

## Déjà

**1 (homogeneous) sample:** "Survival" / "Time-to-event" data:

- [equivalent] Functions:  $S[t]$  , hazard  $h[t]$  , pdf[ $t$ ]
- Links: e.g.  $S[t] = \exp[-\int_0^t h[u] du]$  , integral from  $u=0$  to  $u=t$
- Summaries of these functions (e.g.  $T_{25}$ ,  $T_{50}$ ,  $S[ T ]$  )
- *Non-Parametric / Semi-Parametric* Estimation (point & interval) of  $S[t]$  ,  $h[t]$  and pdf[ $t$ ]  
 --- Lifetable [fixed-] --- K-M / N-A [data-determined intervals]
- *Censored data not necessarily " time - to - event "*  
*Y = PSA levels < detection limit, salaries in intervals, distance travelled on set of tires, pages on single ink cartridge, etc.*
- *'1 (homogeneous) sample' structure*  
 => think of as "intercept-only" regression model

**Comparison of 2 Survival/Hazard Curves or Distributions**

- think of as regression model with single binary  $X$
- Risk Sets (match on time of event)
- Adjusted comparisons (non-regression methods)

**Not covered: Parametric models for Lifetime Distributions**

**SAS LIFEREG** procedure fits parametric models to failure time data that can be right, left, or interval censored. The models for the response variable consist of a linear effect composed of the covariates and a random disturbance term. The distribution of the random disturbance can be taken from a class of distributions that includes the extreme value, normal, logistic, and, by using a log transformation, the exponential, Weibull, lognormal, loglogistic, and gamma distributions.

**Stata streg** performs maximum likelihood estimation of parametric regression survival-time models. Survival models currently supported are exponential, Weibull, Gompertz, lognormal, log-logistic and generalized gamma. Also see help `stcox` for estimation of proportional hazards models.

**R survival package:** Regression for a Parametric Survival Model: These are all time-transformed location models, with the most useful case being the accelerated failure models that use a log transformation.

**(Parametric) Regression Models for Rates**

- Model (event) *rates* or *hazards*
- (i) Models with 'multiplicative' rates/hazards  
 $\log[\text{rate}]$  or  $\log[\text{hazard}] = B_0 + B_t t + \dots B_1 X_1 + B_2 X_2 \dots$   
 Rates/hazards PROPORTIONAL (rate ratio parameter constant over time-bands and covariate patterns... if no product terms for 'effect modification/interaction')

In Generalized Linear Model, model the *numbers* of events, with log link .. and log(PT) as offset

$\exp[ B ]$ : rate *ratio* (RR) contrasting rates for two  $X$  (or  $t$ ) values 1 unit apart

- (ii) Models with additive rates/hazards

rate or hazard =  $B_0 + B_t t + \dots B_1 X_1 + B_2 X_2 \dots$

Not as 'natural'. See pp59-- from Chapter 2 of Volume I of Breslow and Day (in Resources) for empirical evidence for proportional rate models (constant rate ratio models) over to additive rates models (constant rate difference models) in cancer epidemiology.

N.B.: B&D use the term "relative Risk" very loosely, when in fact they mean "relative Rates" or rate Ratios.

In Generalized Linear Model, model the *numbers* of events, with identity link .. no intercept (no cases if denominator is zero) and (as regressors) product of PT denominator with each regressor in the rate model.

$B$ : rate *difference* contrasting rates for two  $X$  (or  $t$ ) values 1 unit apart

## New

## Semi-Regression Models

- Restrict to multiplicative rate/hazard models
- Avoid modelling the nuisance ("t") part  
don't fit parameters that (a) are not our focus (b) waste "d.f."
- Use risksets & conditioning to reduce # parameters and avoid having to model "t".
- Choice of Time-scale and "Time-zero" is important  
(has implications for risksets)
- Models, and conditioning as a way of eliminating parameters, applicable to matched case-control studies and even to c-c and other (e.g. consumer choice\*) studies with no 'time' element ["conditional logistic regression"]

(\* Daniel McFadden shared the Nobel Prize for his development of theory and methods for analyzing discrete choice in Economics:  
<http://www.nobel.se/economics/laureates/2000/mcfadden-autobio.html>)

$$\log[\text{rate}] \text{ or } \log[\text{hazard}] = U[t] + \dots B_1 X_1 + B_2 X_2 \dots$$

or

$$[t | X_1 X_2 \dots] = \text{hazard} = \exp[U[t]] \times \exp[B_1 X_1 + B_2 X_2 \dots]$$

where "U[t]" stands for "*unspecified* function of t" for log of hazard function for persons with covariate pattern  $\{X_1=0, X_2=0, \dots\}$ . Note that  $\exp[U[t]]$  is often denoted by  ${}_0[t]$  and often called the "baseline" hazard function.

N.B.: the word "baseline" does *not* refer to measurements (covariates) recorded at  $T=0$ . Rather, it refers to the "reference" category or covariate pattern, against which all other categories or covariate patterns are compared.

Thus, it has the same meaning as the "corner" or "point of departure" category used by Clayton and Hills (eg. "40-49 year olds, unexposed" in the regression example in Table 22.6 p 221 of Clayton and Hills.

The curve  ${}_0[t] = \exp[U[t]]$  is the "intercept curve".

