

More on the Cox PH model

I. Confidence intervals and hypothesis tests

- Two methods for confidence intervals
- Wald tests and likelihood ratio tests
- Interpretation of parameter estimates
- An example with real data from an AIDS clinical trial

II. Predicted survival under proportional hazards

III. Predicted medians and P-year survival

I. Constructing Confidence intervals and tests for the Hazard Ratio (see H & L 4.2, Collett 3.4):

Many software packages provide estimates of β , but the hazard ratio $HR = \exp(\beta)$ is usually the parameter of interest.

We can use the delta method to get standard errors for $\exp(\hat{\beta})$:

$$Var(\widehat{HR}) = Var(\exp(\hat{\beta})) = \exp(2\hat{\beta})Var(\hat{\beta})$$

Constructing confidence intervals for $\exp(\beta)$ Two options: (assuming that β is a scalar)

- I. Using $se(\exp \hat{\beta})$ obtained above via the delta method as $se(\exp \hat{\beta}) = \sqrt{[Var(\exp(\hat{\beta}))]}$, calculate the endpoints as:

$$[L, U] = [\widehat{OR} - 1.96 se(\widehat{OR}), \widehat{OR} + 1.96 se(\widehat{OR})]$$

- II. Form a confidence interval for $\hat{\beta}$, and then exponentiate the endpoints.

$$[L, U] = [e^{\hat{\beta} - 1.96 se(\hat{\beta})}, e^{\hat{\beta} + 1.96 se(\hat{\beta})}]$$

Which approach do you think would be the most preferable?

Hypothesis Tests:

For each covariate of interest, the null hypothesis is

$$H_o : HR_j = 1 \Leftrightarrow \beta_j = 0$$

A Wald test¹ of the above hypothesis is constructed as:

$$Z = \frac{\hat{\beta}_j}{se(\hat{\beta}_j)} \quad \text{or} \quad \chi^2 = \left(\frac{\hat{\beta}_j}{se(\hat{\beta}_j)} \right)^2$$

This test for $\beta_j = 0$ assumes that all other terms in the model are held fixed.

Note: if we have a factor A with a levels, then we would need to construct a χ^2 test with $(a - 1)$ df, using a test statistic based on a quadratic form:

$$\chi_{(a-1)}^2 = \widehat{\boldsymbol{\beta}}_A' Var(\widehat{\boldsymbol{\beta}}_A)^{-1} \widehat{\boldsymbol{\beta}}_A$$

where $\boldsymbol{\beta}_A = (\beta_2, \dots, \beta_a)'$ are the $(a - 1)$ coefficients corresponding to Z_2, \dots, Z_a (or Z_1, \dots, Z_{a-1} , depending on the reference group).

¹The first follows a normal distribution, and the second follows a χ^2 with 1 df. STATA gives the Z statistic, while SAS gives the χ_1^2 test statistic (the p-values are also given, and don't depend on which form, Z or χ^2 , is provided)

Likelihood Ratio Tests:

Suppose there are $(p + q)$ explanatory variables measured:

$$Z_1, \dots, Z_p, Z_{p+1}, \dots, Z_{p+q}$$

and proportional hazards are assumed.

Consider the following models:

- **Model 1:** (contains only the first p covariates)

$$\frac{\lambda_i(t, \mathbf{Z})}{\lambda_0(t)} = \exp(\beta_1 Z_1 + \dots + \beta_p Z_p)$$

- **Model 2:** (contains all $(p + q)$ covariates)

$$\frac{\lambda_i(t, \mathbf{Z})}{\lambda_0(t)} = \exp(\beta_1 Z_1 + \dots + \beta_{p+q} Z_{p+q})$$

These are *nested* models. For such nested models, we can construct a **likelihood ratio** test of

$$H_0 : \beta_{p+1} = \dots = \beta_{p+q} = 0$$

as:

$$\chi_{LR}^2 = -2 [\log(\hat{L}(1)) - \log(\hat{L}(2))]$$

Under H_o , this test statistic is approximately distributed as χ^2 with q df.

