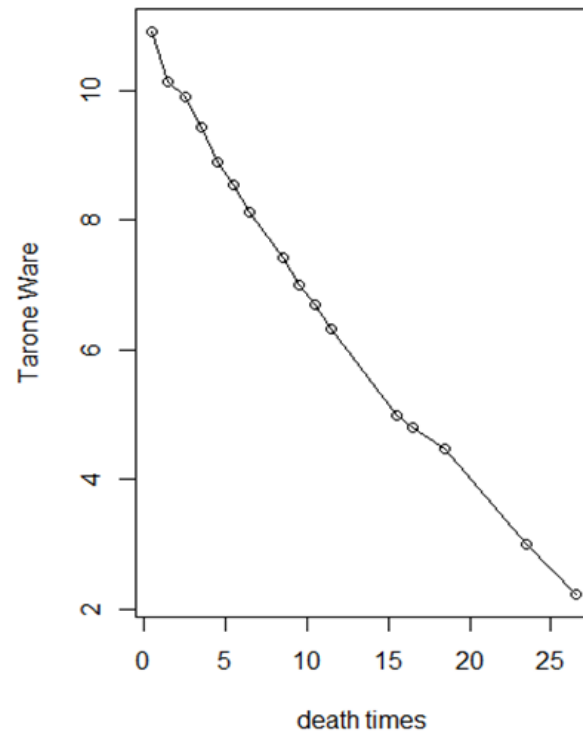
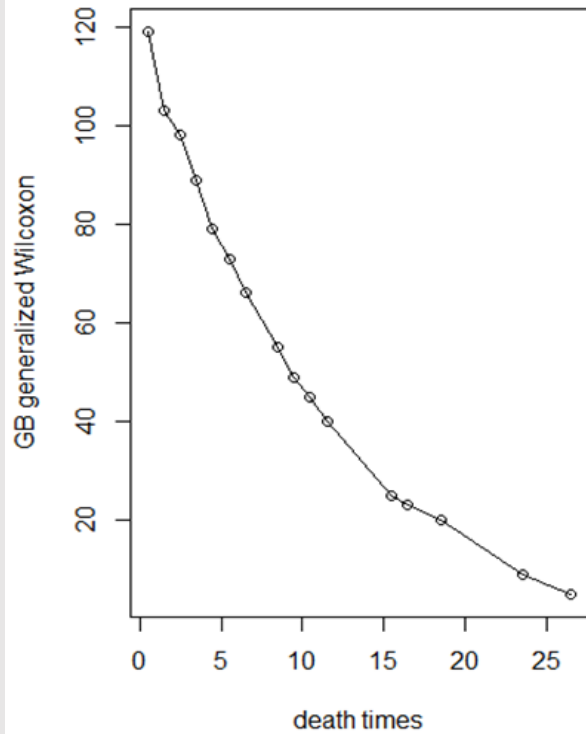
F-H: The weight at $k_1 = 1$ and thereafter is: $(S_{k_{i-1}})^p (1 - S_{k_{i-1}})^q$

In R default p (early)= q (late)=1

Weight	Test
1	log-rank
n_k (#pts at risk at k)	Gehan-Breslow generalized Wilcoxon
$\sqrt{n_k}$	Tarone-Ware
S1: \hat{S}_{k-1}	Peto-Peto's modified survival estimate
S2: \hat{S}'_{k-1}	modified Peto-Peto (by Andersen)
FH	Fleming-Harrington

Different *weights* are applied to differences between expected and observed deaths to *emphasize* certain times more than others...

early differences in survival times



How should weights be chosen?

For scientific inference it is not reasonable to look at the survival curves first, then choose weights.

A priori knowledge: is there a reason to believe we will have non PH?

-> If not, go with the logrank test

-> If so, consider what survival differences are most clinically meaningful (early vs late)

-> Childhood cancer (**late** differences)

-> Late stage lung cancer remission (**early** differences)