 $S_x[t]$  in terms of ( $h_x[t]$  and)  $S_0[t]$ 

- Remember general law:  $S[t] = \exp[-H[t]]$ ,  
where  $H[t]$  is the integrated or "cumulative" hazard.
- Relationship b/w  $S[t]$  for  $x=1$  and  $S[t]$  for  $x=0$  ["corner"]

if  $[t | x=1] = {}_0[t] \times \exp[-B \times 1] = {}_0[t] \times \text{HazardRatio}$   
then

$$\text{integral of } {}_1[t] = \text{HazardRatio} \times \text{intergral of } {}_0[t]$$

so,

$$\begin{aligned} S[t | x = 1] &= \exp[-H[t | x = 0] \times \text{HazardRatio}] \\ &= \{ \exp[-H[t | x = 0]] \}^{\text{HazardRatio}} \\ &= \{ S_0[t] \}^{\text{HazardRatio}} \end{aligned}$$

i.e., **S curve is constant power of "baseline" curve**

## Test of Proportionality

Two  $\log[-\log[S]]$  functions (for  $x=1$  &  $x=0$ ) should be parallel

- $H[t]$  is the integrated or "cumulative" hazard
- $-\log[S] = H[t]$ , so  $-\log[S_1[t]] = HR \times \{-\log[S_0[t]]\}$

2  $-\log[S]$  curves should be proportional  
(easier to judge if these parallel than hazards proportional)

- use as test of proportionality assumption
- hazard functions may not be stable enough  
(so cannot assess whether 2  $h[t]$  curves are proportional)

## Choice of time scale for Cox model

- the one over which the hazard function is the most difficult to model .. avoid this challenge: *match* risksets on this scale.

```

Framingham study: TIME Scale = YEAR_of_research_grant ...
FU_AGE | risk set (vertical) based on deaths in calendar (project) year
88 + ^
87 + ^
86 + time scale is 'rough', because of 2-year cycles ^
85 + ^
84 + ^
83 + ^
82 + ^
81 + ^
80 + ^
79 + ^
78 + ^
77 + ^
76 + ^
75 + ^
74 + ^
73 + ^
72 + ^
71 + ^
70 + ^
69 + ^
68 + ^
67 + ^
66 + ^
65 + ^
64 + ^
63 + ^
62 + ^
61 + ^
60 + ^
59 + ^
58 + ^
57 + ^
56 + ^
55 + ^
54 + ^
53 + ^
52 + ^
51 + ^
50 + ^
49 + ^
48 + ^
47 + ^
46 + ^
45 + ^
44 + ^
43 + ^
-----
FU_YEAR (Since 1948)
1 3 5 7 9 11 13 15 17 19 21 23 25 27 29

```

Testing Global Null Hypothesis: BETA=0  
-2LOGL W/out:24098.1 With:23968.1 Covariates; Chi-Sq(1) 130; p=0.0001

Maximum Likelihood Estimates

Variable	DF	Parameter Estimate	Standard Error	Wald Chi-Sq	Pr > Chi-Sq	Risk (hazard) Ratio
I_MALE	1	0.583	0.051	129.3	0.0001	1.79

Risk-sets "1" "3", ... candidates for deaths in FU\_YEAR "1" "3" ...  
(each set has persons with a range of ages)

```

TIME Scale = AGE .(NOTE how delayed entry is specified)
FU_AGE | risk set (horizontal) based on deaths at a particular age
88 + >
87 + >
86 + >
85 + >
84 + >
83 + >
82 + >
81 + >
80 + >
79 + >
78 + >
77 + >
76 + >
75 + >
74 + >
73 + >
72 + >
71 + >
70 + >
69 + >
68 + >
67 + >
66 + >
65 + >
64 + >
63 + >
62 + >
61 + >
60 + >
59 + >
58 + >
57 + >
56 + >
55 + >
54 + >
53 + >
52 + >
51 + >
50 + >
49 + >
48 + >
47 + >
46 + >
45 + >
44 + >
43 + >
-----
FU_YEAR
1 3 5 7 9 11 13 15 17 19 21 23 25 27 29

```

Testing Global Null Hypothesis: BETA=0  
-2LOGL W/out:22819.8 With:22662.7 Covariates; Chi-Sq(1) 157; p=0.0001

Maximum Likelihood Estimates

Variable	DF	Parameter Estimate	Standard Error	Wald Chi-Sq	Pr > Chi-Sq	Risk (hazard) Ratio
I_MALE	1	0.643	0.051	156.3	0.0001	1.90

Risk-sets "68" "69" ... candidates for death at age 68, 69, ...  
**Mortality rates vary much more (and in more complex way) over 20 years of age, than over 20 calendar years => "t"=age**

## Topics

- Fitting proportional hazards model to data
- estimating HR by (Partial) Likelihood approach
- "Information": how sharp is curvature of LogL fn
- Estimating HR via SAS PROC PHREG / Stata / R
- Estimating  $h_0(t)$  and  $S_0(t)$  [the "corner"]
- Estimating  $S_{\underline{X}}(t)$  [ $\underline{X}$  = a specified covariate pattern]
- Split Records  
(also a way to handle time-dep. covariates)
- Estimation for stratified survival data

**Readings** [ <http://www.epi.mcgill.ca/hanley/c681/cox> ]

Clayton&Hills, Ch 30, sections 4-6

Collett Textbook, Chapter 3/4

Kleinbaum's 'Self-Learning' textbook, Chapter 3/4

Pair of expository articles by JH

## Proportional hazards model

Simplest case (1 covariate  $z$ , 2 levels/groups which we will refer to as 0 and 1)

Compared with reference individuals (group 0), who have a hazard  $h_0[t]$  at time  $t$ , those in group 1 have a hazard that is a constant times  $h_0[t]$ , i.e.

$$\frac{h_1[t]}{h_0[t]} = \text{constant}$$

{Selvin uses 'c' and Collett uses ' ' for HR, the hazard ratio}.