Some examples using the Stata stcox command:

Model 1:

```
. use mac

. stset mactime macstat

. stcox karnof rif clari, nohr

        failure _d:  macstat
analysis time _t:  mactime
```

Cox regression -- Breslow method for ties

No. of subjects =	1151	Number of obs =	1151
No. of failures =	121		
Time at risk =	489509		
		LR chi2(3) =	32.01
Log likelihood =	-754.52813	Prob > chi2 =	0.0000

_t						
_d	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
-----+-----						
karnof	-.0448295	.0106355	-4.215	0.000	-.0656747	-.0239843
rif	.8723819	.2369497	3.682	0.000	.4079691	1.336795
clari	.2760775	.2580215	1.070	0.285	-.2296354	.7817903

Model 2:

```
. stcox karnof rif clari cd4, nohr

        failure _d:  macstat
analysis time _t:  mactime
```

Cox regression -- Breslow method for ties

No. of subjects =	1151	Number of obs =	1151
No. of failures =	121		
Time at risk =	489509		
		LR chi2(4) =	63.74
Log likelihood =	-738.66225	Prob > chi2 =	0.0000

_t						
_d	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
-----+-----						
karnof	-.0368538	.0106652	-3.456	0.001	-.0577572	-.0159503
rif	.880338	.2371111	3.713	0.000	.4156089	1.345067
clari	.2530205	.2583478	0.979	0.327	-.253332	.7593729
cd4	-.0183553	.0036839	-4.983	0.000	-.0255757	-.0111349

Notes:

- If we omit the **nohr** option, we will get the estimated hazard ratio along with 95% confidence intervals using Method II (i.e., forming a CI for the log HR (beta), and then exponentiating the bounds)

_t						
_d	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
karnof	.9638171	.0102793	-3.456	0.001	.9438791	.9841762
rif	2.411715	.5718442	3.713	0.000	1.515293	3.838444
clari	1.28791	.3327287	0.979	0.327	.7762102	2.136936
cd4	.9818121	.0036169	-4.983	0.000	.9747486	.9889269

- We can also compute the hazard ratio ourselves, by exponentiating the coefficients:

$$HR_{cd4} = \exp(-0.01835) = 0.98$$

Why is this HR so close to 1, and yet still highly significant?

What is the interpretation of this HR?

- The likelihood ratio test for the effect of CD4 is twice the difference in minus log-likelihoods between the two models:

$$\chi^2_{LR} = 2 * (754.533 - (738.66)) = 31.74$$

How does this test statistic compare to the Wald χ^2 test?

- In the mac study, there were three treatment arms (rif, clari, and the rif+clari combination). Because we have only included the **rif** and **clari** effects in the model, the combination therapy is the “reference” group.
- We can conduct an overall test of treatment using the **test** command in Stata:

```
. test rif clari

( 1)  rif = 0.0
( 2)  clari = 0.0

             chi2( 2) =    17.01
             Prob > chi2 =    0.0002
```

for a 2 df Wald chi-square test of whether both treatment coefficients are equal to 0. This **test** command can be used to conduct an overall test for any number of effects.

- The **test** command can also be used to test whether there is a difference between the **rif** and **clari** treatment arms:

```
. test rif=clari

( 1)  rif - clari = 0.0

             chi2( 1) =    8.76
             Prob > chi2 =    0.0031
```

Some examples using SAS PROC PHREG

```
proc phreg data=alloi;
  model dthtime*dthstat(0)=mlogrna cd4grp1 cd4grp2 combther
    / risklimits;
  cd4level: test cd4grp1, cd4grp2;
  title1 'Proportional hazards regression model for time to Death';
  title2 'Baseline viral load and CD4 predictors';

proc phreg data=alloi;
  model dthtime*dthstat(0)=mlogrna cd4grp1 cd4grp2 combther decrs8 incrs8
    / risklimits;
  cd4level: test cd4grp1, cd4grp2;
  wk8resp: test decrs8, incrs8;
```

Notes:

- The “risklimits” option on the model statement provides 95% confidence intervals using Method II from page 2. (i.e., forming a CI for the log HR (beta), and then exponentiating the bounds)

- The “test” statement has the following form:

Label: test varname1, varname2, ..., varnamek;

for a k df Wald chi-square test of whether the k coefficients are all equal to 0.