Renyi Type Tests

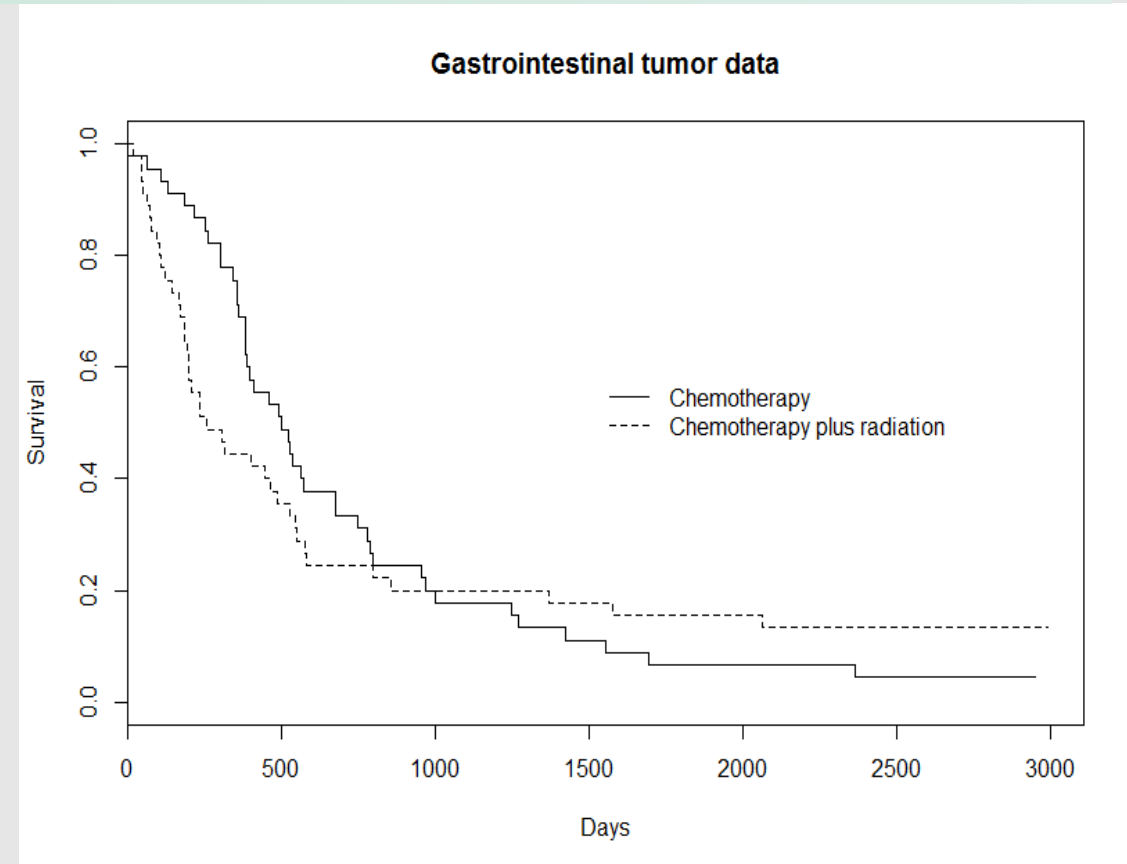
alternatives to log-rank test

A class of tests with power to detect crossing hazards.

A clinical trial of chemotherapy vs chemotherapy + radiotherapy in the treatment of locally unresectable gastric cancer. 45 patients were randomized to each of the two arms and followed for about 8 years.

During the initial 1000 days, chemotherapy showed a **higher** survival rate, whereas chemotherapy+radiation was associated with an increased number of early deaths, which may be attributable to the progression of tumors within the radiation field or to complications.

However, chemotherapy+radiation appear to offer **better prospects for long-term** survival during the late follow-up period.



Renyi Type Tests

alternatives to log-rank test

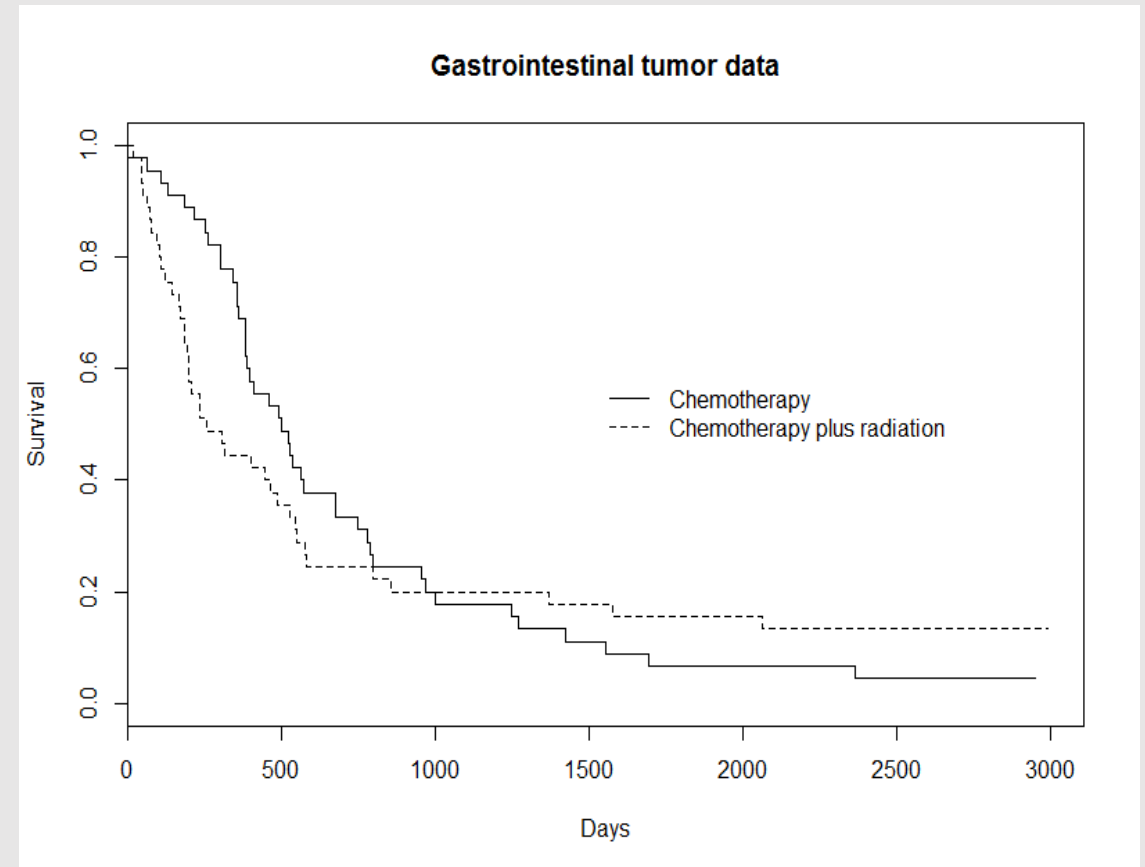
Renyi tests are based on **sequence** of test statistics; greater power to detect crossing hazard rates.

Value of the test statistic [for some weight function] **at each death time** is computed.

When the hazard rates cross, the **absolute value** of these sequential evaluations will have a **maximum value** at some **time point** [prior to the largest death time].

When this value is 'too large', the null hypothesis is rejected.