Equivalently, one can write

$$h_1[t] = \text{HR} \times h_0[t]$$

The hazard ratio HR will be a number between 0 and infinity. To make it easier to fit this parameter without having to constrain it within these bounds, it helps to re express HR as

$$\text{HR} = e^{\quad} \quad \{ \text{or } \ln[\text{HR}] = \quad \}$$

so that the model becomes

$$h_1[t] = e^{\quad} \times h_0[t]$$

or

$$\ln[h_{z=1}[t]] = \ln[h_{z=0}[t]] + \quad \times (z=1) .$$

One can think of the  $\ln[h_0[t]]$  as the intercept and  $z$  as the indicator variable for group in a regression. Note that the 'intercept' here is a full hazard curve over  $t$ ; Unlike the case of other regressions, here the intercept may be of interest. However we may not have enough data to estimate it well, especially if, as is often the case, it varies considerably over  $t$ , or we do not have many events.

**Relationship between S[t] for z=1 versus S[t] for z=0 ["corner"]**

(*déjà, but in less detail.. can skip this page for now*)

If  $h_1(t) = e \cdot h_0(t)$ , and if  $H[t]$  is the integrated hazard,

then the integrated (or "cumulative") hazard for  $z=1$  is

$$H_1(t) = e \cdot H_0(t),$$

so that the survival functions are

$$\begin{aligned} S_1(t) &= e^{-H_1(t)} \\ &= e^{-e \cdot H_0(t)} = [e^{-H_0(t)}]^e \\ &= [S_0(t)]^e = [S_0(t)]^{HR} \end{aligned}$$

Thus, the **log[-log]** functions should be **parallel**,

$$\begin{aligned} \log[-\log[S_1(t)]] & \\ &= \log[-\log[S_0(t)]^{HR}] \\ &= \log[HR \cdot -\log[S_0(t)]] \\ &= \log[HR] + \log[-\log[S_0(t)]] \end{aligned}$$

and **separated by the quantity log[HR]**.

Thus one can **visually estimate**  $b = \log[HR]$  from  $\log[-\log S]$  plots. If limited data, the hazard functions may be too unstable to use.

The "baseline" hazard function  $h_0(t)$  can be from some parametric family [e.g.  $h_0(t) = \text{constant}$  {negative exponential distribution of failure times}, Weibull, ...] or can be unspecified. In the latter case, the mixture of a parametric form for HR and a 'free' form for  $h_0(t)$  is why the model is called "semi-parametric".

More general case (1 covariate  $z$ , with possibly several levels or possibly continuous; or several covariates, continuous/discrete/mixed)

For short, refer to set of covariates  $\{z_1, z_2, \dots, z_k\}$  as  $\mathbf{z}$ ; without loss of generality, refer to a reference group of individuals as having  $\{z_1=0, z_2=0, \dots, z_k=0\}$  as  $\mathbf{z}=\mathbf{0}$ .

Compared with reference individuals (group with  $\mathbf{z}=\mathbf{0}$ ), who have a hazard  $h_0[t]$  at time  $t$ , those with covariate values  $\{z_1, z_2, \dots, z_k\}$  have a hazard that is some multiple times  $h_0[t]$ , where the multiple depends only on  $\mathbf{z}$  i.e.

$$\frac{h_{\mathbf{z}}[t]}{h_0[t]} = HR(\mathbf{z})$$

or

$$h_{\mathbf{z}}(t) = HR(\mathbf{z}) \times h_0(t)$$

Most often,  $HR(\mathbf{z})$  is taken as log-linear i.e. the log of  $HR(\mathbf{z})$  is taken as linear in the  $k$  parameters  $\{z_1, z_2, \dots, z_k\}$  i.e.

$$\log[HR(\mathbf{z})] = \beta_1 z_1 + \beta_2 z_2 + \dots + \beta_k z_k$$

or

$$HR(\mathbf{z}) = \exp\{\beta_1 z_1 + \beta_2 z_2 + \dots + \beta_k z_k\}.$$

Since

$$\exp\{\beta_1 z_1 + \beta_2 z_2\} = \exp\{\beta_1 z_1\} \times \exp\{\beta_2 z_2\},$$

we can rewrite model as

$$h_{\mathbf{z}}(t) = HR(z_1) \times HR(z_2) \times \dots \times HR(z_k) \times h_0(t)$$

or

$$S_{\mathbf{z}}(t) = [S_0(t)]^{e^{\mathbf{z}}} = [S_0(t)]^{HR_1 + HR_2 + \dots + HR_k}$$

where  $HR_i$  is shorthand for  $\exp\{\beta_i z_i\}$ , same for  $HR_2$  etc.