- We can use the same approach described by Freedman to assess the effects of intermediate endpoints (incrs8, decrs8) on the treatment effect (i.e., assess their use as surrogate markers). The percentage of treatment effect explained, γ , is estimated by:

$$\hat{\gamma} = 1 - \frac{\hat{\beta}_{trt,M2}}{\hat{\beta}_{trt,M1}}$$

where M1 is the model without the intermediate endpoint and M2 is the model with the marker.

OUTPUT FROM PROC PHREG (Model 1)

Proportional hazards regression model for time to Death
Baseline viral load and CD4 predictors

Data Set: WORK.ALLOI
Dependent Variable: DTHTIME Time to death (days)
Censoring Variable: DTHSTAT Death status (1=died,0=censored)
Censoring Value(s): 0
Ties Handling: BRESLOW

Summary of the Number of
Event and Censored Values

Total	Event	Censored	Percent Censored
690	89	601	87.10

Testing Global Null Hypothesis: BETA=0

Criterion	Without Covariates	With Covariates	Model Chi-Square
-2 LOG L	1072.543	924.167	148.376 with 4 DF (p=0.0001)
Score	.	.	189.702 with 4 DF (p=0.0001)
Wald	.	.	127.844 with 4 DF (p=0.0001)

Analysis of Maximum Likelihood Estimates

Variable	DF	Parameter Estimate	Standard Error	Wald Chi-Square	Pr > Chi-Square
MLOGRNA	1	0.833237	0.17808	21.89295	0.0001
CD4GRP1	1	2.364612	0.32436	53.14442	0.0001
CD4GRP2	1	1.171137	0.34434	11.56739	0.0007
COMBTHER	1	-0.497161	0.24389	4.15520	0.0415

OUTPUT FROM PROC PHREG, continued

Output from “risklimits” and “test” statements

Analysis of Maximum Likelihood Estimates				
Conditional Risk Ratio and 95% Confidence Limits				
Variable	Risk Ratio	Lower	Upper	Label
MLOGRNA	2.301	1.623	3.262	log baseline rna (roche assay)
CD4GRP1	10.640	5.634	20.093	CD4<=100
CD4GRP2	3.226	1.643	6.335	100<CD4<=200
COMBTHER	0.608	0.377	0.981	Combination therapy with AZT/ddI/ddC/Nvp

Linear Hypotheses Testing			
Label	Wald Chi-Square	DF	Pr > Chi-Square
CD4LEVEL	55.0794	2	0.0001

OUTPUT FROM PROC PHREG, (Model 2)

Proportional hazards regression model for time to Death					
Baseline viral load and CD4 predictors					
Data Set: WORK.ALLOI					
Dependent Variable: DTHTIME		Time to death (days)			
Censoring Variable: DTHSTAT		Death status (1=died,0=censored)			
Censoring Value(s): 0					
Ties Handling: BRESLOW					
Summary of the Number of Event and Censored Values					
	Total	Event	Censored	Percent Censored	
	690	89	601	87.10	
Testing Global Null Hypothesis: BETA=0					
Criterion	Without Covariates	With Covariates	Model Chi-Square		
-2 LOG L	1072.543	912.009	160.535 with 6 DF (p=0.0001)		
Score	.	.	198.537 with 6 DF (p=0.0001)		
Wald	.	.	132.091 with 6 DF (p=0.0001)		
Analysis of Maximum Likelihood Estimates					
Variable	DF	Parameter Estimate	Standard Error	Wald Chi-Square	Pr > Chi-Square
OGRNA	1	0.893838	0.18062	24.48880	0.0001
AGRP1	1	2.023005	0.33594	36.26461	0.0001
AGRP2	1	1.001046	0.34907	8.22394	0.0041
MBTHER	1	-0.456506	0.24687	3.41950	0.0644
CRS8	1	-0.410919	0.26383	2.42579	0.1194
CRS8	1	-0.834101	0.32884	6.43367	0.0112