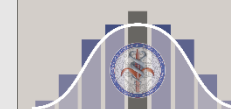
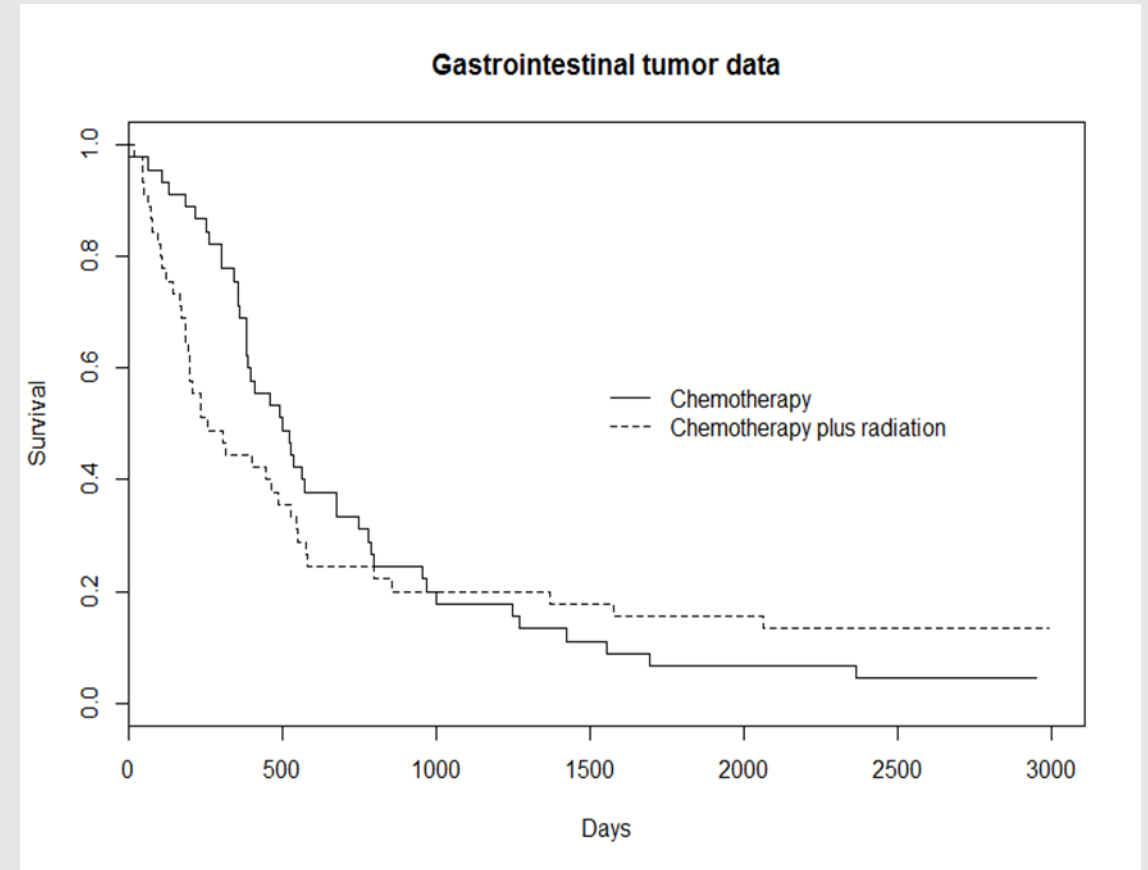
To adjust for the fact that **multiple test statistics are estimated on the same set of data**, a correction is made to the critical value of the test.



N.B.: a significant Renyi test statistic do not indicate that the **global** survival rate in the chemotherapy only group is significantly **greater** than the chemotherapy + radiation.

It does mean that **there is a significant difference in the survival rates between the two groups at some time point**, but does not imply that one group is superior to another.

```
##          Q          Var          Z    pNorm
## 1      2.1463e+00  1.9862e+01  0.48159 0.630098
## n      4.9100e+02  6.0322e+04  1.99913 0.045594 *
## sqrtN  4.3629e+01  9.8798e+02  1.38803 0.165129
## S1     5.4126e+00  7.2723e+00  2.00710 0.044739 *
## S2     5.3864e+00  7.0411e+00  2.02993 0.042364 *
## FH_p=1_q=1 -8.9383e-02  7.1790e-01 -0.10549 0.915985
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##          maxAbsZ          Var          Q    pSupBr
## 1          9.8049      19.8617  2.2001  0.0556044 .
## n         725.0000  60322.4412  2.9519  0.0063169 **
## sqrtN     84.1532   987.9774  2.6773  0.0148437 *
## S1         7.9752     7.2723  2.9574  0.0062054 **
## S2         7.8688     7.0411  2.9654  0.0060450 **
## FH_p=1_q=1  1.3396     0.7179  1.5811  0.2277168
```



The Cox model assumes that the **hazards are proportional** (PH), which means that the hazard ratio is constant over time with different predictor or covariate levels.

This PH assumption in any covariate is quite a strong assumption. Considering the complexity of biological and physiological responses and associations, this assumption has rarely a solid justification.

If PH doesn't exactly hold for a particular covariate but we fit the PH model anyway, then what we are getting is sort of an **average HR**, averaged over the event times.

The two most common ways to assess the PH assumption are:

- Visual assessment by means of *the log-cumulative hazard plot*
 - Testing of *scaled Schoenfeld residuals*

Eventually, if the non-PH variable is a categorical one it could make sense using a stratified approach

The Stratified Cox Model

Suppose a confounder C has k levels on which we would like **to stratify** when comparing $h(t | E)$ and $h(t | \text{not } E)$ where E is an indicator of “exposure”.

$$h_i(t|E) = h_{0i}(t)\exp(E\beta)$$

$$i = 1, \dots, k$$

1. A [non-parametric] baseline hazard is estimated **within** each stratum (solve ev. non PH hazard)
2. If the confounder is controlled using stratification, there is no way to estimate an hazard ratio comparing two levels of the confounder.
3. Stratification generally requires **more data** to obtain the same precision in coefficient estimates