**Important** to have the "corner" covariate pattern near the actual  $\mathbf{z}$  values (so, might want to '**center**' the  $\mathbf{z}$  values first, *before* fitting).

Not precluded from using **products** or **powers** of the  $z$ 's.

### Fitting proportional hazards model: Risksets

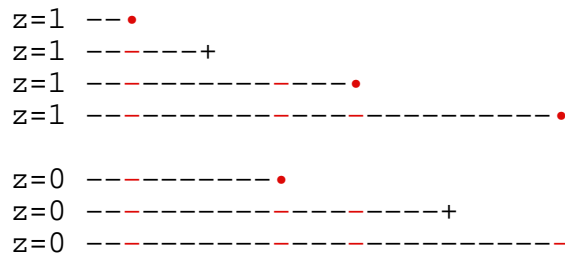
Our prime interest is in estimating the parameters of HR; we will also, as a secondary objective, estimate  $h_0(t)$ . The keys to the estimation are the Risk Sets, the collections of candidates for (individuals at risk just before) each distinct failure time (event)

Simplest case (1 covariate  $z$ , 2 levels or Tx groups which we will distinguish using indicator variable  $z=0$  and  $z=1$ ). In e.g. below, a  $\bullet$  denotes a failure (event), a  $+$  denotes a censored observation; and time runs from left to right [*note*: to estimate HR function we do not need the failure & censoring times themselves, only their *order* with respect to  $z$ ].

Raw data..( 7 individuals).



It is easier to lay them out as separate time lines [in the 'early days' before computers, some investigators would represent survival data on their patients using lines of thread along a wall].



Riskset # 1 2 3 4 5

Cox argued that since there are no failures (events) between the  $\bullet$ 's, we do not know much about the hazards in these gaps [unless we want to posit parametric form for  $h_0(t)$  or  $S_0(t)$ ]. In any case our prime interest is in HR, and so we will concentrate just on these risk sets.

### Estimating HR by (Partial) Likelihood approach

It helps to lay it out the 5 risk sets as follows (note that in the 5th riskset there is 'no contest') ...

$o=d_1$	<b>1</b>	0	<b>1</b>	<b>1</b>	-
$s_1$	3	2	1	0	-
$n_1$	4	2	2	1	
$d_0$	0	<b>1</b>	0	0	<b>1</b>
$s_0$	3	2	2	1	0
$n_0$	3	3	2	1	1

In the Maximum Likelihood method, we find that value of the HR which maximizes the likelihood of the **observed data pattern** (the **sequence** is indicated in **bold** above) The likelihood is a function of HR.

To construct the Likelihood function, we need a probability model for each table (i.e., for the outcome in each riskset) and an assumption regarding the separate tables. In the calculation of a variance for the MH statistic (log rank test) we already assumed that the 2x2 tables were realizations of hypergeometric (urn sampling) models and that the tables could be treated as if they were independent of each other. We could do the same here to set up a likelihood.

For each risk set, we ask

*"Given that the event occurred, what is the chance that it occurred to the individual it happened to, rather than to someone else in the risk set?"*

Consider a risk set where the event happened at t to a person with z=1.

If the hazard for persons with z=1 is  $HR \times h_0(t)$  and  $1 \times h_0(t)$  for those with z=0, and if in the risk set there are  $n_1$  and  $n_0$  persons respectively, then the [conditional] probability that the event happened to that particular person with z=1 out of the  $n_1$  and  $n_0$  'at risk' is

$$\frac{HR \times h_0[t]}{n_1 \times HR \times h_0[t] + n_0 \times 1 \times h_0[t]}$$

which simplifies to

$$\frac{HR}{n_1 \times HR + n_0 \times 1}$$

Conversely, in a risk set where the event happened to a person with z=0. then the [conditional] chance that the event happened to that particular person with z=0 out of the  $n_1$  and  $n_0$  'at risk' is

$$\frac{1}{n_1 \times HR + n_0 \times 1}$$

Thus, for the example above, the product of the probabilities of the observed outcome (likelihood) in each of the 4 informative risksets is

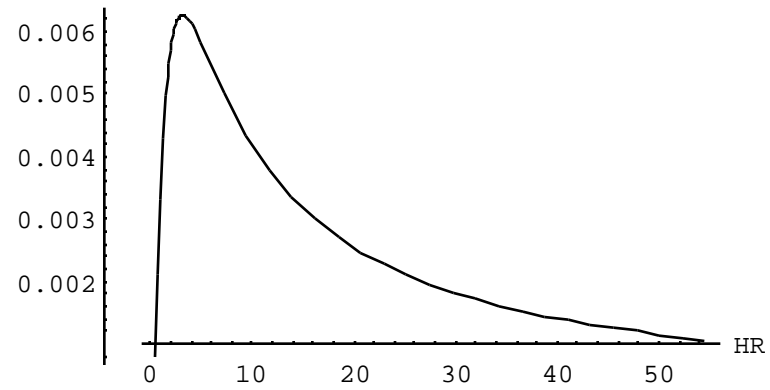
$$L = \frac{HR}{4HR+3} \times \frac{1}{2HR+3} \times \frac{HR}{2HR+2} \times \frac{HR}{HR+1}$$

This likelihood  $L(HR) = \text{prob}(\text{data} | HR)$  can be evaluated for a range of HR values in order to find the value  $\hat{HR}_{ML}$  which maximises L. e.g.

HR	1/2	1	2	4	8	16
$L \times 10^3$	1.4	3.6	5.8	6.1	4.8	3.0

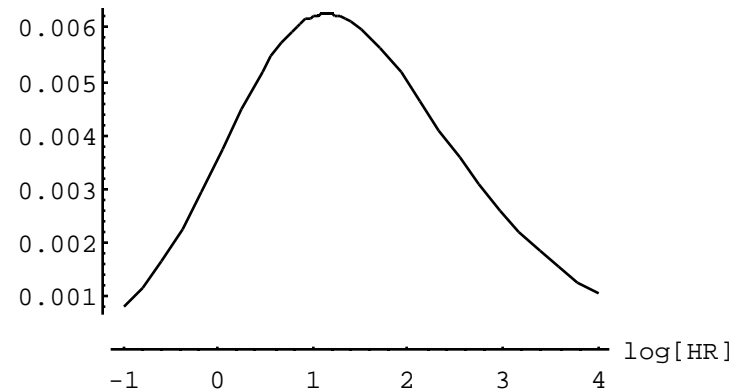
The function L & derived functions are shown graphically on next page.

Likelihood



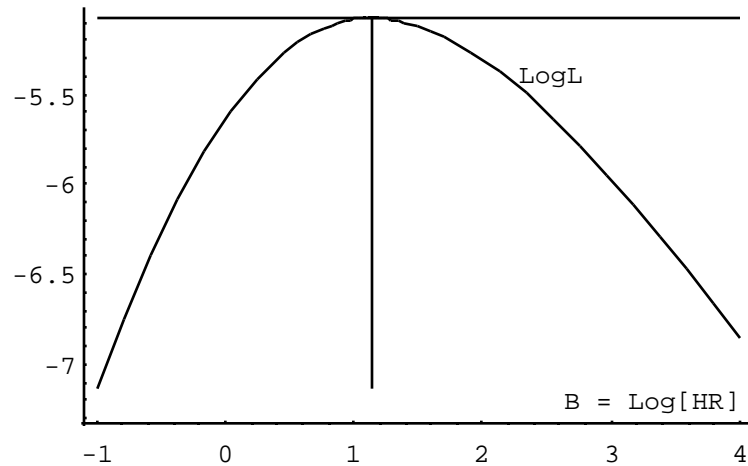
Or with the parameter  $B = \text{Log}[HR]$  ...

Likelihood

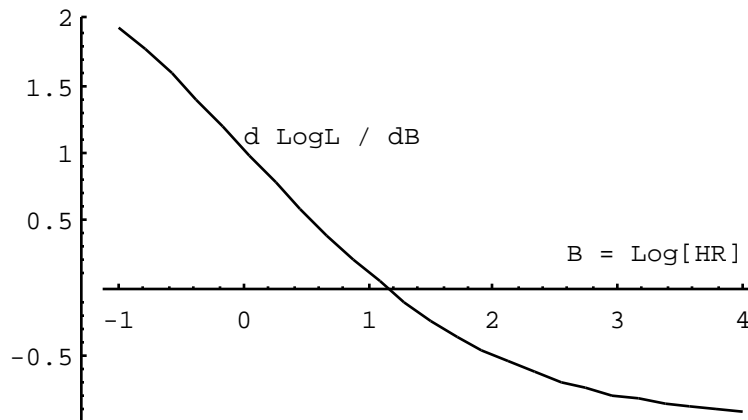




or in the log Likelihood scale...



The Derivative of the log Likelihood ...



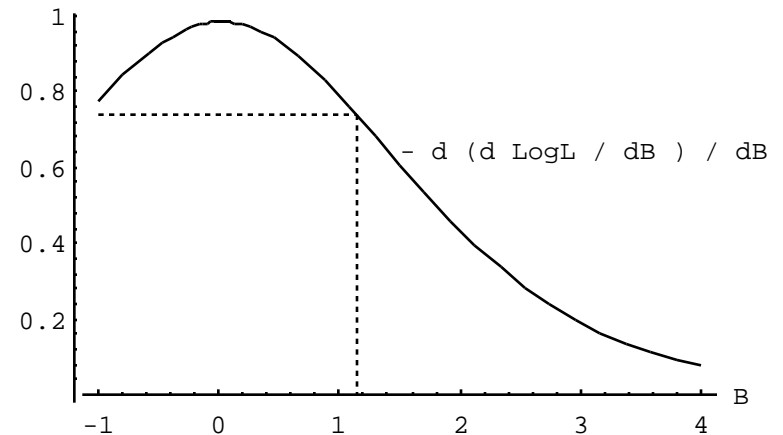
Tangent to logL curve is zero at B = 1.14 (we call this B\_hat or b);

So...  $\hat{HR}_{ML} = \exp[b] = 3.14.$

**Uncertainty / Information concerning log [HR]**

The 'sharpness' or 'flatness' of the logL(HR) curve in the vicinity of B = 1.14 gives an indication of how sensitive logL is to changes in log[HR] i.e. of how well or badly other values of log[HR] would do in producing a large likelihood. This can be measured by the 2nd derivative of logL (or if you like by the tangent to the 1st derivative curve) with respect to B. Note that the L curve increases until B = 1.14 then decreases. Thus the slope  $d\log L/dB$  goes from positive to negative over this range. ie the 2nd derivative is negative. Since we are simply interested in the curvature we use the negative of the 2nd derivative; it will be a big positive quantity when the curvature is very sharp, and a small positive quantity when the curvature is very slow.

The plot below shows that the curvature of logL is quite small (approximately 0.7412 at B = 1.14). This negative of the 2nd derivative of the log likelihood, evaluated at the ML estimate, is called the "**Information**" in the data. Its reciprocal is a good measure of the variance of the ML estimate of B.



We usually work with B =log[HR], since the sampling variability of b is more symmetric. The  $I[\ ]$  calculated at  $b = 1.14$  is approximately 0.7412, yielding  $SE[b] = (1/0.7412) = 1.16$ , yielding a 95% CI for  $HR=\exp[B]$  of {0.3 to 31}.The 4 informative risk sets provide just a small amount of information about log[HR] and our confidence in values near the ML estimate is low.

## Estimating HR via SAS PROC PHREG

```
DATA a;
INPUT event time tx ; /* Note arbitrary times */
LINES; /* only ORDER matters */
1 2 1 /* event=0 stands for censored obsn. */
0 4 1
1 6 0
1 8 1
0 10 0
1 12 1
1 14 0
;

title null model; proc phreg ; model time*event(0) = ;
Dependent Variable: TIME Number of Event & Censored Values
Censoring Variable: EVENT
Censoring Value(s): Total Event Censored %Censored
Ties Handling: BRESLOW 7 5 2 28.57
NOTE: No explanatory variables in this model. -2 LOG L = 11.27
JH: LOG L = log{1/7} x {1/5} x {1/4} x {1/2} = LOG[1/280] = -5.63
title model with tx; proc phreg data=a;
model time*event(0) = tx / RISKLIMITS;
Testing Global Null Hypothesis: BETA=0
Without With Covariates Model Chi-Square
-2 LOG L 11.27 10.15 1.12 with 1 DF (p=0.29)
ML Estimates
Parameter Standard Wald Pr > Risk* 95% CL
Variable Estimate Error Chi-Sq Chi-Sq Ratio Lower Upper
TX 1.14 1.16** 0.9685 0.33 3.14 0.32 30.6
```

\* Technically speaking, should be called Hazard Ratio; Obtained as  $\exp[1.14]$   
 \*\* See 2nd Derivative graph on left:  $SE[b] = \sqrt{\text{var}} = \sqrt{1/\text{Information}}$

## Estimating HR via Stata

```
1. input event time tx
1 2 1
0 4 1
1 6 0
1 8 1
0 10 0
1 12 1
1 14 0
end
2. stset time , failure(event)
```

7 obs., representing 5 failures in single record/single failure data  
 56 total analysis time at risk, at risk from t = 0  
 earliest observed entry t = 0 last observed exit t = 14

```
* null model
stcox, estimate
failure _d: event
analysis time _t: time
Iteration 0: log likelihood = -5.6347896
Log likelihood = -5.63 Prob > chi2 =
* model with tx.. gives beta_hats, not HR_hats
stcox tx, nohr
Iteration 0: log likelihood = -5.074435
LR chi2(1) = 1.12
Log likelihood = -5.07 Prob > chi2 = 0.2898
_t | Coef. Std. Err. z P>|z| [95% Conf. Int]
-----+-----
tx | 1.143 1.161 0.98 0.325 -1.13 3.41
-----+-----
* model with tx.. gives HR_hats, not beta_hats
stcox tx
_t | HazRatio Std. Err. z P>|z| [95% Conf. Int]
-----+-----
tx | 3.14 3.64 0.98 0.325 .32 30.55
-----+-----
```

## Estimating HR via survival package in R

```
require(survival); event=c(1,0,1,1,0,1,1);
time=c(2,4,6,8,10,12,14); tx =c(1,1,0,1,0,1,0);
fit=coxph( Surv(time, event) ~ tx); summary(fit)
coef exp(coef) se(coef) z p
tx 1.14 3.14 1.16 0.984 0.33
hr
exp(coef) exp(-coef) lower .95 upper .95
tx 3.14 0.319 0.322 30.6
Likelihood ratio test= 1.12 on 1 df, p=0.29
Wald test = 0.97 on 1 df, p=0.325
Score (logrank) test = 1.07 on 1 df, p=0.3
```

**Estimating  $h_0(t)$  and  $S_0(t)$**  [see Collett §3.8]

Once one has estimated HR, using  $\hat{HR}_{ML} = \exp[\hat{\lambda}_{ML}]$ , one can estimate the baseline hazard (and S) via a procedure similar to the Kaplan-Meier product method. One might have expected this type of non-parametric approach, since no form is specified for  $h_0(t)$ .

One uses all the events (in both groups) even though the estimate is supposed to represent individuals with  $z=0$ . The reason for this is that in a dataset with continuous covariates, there may be nobody with the specific configuration of  $z$ 's that one considers the 'reference' population.

As with all modelling and regression, we are being 'synthetic' and borrowing strength from all the data. As Collett explains, the derivation is complex, but one can get some sense of the logic from the 2-sample case where there is one event at a time [see Collett's 'particular case' following his equation 3.16]. The way JH thinks of it is to imagine a two sample situation where we were given 2 samples of death times, 1 for males and 1 for females (reference group), and told that the ratio (HR) of the death rates in the population of males and females was say 2. Would you just estimate a K-M curve for females using the data for females and call it your best estimate of the 'reference' of female S or would you try to use all the data, including the deaths from males, to estimate a better K-M curve for females?

Cox [and later Kalbfleisch and Prentice] take the 'synthetic' approach. One estimates a quantity Collett calls  $\hat{S}_0$  for each riskset. This is the 'conditional probability of survival'; estimates of these various success probabilities are multiplied together to give the unconditional probability

$\hat{S}_1$  of surviving past the time of the 1st riskset,  $\hat{S}_2$  for surviving past the 2nd, etc as in the K-M approach.

If there are multiple events per riskset, one must iteratively solve equation 3.16 for  $\hat{S}_1$ . If there is only one,  $\hat{S}_1$  can be calculated directly as

$$\hat{S}_1 = \left\{ 1 - \frac{\mathbf{HR}}{\sum \mathbf{HR}} \right\}^{1/\mathbf{HR}}$$

where  $\mathbf{HR}$  is *the calculated HR for the individual who suffered the event*, and the summation of the HR's for all the persons in the risk set.

To go back to our example of males and females and a HR of 2 for males relative to a "1" for females: suppose the risk set had 100 men and 50 women. From a hazard point of view, one can think of this as

$$100 \times 2 + 50 \times 1 = 250 \text{ "women equivalents"}$$

at risk. Now if the one event occurs to a woman, that is like saying that we had a failure of 1/250 and thus

$$\hat{S}_1 = \left\{ 1 - \frac{1}{\mathbf{HR}} \right\}^{1/1} = \left\{ 1 - \frac{1}{250} \right\}^{1/1} = \frac{249}{250}$$

If however the one event occurs to a man, that is like saying that we had a failure of 2/250 (or a success of 248/250) in two trials, so that in 1 trial of 250, we should have a success of

$$\hat{S}_1 = \left\{ 1 - \frac{2}{\mathbf{HR}} \right\}^{1/2} = \left\{ 1 - \frac{2}{250} \right\}^{1/2} = \frac{248.998}{250}$$

Obviously, in smaller n's the differences would be more dramatic. For example, with the data above, we had  $\hat{HR}_{ML} = 3.14$  which for simplicity we will round to  $\hat{HR}_{ML} = 3$ . Thus in the 5 risksets, we had the following structure

o=d <sub>1</sub>	1	0	1	1	-
s <sub>1</sub>	3	2	1	0	-
n <sub>1</sub>	4	2	2	1	
d <sub>0</sub>	0	1	0	0	1
s <sub>0</sub>	3	2	2	1	0
n <sub>0</sub>	3	3	2	1	1

or in "z=0 equivalents", replacing each person in z=1 group with HR=3

o=d <sub>1</sub>	3	0	3	3	-
s <sub>1</sub>	9	6	3	0	-
n <sub>1</sub>	12	6	6	3	
d <sub>0</sub>	0	1	0	0	1
s <sub>0</sub>	3	2	2	1	0
n <sub>0</sub>	3	3	2	1	1

or, summing and putting all individuals into "z=0 equivalents"...

d <sub>0</sub>	3	1	3	3	1
s <sub>0</sub>	12	8	5	1	0
n <sub>0</sub>	15	9	8	4	1

Then the  $\hat{S}$ 's are estimated as

$$\left\{ \frac{12}{15} \right\}^{1/3} \quad \frac{8}{9} \quad \left\{ \frac{5}{8} \right\}^{1/3} \quad \left\{ \frac{1}{4} \right\}^{1/3} \quad 0$$

or 0.93 0.89 0.86 0.63 0

yielding a Product Limit estimate of the S[ ] function for the z=0 group:

1.00	<u>0.93</u>	<u>0.83</u>	<u>0.71</u>	<u>0.44</u>	0
------	-------------	-------------	-------------	-------------	---

For the z= 1 group, the corresponding estimate is

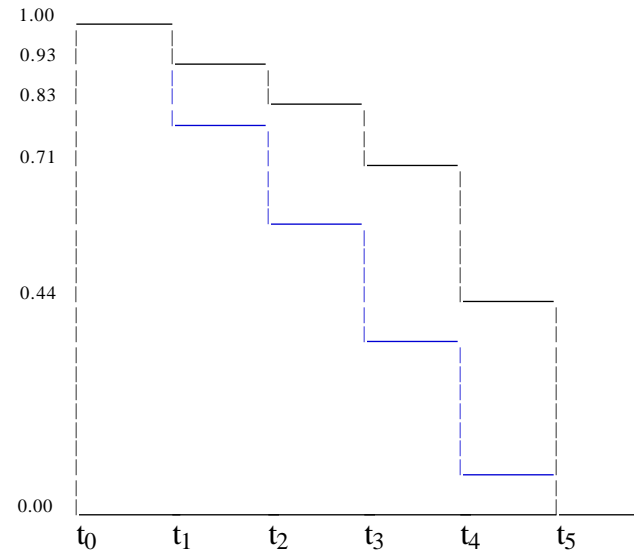
1.00	<u>0.93<sup>3</sup></u>	<u>0.83<sup>3</sup></u>	<u>0.71<sup>3</sup></u>	<u>0.44<sup>3</sup></u>	<u>0</u>
------	-------------------------	-------------------------	-------------------------	-------------------------	----------

or

1.00	<u>0.80</u>	<u>0.57</u>	<u>0.36</u>	<u>0.09</u>	0
------	-------------	-------------	-------------	-------------	---

or (roughly)

$$\hat{S}_0(t) \text{ (estimated from 5 risk sets)} \quad \text{and} \quad \hat{S}_1(t) = \{ \hat{S}_0(t) \}^3$$



Following is the estimate of S<sub>0</sub>(t) and S<sub>1</sub>(t) produced by PHREG.

First one must set up a file with the covariate patterns for which one wants survival (and other) curves.. Here there is just one covariate z, with 2 values, so there are only 2 possible covariate patterns.

```

1.          2.
DATA cov_vals;          title model with tx, and obtain curves;
INPUT          tx ;      proc phreg data=a;
LINES;          model time*event(0) = tx / RISKLIMITS;
                  0          BASELINE OUT = curves
                  1          COVARIATES = cov_vals
;                  SURVIVAL = SURVIVAL
RUN;              LOGSURV = LOGSURV
                  LOGLOGS = LOGLOGS;
    
```

**in R**

```

fit = coxph( );
plot( survfit(fit) )
baseline(fit)
    
```

\* use newdata in survfit to specify covariate patterns.

```
PROC PRINT DATA=curves ROUND; RUN;
```

TX	TIME	SURVIVAL	LOGSURV	LOGLOGS	
0	0	1.00	0.00	.	
0	2	0.93	-0.07	-2.63	*
0	6	0.83	-0.19	-1.68	**
0	8	0.71	-0.34	-1.08	***
0	12	0.45	-0.79	-0.23	****
0	14	0.00	.	.	
1	0	1.00	0.00	.	(difference)#
1	2	0.80	-0.23	-1.49	*(1.14)
1	6	0.56	-0.58	-0.54	** (1.14)
1	8	0.35	-1.06	0.06	*** (1.14)
1	12	0.08	-2.48	0.91	**** (1.14)
1	14	0.00	.	.	
0.57*	0	1.00	0.00	.	
0.57	2	0.87	-0.14	-1.98	
0.57	6	0.70	-0.36	-1.03	
0.57	8	0.52	-0.65	-0.43	
0.57	12	0.22	-1.52	0.42	
0.57	14	0.00	.	.	

\* 0.57 is the average, in the dataset, of the z values

# It is not a coincidence that there is a constant difference of 1.14 between the two FITTED  $\log[-\log[S]]$  curves: this is a **consequence** of the proportional hazards assumption..

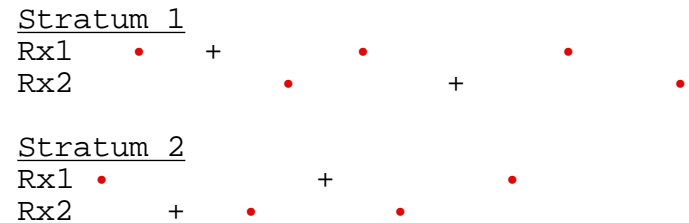
**Plotting the EMPIRICAL  $\log[-\log[S]]$  curves to see if they are reasonably parallel allows a visual check on the proportional hazards assumption.**

**Stata**

- \* store fitted survival for baseline group into new variable called **stcox tx, basesurv(s)**
- \* generate corresponding curve for tx=1 .. ^ = 'to power of' )  
**gen s\_1 = s^(exp(1.14))**
- \* graph  $-\log[S]$  i.e., cumulative hazard curves (na = Nelson-Aalen)  
**sts graph, na by(tx)**
- \* graph  $-\log[-\log[S]]$  versus time, so check if parallel  
**stspplot, by(tx)**
- \* **stcoxkm** plots Kaplan-Meier observed survival curves and compares them to the Cox predicted curves for the same variable.

**ML estimates for stratified survival data**  
(exercise.. follow *Fig 3 in part II of JH's expository article*)

Consider the following pattern of observations for two treatments where • denotes a failure (event) and + denotes a censored observation and time runs from left to right. The observations are in 2 strata.



The above calculations used the data for stratum 1.

For the second stratum

- a set up the risk sets.
- b set up the likelihood contribution from each set and the overall likelihood for the stratum (follow e.g. of stratum 1)
- c calculate the likelihood for several values of
- d draw a smooth sketch of the likelihood function (the numbers may be so tiny that you prefer to plotting the log of the likelihood function)
- e at what value (approx) of is the function a maximum?
- f calculate numerically the 1st and 2nd derivatives of the log likelihood function in the neighbourhood of  $\hat{\theta}$

**Multiply the likelihood (or add the log Likelihoods) from the 1st stratum and the likelihood from b to produce the overall likelihood (or log Likelihood) from the 2 strata combined. Then maximize the combined likelihood (or log likelihood).**

**Individuals from different strata cannot be in same riskset** (but, if strata too fine, may have uninformative risksets)