OUTPUT FROM PROC PHREG, continued

Output from “risklimits” and “test” statements

Analysis of Maximum Likelihood Estimates			
Conditional Risk Ratio and 95% Confidence Limits			
Variable	Risk Ratio	Lower	Upper Label
MLOGRNA	2.444	1.716	3.483 log baseline rna (roche assay)
CD4GRP1	7.561	3.914	14.606 CD4<=100
CD4GRP2	2.721	1.373	5.394 100<CD4<=200
COMBTHER	0.633	0.390	1.028 Combination therapy with AZT/ddI/ddC/Nvp
DECRS8	0.663	0.395	1.112 Decrease>=0.5 log rna at week 8?
INCRS8	0.434	0.228	0.827 Increase>=50 CD4 cells, week 8?

Linear Hypotheses Testing			
Label	Wald Chi-Square	DF	Pr > Chi-Square
CD4LEVEL	37.6833	2	0.0001
WK8RESP	10.4312	2	0.0054

The percentage of treatment effect explained by including the RNA and CD4 response to treatment by Week 8 is:

$$\hat{\gamma} = 1 - \frac{-0.456}{-0.497} \approx 0.08$$

or 8%. The percentage of treatment effect on time to first opportunistic infection or death is much higher (about 24%).

II. Predicted Survival using PH

The Cox PH model says that $\lambda_i(t, \mathbf{Z}) = \lambda_0(t) \exp(\boldsymbol{\beta}\mathbf{Z})$. What does this imply about the survival function, $S_z(t)$, for the i -th individual with covariates \mathbf{Z}_i ?

For the baseline (reference) group, we have:

$$S_0(t) = e^{-\int_0^t \lambda_0(u) du} = e^{-\Lambda_0(t)}$$

This is by definition of a survival function (see intro notes).

For the i -th patient with covariates \mathbf{Z}_i , we have:

$$\begin{aligned}
 S_i(t) &= e^{-\int_0^t \lambda_i(u) du} = e^{-\Lambda_i(t)} \\
 &= e^{-\int_0^t \lambda_0(u) \exp(\boldsymbol{\beta}\mathbf{Z}_i) du} \\
 &= e^{-\exp(\boldsymbol{\beta}\mathbf{Z}_i) \int_0^t \lambda_0(u) du} \\
 &= \left[e^{-\int_0^t \lambda_0(u) du} \right]^{\exp(\boldsymbol{\beta}\mathbf{Z}_i)} \\
 &= [S_0(t)]^{\exp(\boldsymbol{\beta}\mathbf{Z}_i)}
 \end{aligned}$$

(This uses the mathematical relationship $[e^b]^a = e^{ab}$)

Say we are interested in the survival pattern for single males in the nursing home study. Based on the previous formula, if we had an estimate for the survival function in the reference group, i.e., $\hat{S}_0(t)$, we could get estimates of the survival function for any set of covariates \mathbf{Z}_i .

How can we estimate the survival function, $S_0(t)$?

We could use the KM estimator, but there are a few disadvantages of that approach:

- It would only use the survival times for observations contained in the reference group, and not all the rest of the survival times.
- It would tend to be somewhat choppy, since it would reflect the smaller sample size of the reference group.
- It's possible that there are no subjects in the dataset who are in the "reference" group (ex. say covariates are age and sex; there is no one of age=0 in our dataset).