$$h_i(t|X_i) = h_0(t)\exp(X_i\beta)$$

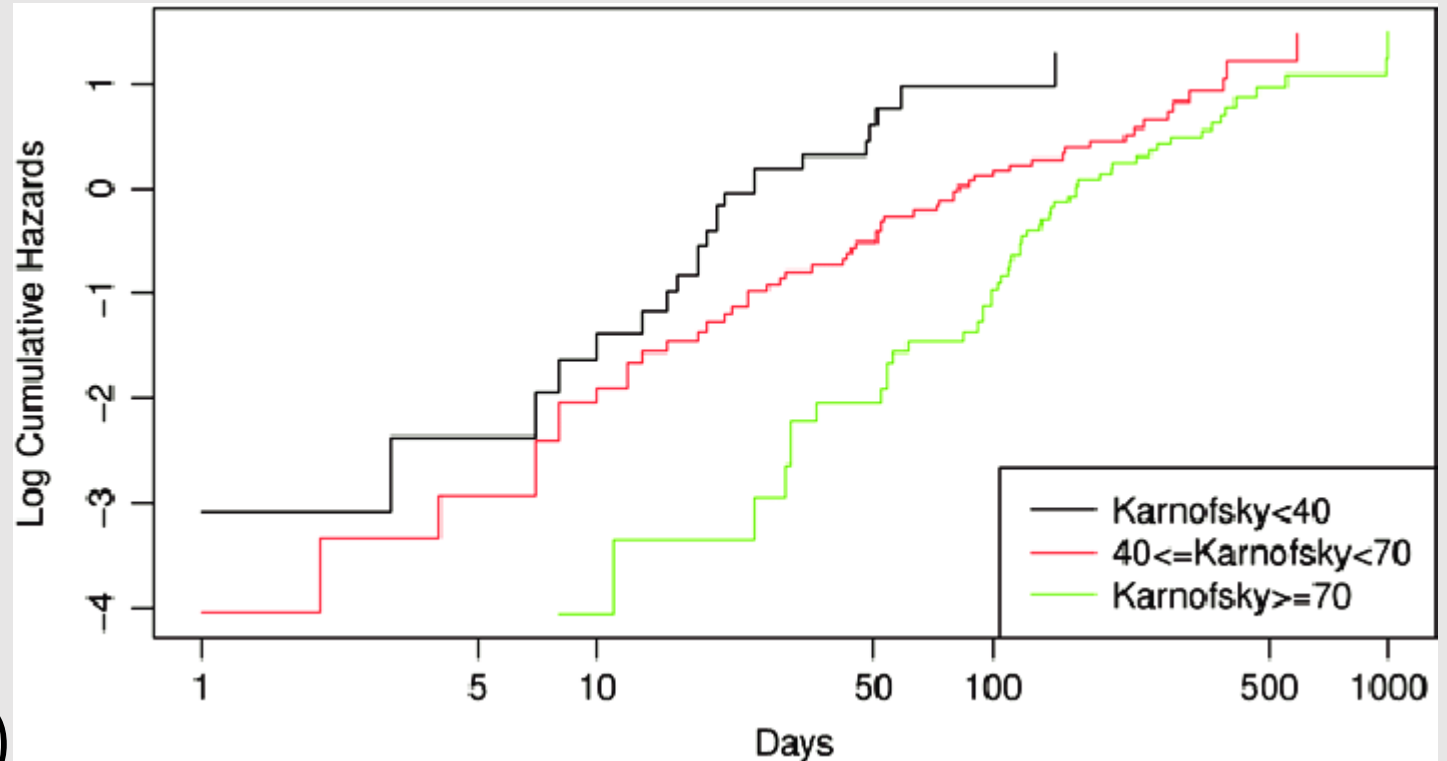
If the **log-cumulative hazards** for individuals with different values of X are plotted against time, the curves will be **parallel** if the PH assumption is valid.

$$\int_0^t h_i(u)du = \exp(X_i\beta) \int_0^t h_0(u)du$$

$$H_i(t|X_i) = \exp(X_i\beta)H_0(t)$$

Cumulative hazard functions

$$\log(H_i(t|X_i)) = X_i\beta + \log(H_0(t))$$



- Values of X need to be categorical/grouped
- Just a visual appreciation

Schoenfeld residuals

Time-varying residuals from the model are added to the corresponding **time-invariant** coefficient estimate β and smoothed. The result is a **plot** of an estimate of the regression coefficient for the covariate **over time**. If the plot is **reasonably flat** (there is here a formal test), the PH assumption holds.

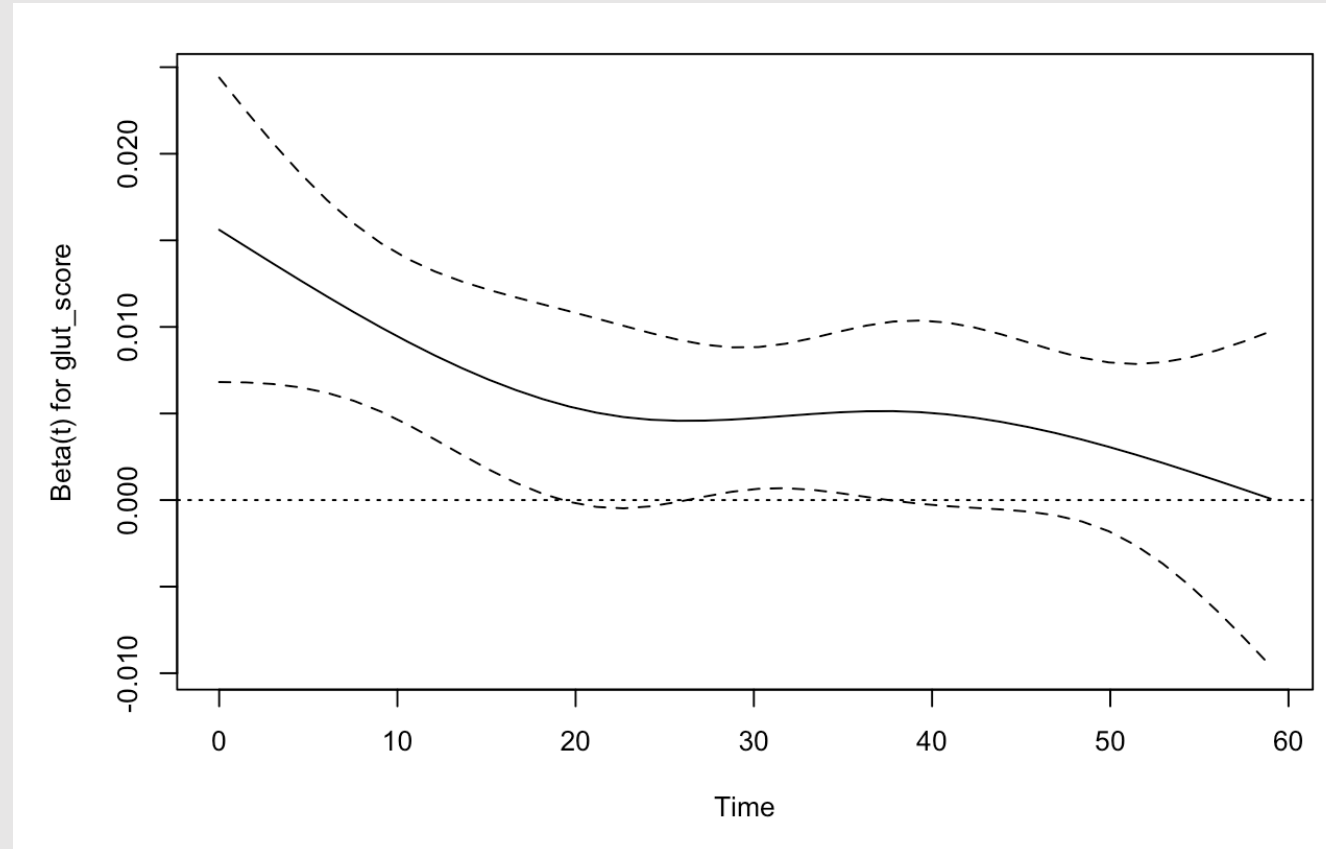
$$s_{k,j}$$

Schoenfeld residual for covariate X_j at time t_k

$$E(s_{k,j}) + \hat{\beta}_j \approx \beta_j(t_k)$$

The Schoenfeld residuals are the differences between that individual's covariate values at the event time and the corresponding risk-weighted average of covariate values among all those at risk at that time.

The word "residual" thus makes sense, as it's the difference between an observed covariate value and what you might have expected based on all those at risk at that time.



- `library(ISwR)`
- `data(melanom)`

status: indicator of the patient's status by the end of the study:

1="dead from malignant melanoma"

2= "alive"

3= "dead from other causes"

days: observation time in days

ulc: 1=present (tumor ulcerated) 2 = absent

thick: tumor thickness


sex: 1 for women and 2 for men

	no	status	days	ulc	thick	sex
1	789	3	10	1	676	2
2	13	3	30	2	65	2
3	97	2	35	2	134	2
4	16	3	99	2	290	1
5	21	1	185	1	1208	2
6	469	1	204	1	484	2
7	685	1	210	1	516	2
8	7	1	232	1	1288	2
9	932	3	232	1	322	1
10	944	1	279	1	741	1
11	558	1	295	1	419	1
12	612	3	355	1	16	1
13	2	1	386	1	387	1
14	233	1	426	1	484	2
15	418	1	469	1	242	1
16	765	3	493	1	1256	2
17	777	1	529	1	580	2
18	61	1	621	1	706	2
19	67	1	629	1	548	2
20	819	1	659	1	773	2
21	10	1	667	1	1385	1
22	15	1	718	1	234	2


Consider a model with the single regressor sex:

```
mod.sex <- coxph(Surv(days,status==1)~sex)
summary(mod.sex)
```

```
##      coef exp(coef) se(coef)      z Pr(>|z|)
## sex 0.6622  1.9390  0.2651 2.498  0.0125 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##      exp(coef) exp(-coef) lower .95 upper .95
## sex      1.939      0.5157   1.153   3.26
##
## Concordance= 0.59 (se = 0.033 )
## Rsquare= 0.03 (max possible= 0.937 )
## Likelihood ratio test= 6.15 on 1 df,  p=0.01314
## Wald test               = 6.24 on 1 df,  p=0.01251
## Score (logrank) test = 6.47 on 1 df,  p=0.01098
```



Males (=2) have an hazard nearly twice than women (=1)



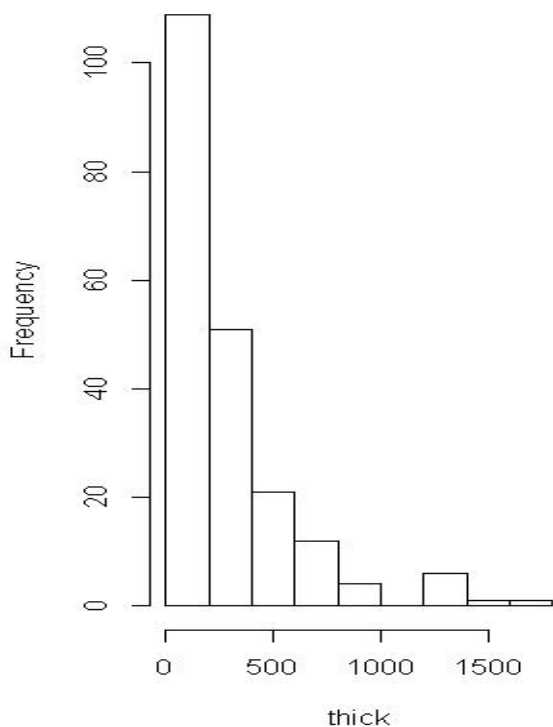
These tests are all equivalent in large samples but may differ somewhat in small-sample cases

A more elaborate example, involving also a continuous covariate:

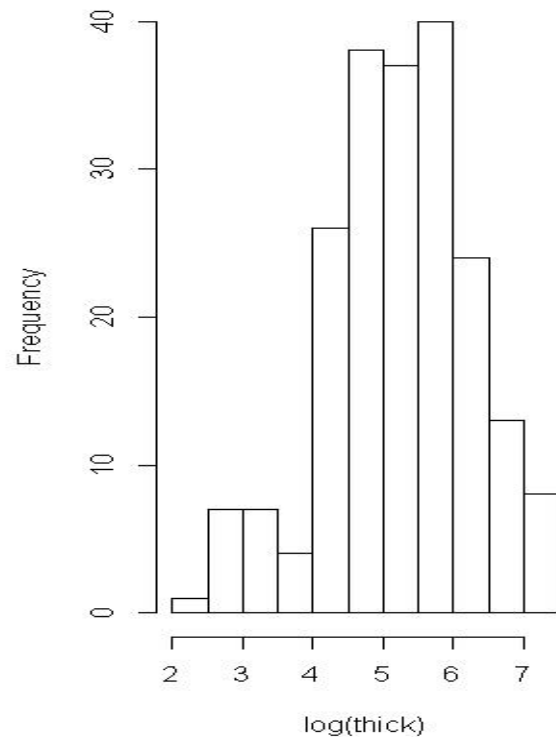
```
mod.cov <- coxph(Surv(days,status==1)~sex+log(thick))
summary(mod.cov)
```

'thick' is the tumor thickness; we use logarithm since the distribution is asymmetric:

Histogram of thick



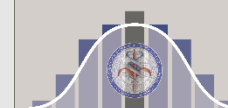
Histogram of log(thick)



HR of $\log(\text{thick})=2.18$
 each 1 point change in $\log(\text{thick})$
 is associated with a 2.2-fold
 increase in a patient's risk

```
##           coef exp(coef) se(coef)      z Pr(>|z|)
## sex          0.4580   1.5809   0.2687  1.705  0.0883 .
## log(thick)  0.7809   2.1834   0.1573  4.963  6.94e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##           exp(coef) exp(-coef) lower .95 upper .95
## sex              1.581    0.6326    0.9337    2.677
## log(thick)       2.183    0.4580    1.6040    2.972
```

Note that taking into account $\log(\text{thick})$ the effect of sex is reduced...



Assessing the PH Assumption (I)

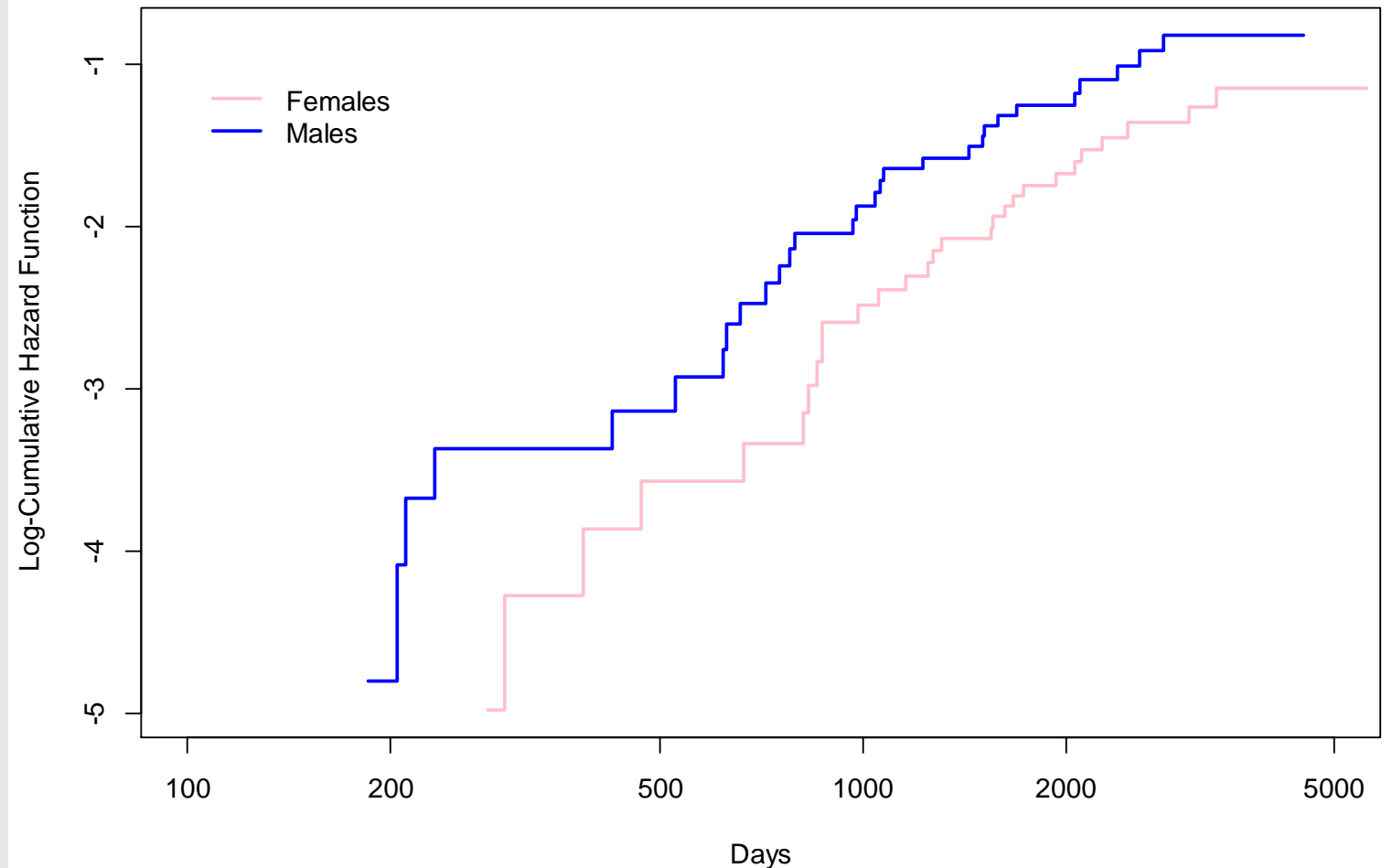
Adjusting for $\log(\text{thick})$ does the effect of gender follow a PH model?

If the PH assumption holds, the log cumulative hazards for the two groups, adjusting for $\log(\text{thick})$, should be roughly parallel...

Conclusion: not strong evidence of non-PH.