Instead, we will use a baseline hazard estimator which takes advantage of the proportional hazards assumption to get a smoother estimate.

$$\hat{S}_i(t) = [\hat{S}_0(t)]^{\exp(\hat{\beta}\mathbf{Z}_i)}$$

Using the above formula, we substitute $\hat{\beta}$ based on fitting the Cox PH model, and calculate $\hat{S}_0(t)$ by one of the following approaches:

- Breslow estimator (Stata)
- Kalbfleisch/Prentice estimator (SAS)

(1) **Breslow Estimator:**

$$\hat{S}_0(t) = \exp^{-\hat{\Lambda}_0(t)}$$

where $\hat{\Lambda}_0(t)$ is the estimated cumulative baseline hazard:

$$\hat{\Lambda}(t) = \sum_{j:\tau_j < t} \left(\frac{d_j}{\sum_{k \in \mathcal{R}(\tau_j)} \exp(\beta_1 Z_{1k} + \dots \beta_p Z_{pk})} \right)$$

(2) **Kalbfleisch/Prentice Estimator**

$$\hat{S}_0(t) = \prod_{j:\tau_j < t} \hat{\alpha}_j$$

where $\hat{\alpha}_j, j = 1, \dots, d$ are the MLE's obtained by assuming that $S(t; Z)$ satisfies

$$S(t; Z) = [S_0(t)]^{e^{\beta Z}} = \left[\prod_{j:\tau_j < t} \alpha_j \right]^{e^{\beta Z}} = \prod_{j:\tau_j < t} \alpha_j^{e^{\beta Z}}$$

Breslow Estimator: further motivation

The Breslow estimator is based on extending the concept of the Nelson-Aalen estimator to the proportional hazards model.

Recall that for a single sample with no covariates, the **Nelson-Aalen Estimator** of the cumulative hazard is:

$$\hat{\Lambda}(t) = \sum_{j:\tau_j < t} \frac{d_j}{r_j}$$

where d_j and r_j are the number of deaths and the number at risk, respectively, at the j -th death time.

When there are covariates and assuming the PH model above, one can generalize this to estimate the cumulative baseline hazard by adjusting the denominator:

$$\hat{\Lambda}(t) = \sum_{j:\tau_j < t} \left(\frac{d_j}{\sum_{k \in \mathcal{R}(\tau_j)} \exp(\beta_1 Z_{1k} + \dots \beta_p Z_{pk})} \right)$$

Heuristic: The expected number of failures in $(t, t + \delta t)$ is

$$d_j \approx \delta t \times \sum_{k \in \mathcal{R}(t)} \lambda_0(t) \exp(z_k \hat{\beta})$$

Hence,

$$\delta t \times \lambda_0(t_j) \approx \frac{d_j}{\sum_{k \in \mathcal{R}(t)} \exp(z_k \hat{\beta})}$$

Kalbfleisch/Prentice Estimator: further motivation

This method is analogous to the Kaplan-Meier Estimator. Consider a discrete time model with hazard $(1 - \alpha_j)$ at the j -th observed death time.

(Note: we use $\alpha_j = (1 - \lambda_j)$ to simplify the algebra!)

Thus, for someone with $z=0$, the survivorship function is

$$S_0(t) = \prod_{j:\tau_j < t} \alpha_j$$

and for someone with $Z \neq 0$, it is:

$$S(t; Z) = S_0(t)^{e^{\beta Z}} = \left[\prod_{j:\tau_j < t} \alpha_j \right]^{e^{\beta Z}} = \prod_{j:\tau_j < t} \alpha_j^{e^{\beta Z}}$$

The likelihood contributions under this model are:

- for someone censored at t : $S(t; Z)$
- for someone who fails at t_j :

$$S(t_{(j-1)}; Z) - S(t_j; Z) = \left[\prod_{k < j} \alpha_k \right]^{e^{\beta z}} [1 - \alpha_j^{e^{\beta Z}}]$$

The solution for α_j satisfies:

$$\sum_{k \in \mathcal{D}_j} \frac{\exp(Z_k \beta)}{1 - \alpha_j^{\exp(Z_k \beta)}} = \sum_{k \in \mathcal{R}_j} \exp(Z_k \beta)$$

(Note what happens when $Z = 0$)

Obtaining $\hat{S}_0(t)$ from software packages

- Stata provides the Breslow estimator of $S_0(t; Z)$, but not predicted survivals at specified covariate values..... you have to construct these yourself
- SAS uses the Kalbfleisch/Prentice estimator of the baseline hazard, and can provide estimates of survival at arbitrary values of the covariates with a little bit of programming.

In practice, they are **incredibly** close! (see Fleming and Harrington 1984, *Communications in Statistics*)

Using Stata to Predict Survival

The Stata command **basesurv** calculates the predicted survival values for the reference group, i.e., those subjects with all covariates=0.

(1) **Baseline Survival:**

To obtain the estimated baseline survival $\hat{S}_0(t)$, follow the example below (for the nursing home data):

```
. use nurshome

. stset los fail

. stcox married health, basesurv(prsurv)

. sort los

. list los prsurv
```

Estimating the Baseline Survival with Stata

	los	prsurv
1.	1	.99252899
2.	1	.99252899
3.	1	.99252899
4.	1	.99252899
5.	1	.99252899
.		
.		
.		
22.	1	.99252899
23.	2	.98671824
24.	2	.98671824
25.	2	.98671824
26.	2	.98671824
27.	2	.98671824
28.	2	.98671824
29.	2	.98671824
30.	2	.98671824
31.	2	.98671824
32.	2	.98671824
33.	2	.98671824
34.	2	.98671824
35.	2	.98671824
36.	2	.98671824
37.	2	.98671824
38.	2	.98671824
39.	2	.98671824
40.	3	.98362595
41.	3	.98362595
.		
.		
.		

Stata creates a predicted baseline survival estimate for every observed event time in the dataset, even if there are duplicates.

(2) Predicted Survival for Subgroups

To obtain the estimated survival $\hat{S}_i(t)$ for any other subgroup (i.e., not the reference or baseline group), follow the Stata commands below:

```
. predict betaz, xb

. gen newterm=exp(betaz)

. gen predsuv=prsurv^newterm

. sort married health los

. list married health los predsuv
```

Predicting Survival for Subgroups with Stata

	married	health	los	predsuv
1.	0	2	1	.9896138
8.	0	2	2	.981557
11.	0	2	3	.9772769
13.	0	2	4	.9691724
16.	0	2	5	.9586483
.....				
300.	0	3	1	.9877566
302.	0	3	2	.9782748
304.	0	3	3	.9732435
305.	0	3	4	.9637272
312.	0	3	5	.9513916
.....				
768.	0	4	1	.9855696
777.	0	4	2	.9744162
779.	0	4	3	.9685058
781.	0	4	4	.9573418
785.	0	4	5	.9428996
.....				
.				
.				
.				
1468.	1	4	1	.9806339
1469.	1	4	2	.9657326
1472.	1	4	3	.9578599
1473.	1	4	5	.9239448
.....				
1559.	1	5	1	.9771894
1560.	1	5	2	.9596928
1562.	1	5	3	.9504684
1564.	1	5	4	.9331349

Using SAS to Predict Survival

The SAS command BASELINE calculates the predicted survival values at the event times for a given set of covariate values.

- (1) To get the estimated baseline survival $\hat{S}_0(t)$, create a dataset with 0's for values of all covariates in the model
- (2) To get the estimated survival $\hat{S}_i(t)$ for any other subgroup (i.e., not the reference or baseline group), create a data set which inputs the baseline values of the covariates for the subgroup of interest.

For either case, we then supply the corresponding dataset name to the BASELINE command under PROC PHREG.

By giving the input dataset several lines, each corresponding to a different combination of covariate values, we can compute predicted survival values for more than one group at once.

(1) Baseline Survival Estimate

(note that the baseline survival function does not correspond to any observations in our sample, since health status values range from 2-5)

```
*** Estimating Baseline Survival Function under PH;
data inrisks;
    input married health;
    cards;
0 0
;

proc phreg data=pop out=survres;
    model los*fail(0)=married health;
    baseline covariates=inrisks out=outph survival=ps/nomean;

proc print data=outph;
title1 'Nursinghome data: Baseline Survival Estimate';
```

Estimating the Baseline Survival with SAS

Nursinghome data: Baseline Survival Estimate

OBS	MARRIED	HEALTH	LOS	PS
1	0	0	0	1.00000
2	0	0	1	0.99253
3	0	0	2	0.98672
4	0	0	3	0.98363
5	0	0	4	0.97776
6	0	0	5	0.97012
7	0	0	6	0.96488
8	0	0	7	0.95856
9	0	0	8	0.95361
10	0	0	9	0.94793
11	0	0	10	0.94365
12	0	0	11	0.93792
13	0	0	12	0.93323
14	0	0	13	0.92706
15	0	0	14	0.92049
16	0	0	15	0.91461
17	0	0	16	0.91017
18	0	0	17	0.90534
19	0	0	18	0.90048
20	0	0	19	0.89635
21	0	0	20	0.89220
22	0	0	21	0.88727
23	0	0	22	0.88270

(2) Predicted Survival Estimate for Subgroup

The following SAS commands will generate the predicted survival probability for each combination of covariates, at every observed event time in the dataset.

```
*** Estimating Baseline Survival Function under PH;
data inrisks;
    input married health;
    cards;
0 2
0 5
1 2
1 5
;

proc phreg data=pop out=survres;
    model los*fail(0)=married health;
    baseline covariates=inrisks out=outph survival=ps/nomean;

proc print data=outph;
title1 'Nursinghome data: predicted survival by subgroup';
```

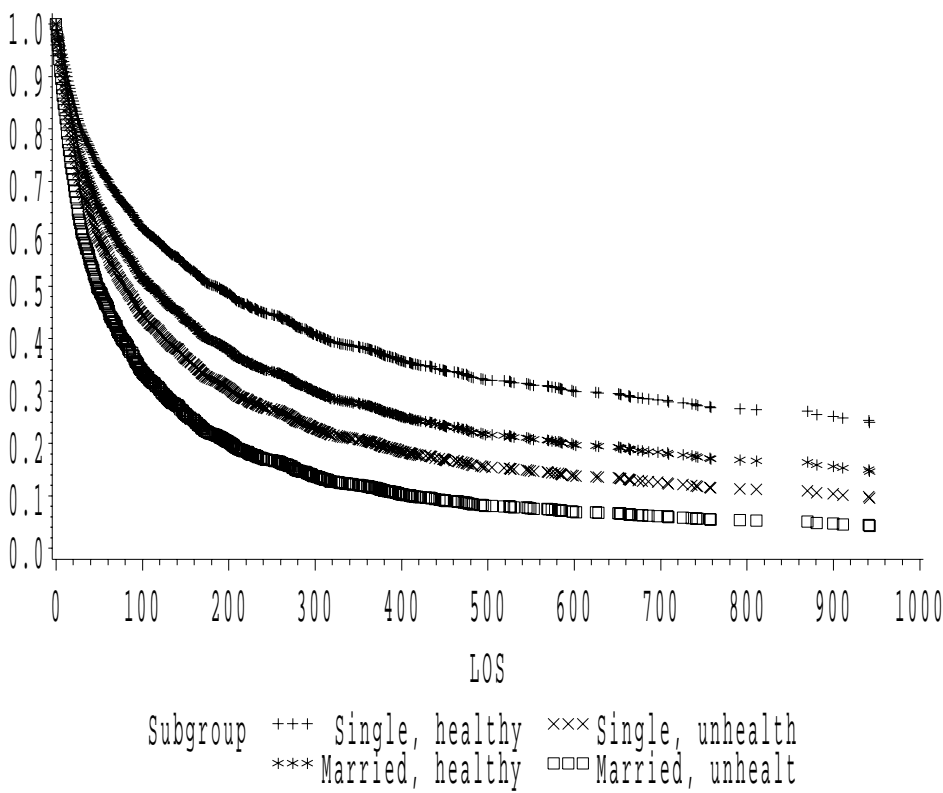
Survival Estimates by Marital and Health Status

Nursinghome data: Predicted Survival by Subgroup

OBS	MARRIED	HEALTH	LOS	PS
1	0	2	0	1.00000
2	0	2	1	0.98961
3	0	2	2	0.98156
4	0	2	3	0.97728
.....				
171	0	2	184	0.50104
172	0	2	185	0.49984
.....				
396	0	5	0	1.00000
397	0	5	1	0.98300
398	0	5	2	0.96988
399	0	5	3	0.96295
.....				
474	0	5	78	0.50268
475	0	5	80	0.49991
.....				
791	1	2	0	1.00000
792	1	2	1	0.98605
793	1	2	2	0.97527
794	1	2	3	0.96955
.....				
897	1	2	108	0.50114
898	1	2	109	0.49986
.....				
1186	1	5	0	1.00000
1187	1	5	1	0.97719
1188	1	5	2	0.95969
1189	1	5	3	0.95047
.....				
1233	1	5	47	0.50519
1234	1	5	48	0.49875

We can get a visual picture of what the proportional hazards assumption implies by looking at these four subgroups

Nursinghome data: Predicted Survival by Subgroup



III. Predicted medians and P-year survival

Predicted Medians

Suppose we want to find the predicted median survival for an individual with a specified combination of covariates (e.g., a single person with health status 5).

Three possible approaches:

- (1) Calculate the median from the subset of individuals with the specified covariate combination (using KM approach)
- (2) Generate predicted survival curves for each combination of covariates, and obtain the medians directly

OBS	MARRIED	HEALTH	LOS	PREDSURV
171	0	2	184	0.50104
172	0	2	185	0.49984
474	0	5	78	0.50268
475	0	5	80	0.49991
897	1	2	108	0.50114
898	1	2	109	0.49986
1233	1	5	47	0.50519
1234	1	5	48	0.49875

Recall that previously we defined the median as the *smallest* value of t for which $\hat{S}(t) \leq 0.5$, so the medians from above would be 185, 80, 109, and 48 days for single healthy, single unhealthy, married healthy, and married unhealthy, respectively.

- (3) Generate the predicted survival curve from the estimated baseline hazard, as follows:

We want the estimated median (M) for an individual with covariates \mathbf{Z}_i . We know

$$S(M; Z) = [S_0(M)]^{e^{\beta Z_i}} = 0.5$$

Hence, M satisfies (multiplying both sides by $e^{-\beta Z_i}$):

$$S_0(M) = [0.5]^{e^{-\beta Z}}$$

Ex. Suppose we want to estimate the median survival for a single unhealthy subject from the nursing home data. The reciprocal of the hazard ratio for unhealthy (health=5) is: $e^{-0.165 \cdot 5} = 0.4373$, (where $\hat{\beta} = 0.165$ for health status)

So, we want M such that $S_0(M) = (0.5)^{0.4373} = 0.7385$

So the median for single unhealthy subject is the 73.8th percentile of the baseline group.

OBS	MARRIED	HEALTH	LOS	PREDSURV
79	0	0	78	0.74028
80	0	0	80	0.73849
81	0	0	81	0.73670

So the estimated median would still be 80 days. Note: similar logic can be followed to estimate other quantiles besides the median.

Estimating P-year survival

Suppose we want to find the P-year survival rate for an individual with a specified combination of covariates, $\hat{S}(P; \mathbf{Z}_i)$

For an individual with $\mathbf{Z}_i = 0$, the P-year survival can be obtained from the baseline survivorship function, $\hat{S}_0(P)$

For individuals with $\mathbf{Z}_i \neq 0$, it can be obtained as:

$$\hat{S}(P; \mathbf{Z}_i) = [\hat{S}_0(P)]^{e^{\widehat{\beta}\mathbf{Z}_i}}$$

Notes:

- Although I say “P-year” survival, the units of time in a particular dataset may be days, weeks, or months. The answer here will be in the same units of time as the original data.
- If $\widehat{\beta}\mathbf{Z}_i$ is positive, then the P-year survival rate for the i -th individual will be lower than for a baseline individual.

Why is this true?