This is a good look at gross departures, but it is far from a formal test...

```
fit1 <- coxph( Surv(days,status==1) ~ log(thick)+ strata(sex))
```



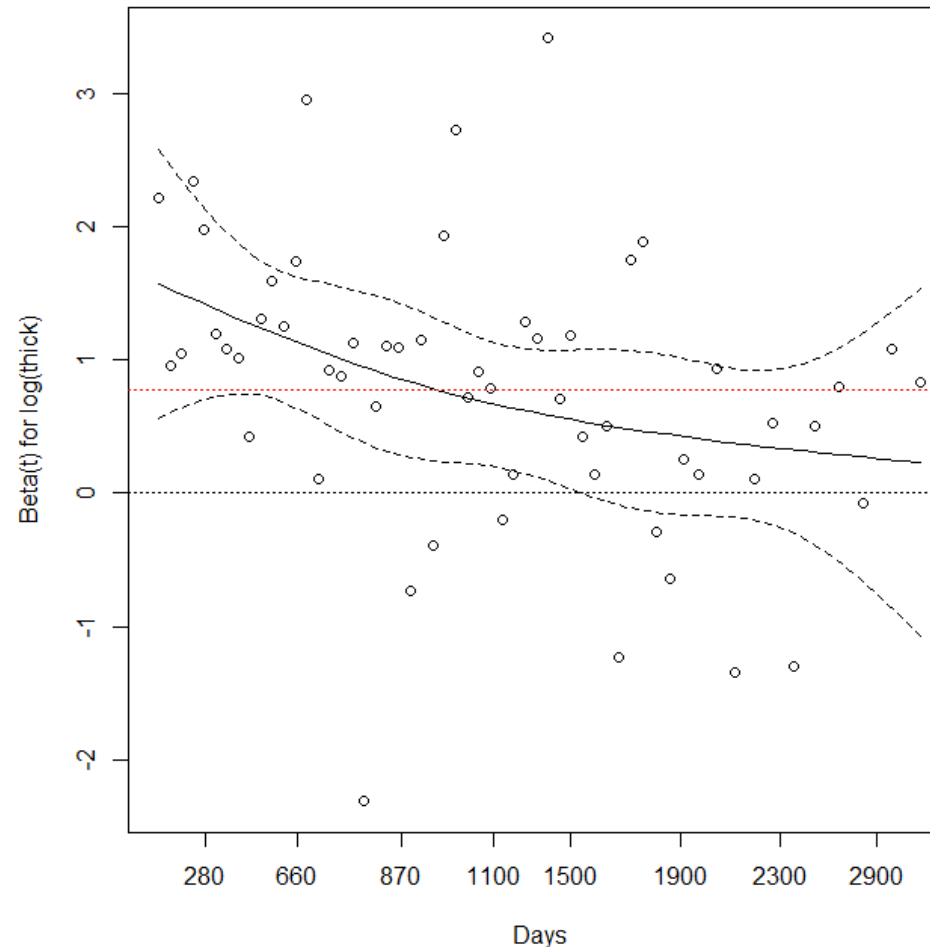
Assessing the PH Assumption (II)

```
check.ph <- cox.zph(mod.cov, transform="km", global=TRUE)
```

Adjusting for gender, does the effect of $\log(\text{thick})$ follow a proportional hazards model?

##		rho	chisq	p
##	sex	-0.102	0.587	0.4436
##	log(thick)	-0.352	5.485	0.0192
##	GLOBAL	NA	6.813	0.0332

The effect of $\log(\text{thick})$ is gradually decreasing with time.



If the PH assumption holds, then the plot of $\beta(t)$ vs time *should* be on a horizontal line.

Cox model's estimate for «overall» log thick effect

Possible solutions to non-proportionality (I):

- **Stratification:** covariates with non PH effects may be used as strata
 - no **direct test** of association with survival;
 - ok for **categorical** covariates, discretization for continuous ones (could be problematic)
 - **less efficient** analyses (usually larger sample size needed)
- **Partition** of the time axis: the PH could be *valid* in some time intervals (**landmark analysis**)

Alternative/Advanced methods:

- **Cox model with time-varying** coefficients: model the dependence of beta on time
Not easy to find the appropriate function... interpretation more complex

- Use a **different approach: Flexible Parametric Survival and Multi-State Models**

Accelerated Failure Time (AFT) Models:

- The `survreg` function in package [survival](#) can fit an accelerated failure time model.
- A modified version of `survreg` is implemented in the [rms](#) package (`psm` function).
- The [eha](#) package also proposes an implementation of the AFT model (function `aftreg`).
- The [NADA](#) package proposes the front end of the `survreg` function for left-censored data.
- The [simexaft](#) package implements the Simulation-Extrapolation algorithm for the AFT model, that can be used when covariates are subject to measurement error.
- A robust version of the accelerated failure time model can be found in [RobustAFT](#).
- The [coarseDataTools](#) package fits AFT models for interval censored data.
- An alternative weighting scheme for parameter estimation in the AFT model is proposed in the [imputeYn](#) package.
- The [AdapEnetClass](#) package implements elastic net regularisation for the AFT model.

Additive Models:

- Both [survival](#) and [timereg](#) fit the additive hazards model of Aalen in functions `aareg` and `aalen`, respectively.
- [timereg](#) also proposes an implementation of the Cox-Aalen model (that can also be used to perform the Lin, Wei and Ying (1994) goodness-of-fit for Cox regression models) and the partly parametric additive risk model of McKeague and Sasieni.
- A version of the Cox-Aalen model for interval censored data is available in the [coxinterval](#) package.
- The [uniah](#) package fits shape-restricted additive hazards models.
- The [addhazard](#) package contains tools to fit additive hazards model to random sampling, two-phase sampling and two-phase sampling with auxiliary information.

Flexible survival models:

- **flexsurv: Flexible parametric models for time-to-event data**
- **rstpm2: Smooth Survival Models, Including Generalized Survival Models**

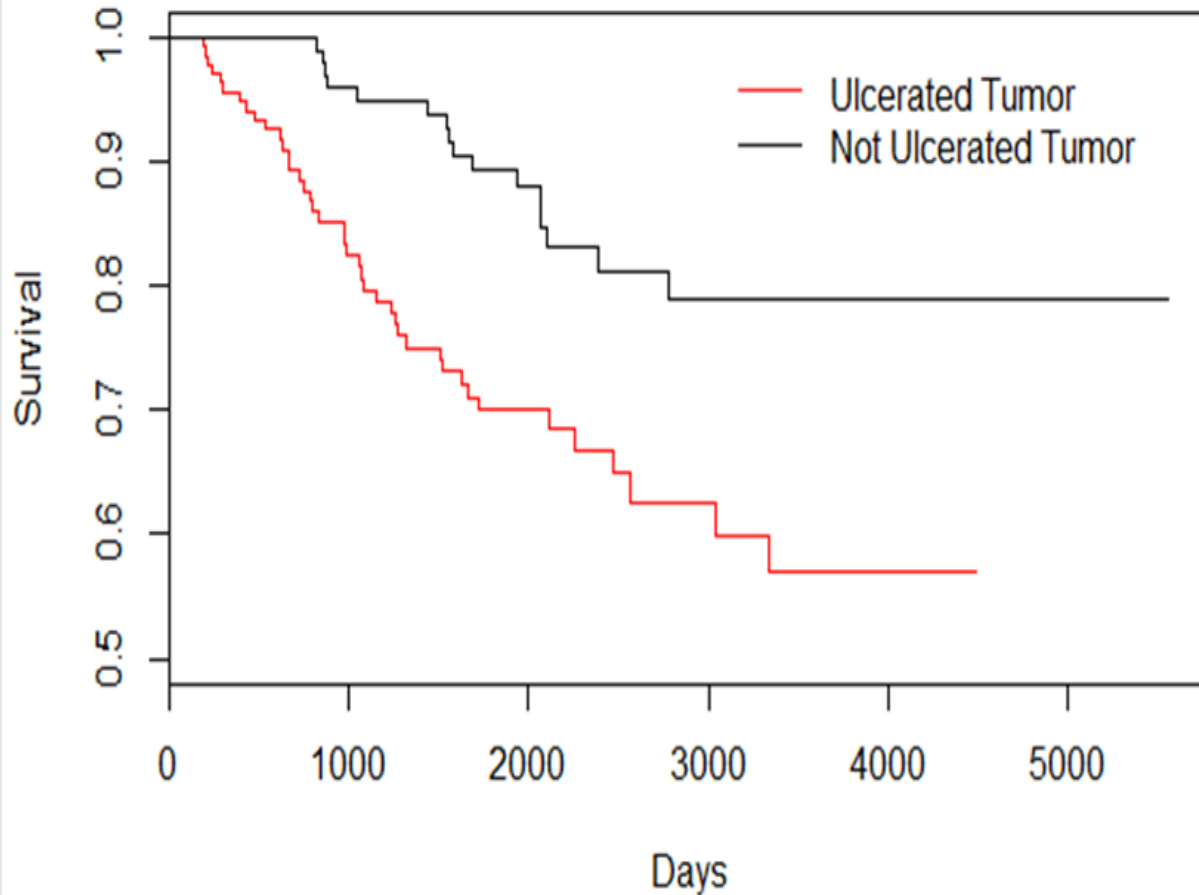
A more elaborate example: binary factor + continuous covariate + **stratification** variable:

```
mod.cov.strat <- coxph(Surv(days,status==1)~sex+log(thick)+strata(ulc))
summary(mod.cov.strat)
```

```
##              coef exp(coef) se(coef)      z Pr(>|z|)
## sex           0.3600   1.4333  0.2702  1.332  0.1828
## log(thick)    0.5599   1.7505  0.1784  3.139  0.0017 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##              exp(coef) exp(-coef) lower .95 upper .95
## sex              1.433     0.6977    0.844    2.434
## log(thick)       1.750     0.5713    1.234    2.483
##
## Concordance= 0.673 (se = 0.058 )
## Rsquare= 0.063 (max possible= 0.9 )
## Likelihood ratio test= 13.3 on 2 df,  p=0.001296
## Wald test            = 12.88 on 2 df,  p=0.001598
## Score (logrank) test = 12.98 on 2 df,  p=0.00152
```

Stratifying by the presence or absence of ulcer, significance of the log(thick) has been reduced and sex is no longer significant.

We can plot survival curves estimated for each strata by using `survfit` on the output of `coxph`:



The default for `survfit` is to generate curves for a *pseudoindividual* for which the covariates are at their mean values.

In the present case, that would correspond to a tumor thickness of 1.86 mm and a gender of 1.39 (!)...

... we have been sloppy in not defining sex as a factor variable, but that would not actually give a different result (HR): `coxph` subtracts the means of the regressors before fitting, so a 1/2 coding is the same as 0/1, which is what a factor with treatment contrasts usually gives.

[But, defining the factor we can define “hypothetical” pts with certain values for the covariates]


```
sex.f <- as.factor(sex)
```



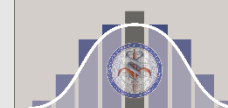
Converting sex into a factor

```
mod.cov.strat.f <- coxph(Surv(days,status==1)~sex.f+log(thick)+strata(ulc))  
summary(mod.cov.strat.f)
```

```
mod.cov.strat.f <- coxph(Surv(days,status==1)~sex.f+log(thick)+strata(ulc))  
summary(mod.cov.strat.f)
```

```
## Call:  
## coxph(formula = Surv(days, status == 1) ~ sex.f + log(thick) +  
##       strata(ulc))  
##  
## n= 205, number of events= 57  
##  
##              coef exp(coef) se(coef)      z Pr(>|z|)  
## sex.f2      0.3600   1.4333  0.2702  1.332  0.1828  
## log(thick)  0.5599   1.7505  0.1784  3.139  0.0017 **  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Now sex.f2
indicates that HR
refers to the
contrast
of level "2"
versus level "1"
for the factor
variable sex,
[the same HR
value as before]



To estimate survival curves for subjects with *certain* values of the covariates, we could use the option `newdata` in `survfit`:

