Estimating and modelling relative survival using SAS

Paul W. Dickman Department of Medical Epidemiology and Biostatistics Karolinska Institutet, Stockholm, Sweden paul.dickman@ki.se

June 2009

Estimating and modelling relative survival using SAS

Downloading the SAS files

• Download all_sas_files.zip from http://www.pauldickman.com/survival/sas/

• If the files are copied to c:\coursetemp\sas\ you won't have to change any library references.

Estimating and modelling relative survival using SAS

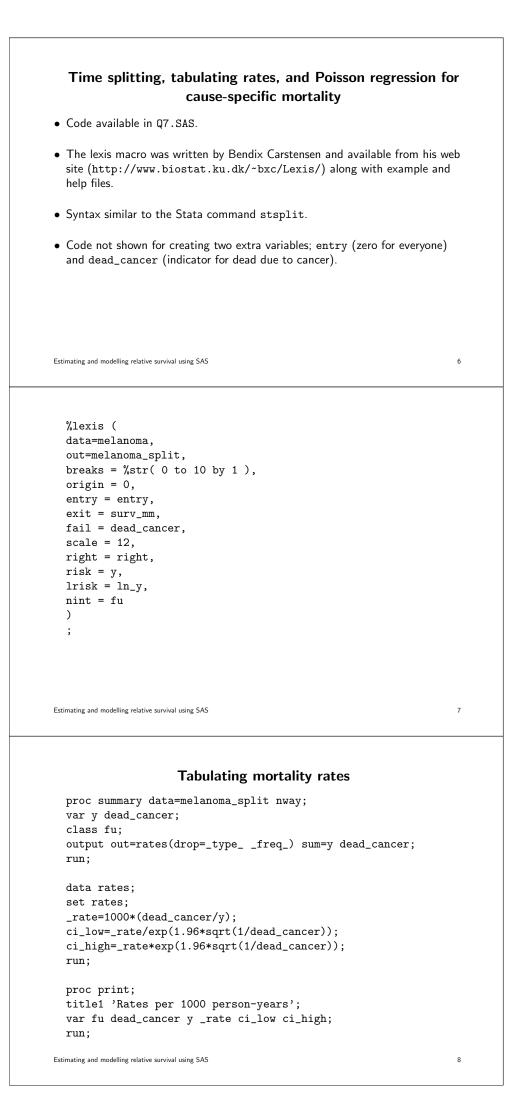
Central exercises

- 1. 'Hand calculation' of life table and Kaplan-Meier estimates
- 5. Estimating expected survival (Ederer I and II)
- 7. Poisson regression, cause-specific mortality
- 9. Cox regression, cause-specific mortality
- 14. Life table estimates of relative survival
- 18. Period analysis
- 20. Modelling relative survival

Estimating and modelling relative survival using SAS

1

SAS PROC LIFETEST	
• The LIFETEST procedure can compute nonparametric estimates of the survivor function (using either the actuarial or Kaplan-Meier method) and test the equality of survival distributions across strata (e.g. using the log-rank test).	st
• The following code produces Kaplan-Meier estimates (and a corresponding plot) of all-cause survival for the sample of 35 patients with colon carcinoma).
<pre>proc lifetest data=survival.colon_sample plots=(s) nocensplot time surv_mm*status(0,4); run;</pre>	;
• Note that we must specify the variable containing survival time, surv_mm, the status variable, status, and the codes in the status variable which define censored observations (0=alive and 4=lost to follow-up).	ie
• The plots=(s) option requests plots of the survivor function.	
Estimating and modelling relative survival using SAS	3
• The nocensplot option suppresses the plot showing censoring times.	
• The following code produces the corresponding actuarial estimates. The width=12 option specifies annual intervals for the life table.	
<pre>proc lifetest data=survival.colon_sample plots=(s) nocensplot method=act width=12; time surv_mm*status(0,4); run;</pre>	
• We could also use surv_yy as the outcome and no width statement.	
• Stratified estimates can be made using the STRATA statement.	
Estimating and modelling relative survival using SAS	4
• The following code estimates the survivor function separately for males and females and calculates several tests for differences in survival between the 2 groups (including the log-rank test).	
• A subsetting WHERE statement is used to restrict the analysis to patients diagnosed during 1985 and later (note that 'ge' is SAS notation for 'greater than or equal to (\geq) ').	
<pre>proc lifetest data=survival.colon_sample plots=(s) nocensplot time surv_mm*status(0,4); where yydx ge 85; strata sex; run;</pre>	;
• The complete SAS code for these analyses is in the file example_lifetest.sas.	
Estimating and modelling relative survival using SAS	5



Output (compare with solutions for 7e)

FU	DEAD_ CANCER	Y	_RATE	CI_LOW	CI_HIGH
1	71	5257.00	13.5058	10.7029	17.0428
2	228	4857.92	46.9337	41.2203	53.4390
3	202	4235.46	47.6926	41.5489	54.7447
4	138	3711.58	37.1809	31.4673	43.9319
5	100	3265.58	30.6224	25.1720	37.2530
6	80	2864.67	27.9265	22.4309	34.7684
7	56	2524.79	22.1800	17.0692	28.8211
8	35	2190.25	15.9799	11.4734	22.2565
9	34	1886.38	18.0240	12.8786	25.2251
10	16	1583.04	10.1071	6.1919	16.4980

Estimating and modelling relative survival using SAS

Poisson regression

ods output parameterestimates=parmest /* parameter estimates */
 type3=type3estimates; /* Type III estimates */

run;

ods output close;

• We save the parameter estimates to a data file and then exponentiate to get the rate ratios and CIs.

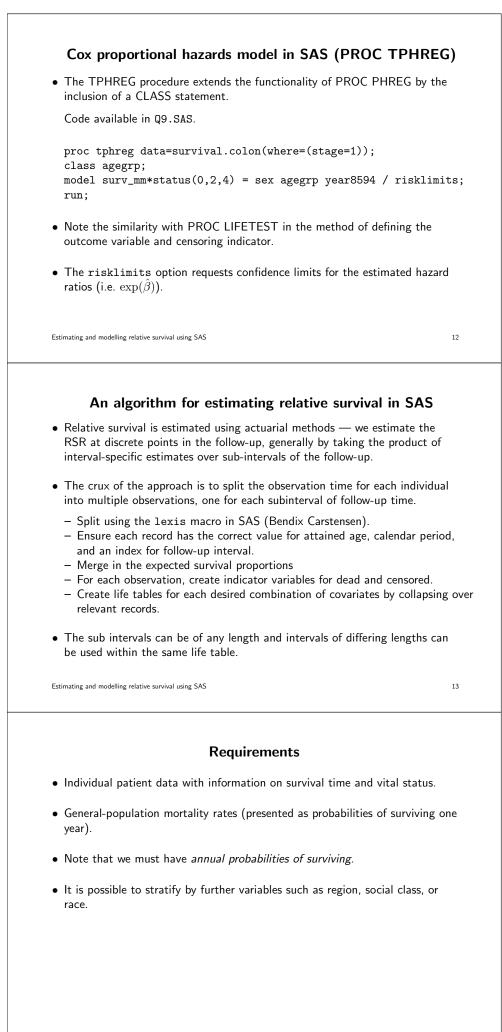
Estimating and modelling relative survival using SAS

10

9

```
data parmest;
set parmest;
if df gt 0 then do;
rr=exp(estimate);
low_rr=exp(estimate-1.96*stderr);
hi_rr=exp(estimate+1.96*stderr);
end;
run;
proc print data=parmest label noobs;
title2 'Estimates for beta and relative risks (rr=exp(beta))';
id parameter; by parameter notsorted;
var level1 estimate stderr rr low_rr hi_rr;
format estimate stderr rr low_rr hi_rr 6.3;
run;
• Results are identical to those obtained using Stata.
```

Estimating and modelling relative survival using SAS



Estimating and modelling relative survival using SAS

 Since the age, the (SEX _Y Standard according 	e patient da YEAR _AGE	be merged ita file mus	0.96429 0.99639 0.99783 0.99842 0.99882 0.99893 0.99913 0.99905 0.99920 0.99931 0.99940 0.99939 0.99920 0.99925 0.99925 0.99914
I I I I I I I I I I Since th age, the (SEX _Y • Standar accordin	1951 1951 1951 1951 1951 1951 1951 1951	2 3 4 5 6 7 8 9 10 11 12 13 14 rvival using SAS	0.99783 0.99842 0.99882 0.99893 0.99913 0.99905 0.99920 0.99931 0.99940 0.99939 0.99920 0.99925 0.99925 0.99914
 Since th age, the (SEX _Y Standar, accordin 	1951 1951 1951 1951 1951 1951 1951 1951	3 4 5 6 7 8 9 10 11 12 13 14 rvival using SAS	0.99842 0.99882 0.99893 0.99913 0.99905 0.99920 0.99931 0.99940 0.99939 0.99920 0.99925 0.99914 with the patient data file by sex, period, and
1 1 1 1 1 1 1 1 1 1 1 1 1 1	1951 1951 1951 1951 1951 1951 1951 1951	4 5 6 7 8 9 10 11 12 13 14 rvival using SAS	0.99882 0.99893 0.99913 0.99905 0.99920 0.99931 0.99940 0.99939 0.99920 0.99925 0.99914 with the patient data file by sex, period, and
I I I I I I I I I I I I I I I I Since th age, the (SEX _Y Standar accordin	1951 1951 1951 1951 1951 1951 1951 1951	5 6 7 8 9 10 11 12 13 14 rvival using SAS	0.99893 0.99913 0.99905 0.99920 0.99931 0.99940 0.99939 0.99920 0.99925 0.99914 with the patient data file by sex, period, and
I I I I I I I I I I I I I I I Since th age, the (SEX _Y Standar accordin	1951 1951 1951 1951 1951 1951 1951 1951	6 7 8 9 10 11 12 13 14 rvival using SAS	0.99913 0.99905 0.99920 0.99931 0.99940 0.99939 0.99920 0.99925 0.99914 with the patient data file by sex, period, and
I I I I I I I I I I I I I I Since th age, the (SEX _Y Standar accordin	1951 1951 1951 1951 1951 1951 1951 1951	7 8 9 10 11 12 13 14 rvival using SAS	0.99905 0.99920 0.99931 0.99940 0.99920 0.99925 0.99914
 Since th age, the (SEX _Y Standar: accordin 	1951 1951 1951 1951 1951 1951 1951 1951	8 9 10 11 12 13 14 rvival using SAS	0.99920 0.99931 0.99940 0.99939 0.99920 0.99925 0.99914 with the patient data file by sex, period, and
 Since th age, the (SEX _Y Standar accordin 	1951 1951 1951 1951 1951 1951 1951 nodelling relative sum	9 10 11 12 13 14 rvival using SAS	0.99931 0.99940 0.99939 0.99920 0.99925 0.99914 with the patient data file by sex, period, and
 1 1 1 1 Estimating and me Since the age, the (SEX _Y Standarraccordin 	1951 1951 1951 1951 1951 1951 nodelling relative sum	10 11 12 13 14 rvival using SAS	0.99940 0.99939 0.99920 0.99925 0.99914 with the patient data file by sex, period, and
1 1 1 1 Estimating and me • Since the age, the (SEX _Y • Standard accordin	1951 1951 1951 1951 nodelling relative sum	11 12 13 14 rvival using SAS be merged tta file mus	0.99939 0.99920 0.99925 0.99914 with the patient data file by sex, period, and
I I I Since th age, the (SEX _Y Standar accordin	1951 1951 1951 modelling relative sum	12 13 14 rvival using SAS be merged ita file mus	0.99920 0.99925 0.99914 with the patient data file by sex, period, and
1 1 Estimating and m Since th age, the (SEX _Y • Standar accordin	1951 1951 modelling relative sum his file will be patient da YEAR _AGE	13 14 rvival using SAS be merged ita file mus	0.99925 0.99914 with the patient data file by sex, period, and
1 Estimating and me Since th age, the (SEX _Y • Standard accordin	1951 modelling relative sum his file will h patient da YEAR _AGE	14 rvival using SAS be merged ita file mus	0.99914 with the patient data file by sex, period, and
 Since th age, the (SEX _Y Standard according 	nis file will h patient da YEAR _AGE	rvival using SAS be merged ita file mus	with the patient data file by sex, period, and
 When estand atta I have u 	ng to age at assification stimating re ained perioe	t diagnosis variables s elative surv d in order t undard that	d cohort life tables) are generally constructed s and calendar period at diagnosis (in addition t such as site, sex, stage, etc.). vival, we also have to keep track of attained ag to merge in the expected probabilities of death t those variables which are updated are prefixed
Estimating and m	odelling relative su	rvival using SAS	

• The variable SURV_MM represents survival time in months and the variable D is the event indicator.

SEX	AGE	YYDX	SURV_MM	D
Male	80	80	8.5	1
Female	77	79	31.5	1
Male	80	92	46.5	0
	Male Female	Male 80 Female 77	Male 80 80 Female 77 79	Male 80 80 8.5 Female 77 79 31.5

```
• The code to split the data into multiple records per patient is as follows
  %split (data=colon, out=colon,
  origin = 0, exit = surv_mm,
  event = d, scale = 1/12,
  cuts = %str( 0 to 15 by 1 )
  );
• The cutpoints must be specified in years (since the program assumes time at
  risk is given in years).
• Since survival time is specified in months, we specify scale = 1/12.
• Other possibilities for specifying the interval lengths are
  cuts = %str( 0,0.5,1,2,3,4,5,10,20,30 )
  or
  cuts = %str( 0 to 10 by 0.5 ).
                                                                         18
Estimating and modelling relative survival using SAS
• If we do not have a variable containing the calculated survival time, but have
  SAS date variables representing the date of diagnosis and date of exit then we
  can specify the following
  %split (data=colon, out=colon,
  entry = dx_date, exit = end_date,
  event = d, scale = 1/365.25,
  cuts = %str( 0 to 15 by 1 )
  );
Estimating and modelling relative survival using SAS
                                                                         19
• After splitting the data we obtain the following
    ID
          SEX
                   AGE
                          YYDX
                                  D
                                       W LEFT
                                                 FU
                                                           Y
                                                                LENGTH
     2
          Male
                    80
                           80
                                  1
                                      0
                                             0
                                                   1
                                                       0.70833
                                                                    1
    99
          Female
                    77
                           79
                                  0
                                      0
                                             0
                                                       1.00000
                                                  1
                                                                    1
    99
          Female
                    77
                           79
                                  0 0 1
                                                  2 1.00000
                                                                    1
    99
          Female
                    77
                           79
                                  1 0 2
                                                3 0.62500
                                                                    1
  4999
          Male
                    80
                           92
                                  0 0 0
                                                1 1.00000
                                                                    1
  4999
          Male
                    80
                           92
                                  0 0 1 2 1.00000
                                                                    1
  4999
          Male
                    80
                           92
                                  0 0
                                             2 3 1.00000
                                                                    1
                                                  4 0.87500
  4999
          Male
                    80
                           92
                                  0 1
                                             3
                                                                    1
• LEFT is the lower cutpoint of the interval (not necessarily an integer).
• FU is the index for the interval (FU=1,2,3,4,...)
• Y is the time at risk during the interval (in years).
Estimating and modelling relative survival using SAS
                                                                         20
```

a 17 != +I	ne ucatii i	ndicate	or.							
	he censori									
	that the v ne attained						e ag	ge an	d year at d	liagnosis
Fosimosing on			using CAS							2
Estimating an	d modelling relat	ive survivai	using SAS							
 We no 'upda 		o creat	te variab	les for a	attained a	age	and	caler	ndar year v	vhich are
									ted proba	
									t as the va this examp	
data	&indivi	4.								
set &	kindivid	;	(+)							
_year	=floor(ag r=floor(
run;										
Estimating an	d modelling relat	ive survival	using SAS							2
Estimating an	d modelling relat	ive survival	using SAS							2
	nd modelling relat									2
				YYDX	_YEAR	D	W	FU	Ŷ	
• This I ID 2	results in t SEX Male	the foll AGE 80	owing: _AGE 80	80	1980	1	0	1	0.70833	LENGTH 1
• This I ID	results in 1 SEX	the foll AGE	owing: _AGE		_					LENGTH
• This I ID 2 99	results in f SEX Male Female	the foll AGE 80 77	owing: _AGE 80 77	80 79	1980 1979	1 0	0 0	1 1	0.70833 1.00000	LENGTH 1 1 1
• This I ID 2 99 99	results in t SEX Male Female Female	the foll AGE 80 77 77	owing: _AGE 80 77 78	80 79 79	1980 1979 1980	1 0 0	0 0 0	1 1 2	0.70833 1.00000 1.00000	LENGTH 1 1 1
• This I ID 2 99 99 99 99 4999 4999	results in f SEX Male Female Female Female Male Male	the foll AGE 80 77 77 80 80	owing: _AGE 80 77 78 79 80 81	80 79 79 79 92 92	1980 1979 1980 1981 1992 1993	1 0 1 0 0	0 0 0 0 0	1 1 2 3 1 2	0.70833 1.00000 1.00000 0.62500 1.00000 1.00000	LENGTH 1 1 1 1 1 1
• This I ID 2 99 99 99 99 4999	results in 1 SEX Male Female Female Female Male	the foll AGE 80 77 77 77 80	owing: _AGE 80 77 78 79 80	80 79 79 79 92	1980 1979 1980 1981 1992	1 0 0 1 0	0 0 0 0	1 1 2 3 1	0.70833 1.00000 1.00000 0.62500 1.00000	LENGTH 1 1 1 1 1 1
• This I ID 2 99 99 99 4999 4999 4999 4999	results in f SEX Male Female Female Male Male Male Male Male	the foll AGE 80 77 77 80 80 80 80 80	owing: _AGE 80 77 78 79 80 81 82 83	80 79 79 92 92 92 92	1980 1979 1980 1981 1992 1993 1994 1995	1 0 1 0 0 0	0 0 0 0 0 0 1	1 2 3 1 2 3 4	0.70833 1.00000 1.00000 0.62500 1.00000 1.00000 1.00000 0.87500	LENGTH 1 1 1 1 1 1 1 1 1 1
 This r ID 2 99 99 4999 4999 4999 4999 4999 4999 4999 	results in f SEX Male Female Female Female Male Male Male Male Male	the foll AGE 80 77 77 80 80 80 80 80 80	owing: _AGE 80 77 78 79 80 81 82 83 83 erge in t	80 79 79 92 92 92 92 92 he expe	1980 1979 1980 1981 1992 1993 1994 1995 cted surv	1 0 1 0 0 0	0 0 0 0 0 1	1 2 3 1 2 3 4 porti	0.70833 1.00000 1.00000 0.62500 1.00000 1.00000 1.00000	LENGTH 1 1 1 1 1 1 1 1 the data
 This is ID 2 99 99 99 4999 4999 4999 4999 4999 The r file ar 	results in f SEX Male Female Female Female Male Male Male Male Male	the foll AGE 80 77 77 80 80 80 80 80 80	owing: _AGE 80 77 78 79 80 81 82 83 83 erge in t	80 79 79 92 92 92 92 92 he expe	1980 1979 1980 1981 1992 1993 1994 1995 cted surv	1 0 1 0 0 0	0 0 0 0 0 1	1 2 3 1 2 3 4 porti	0.70833 1.00000 1.00000 0.62500 1.00000 1.00000 0.87500 ons. Both	LENGTH 1 1 1 1 1 1 1 1 the data

```
data &individ;
length d w fu 4 y ln_y length 5;
merge &individ(in=a) &popmort(in=b);
by sex _year _age;
if a;
/* Need to adjust for interval lengths other than 1 year */
p_star=prob**length;
/* Expected number of deaths */
d_star=-log(p_star)*(y/length);
run;
```

- If the probability of surviving one year from a given date is p then the probability of surviving k years is p^k (provided k is not much larger than 1).
- k may be less than 1. For example, if k = 0.5 then we have $p^{0.5} = \sqrt{p}$.
- Here you can see why we need the population mortality data to be specified in the form of annual probabilities of death and we need the intervals to be specified in units of years.

Estimating and modelling relative survival using SAS

24

• After merging in the probabilities of death (PROB) we have the following

ID	SEX	AGE	_AGE	YYDX	_YEAR	D	W	FU	Y	LENGTH	P_STAR
2	Male	80	80	80	1980	1	0	1	0.70833	1	0.88573
99	Female	77	77	79	1979	0	0	1	1.00000	1	0.94384
99	Female	77	78	79	1980	0	0	2	1.00000	1	0.93809
99	Female	77	79	79	1981	1	0	3	0.62500	1	0.93755
4999	Male	80	80	92	1992	0	0	1	1.00000	1	0.90338
4999	Male	80	81	92	1993	0	0	2	1.00000	1	0.89360
4999	Male	80	82	92	1994	0	0	3	1.00000	1	0.88628
4999	Male	80	83	92	1995	0	1	4	0.87500	1	0.87186

Estimating and modelling relative survival using SAS

25

• The next step is to collapse the data as the first step in constructing life table estimates of survival.

- We are collapsing the data so that we have one observation for each life table interval for each combination of sex, calendar period at diagnosis, and age category at diagnosis.
- Over all of the observations in each class, we sum the number of deaths (and other quantities) and take the average of the expected survival probabilities (giving the Ederer II interval-specific expected survival).

```
Estimating and modelling relative survival using SAS
```

• If we collapse the observations over calendar period at time of diagnosis then the usual life table estimates are obtained. • If we collapse over 'updated' calendar period then we obtain 'period' estimates (i.e. Brenner method). 27 Estimating and modelling relative survival using SAS • The previous step produced 'interval-specific' estimates. Now we calculate the cumulative estimates by multiplying the interval-specific estimates. /* Construct the life table estimates */ data grouped; retain cp cp_star cr 1; set grouped; if fu=1 then do; cp=1; cp_star=1; cr=1; end: n_prime=n-w/2; /* Effective number at risk */ p=1-d/n_prime; /* Interval-specific observed survival */ r=p/p_star; /* Interval-specific relative survival */ /* Cumulative observed survival */ cp=cp*p; cp_star=cp_star*p_star; /* Cumulative expected survival */ cr=cp/cp_star; /* Cumulative relative survival */ ns=n_prime-d; /* Number of survivors */ run; Estimating and modelling relative survival using SAS 28 Sample life table output Skin melanoma diagnosed in Finland 1975-1994 (follow-up to 1995) Life table estimates of patient survival Localised 1985-94, males, aged 75+ at diagnosis Interval- Interval- Interval-Effective specific specific Cumulative number observed expected relative relative I N D W at risk survival survival survival survival 1 200 24 200.0 0.88000 0.88740 0.99166 0.99166 0 0.86596 2 176 38 21 165.5 0.77039 0.88223 0.87324 110.5 0.87558 0.93539 3 117 20 13 0.81900 0.81001 84 11 15 76.5 0.85621 0.86933 0.98491 0.79778 4 5 58 9 8 54.0 0.83333 0.86593 0.96236 0.76775 0.61644 9 3 0.84517 0.72936 0.82787 0.98830 36.5 6 41 14 0.55997 7 18 3 16.5 0.81818 0.55342 12 2 3 10.5 0.80952 0.80681 1.00336 0.55528 8 7 2 2 3 0 2 9 6.0 0.66667 0.80774 0.82535 0.45830 10 2.0 1.00000 0.81035 1.23403 0.56556 Estimating and modelling relative survival using SAS 29

Poisson regression for excess mortality

proc genmod data=&grouped(where=(fu le 5)) order=formatted; fwdlink link = log(_MEAN_-d_star); invlink ilink= exp(_XBETA_)+d_star; class fu sex age yydx; model d = fu sex yydx age / error=poisson offset=ln_y type3; format fu fu. age age. yydx yydx.; run;

Estimating and modelling relative survival using SAS

Estève et al. full likelihood approach

Estimating and modelling relative survival using SAS

31

30

Hakulinen–Tenkanen approach

```
proc genmod data=&grouped(where=(fu le 5)) order=formatted;
fwdlink link = log(-log(_mean_/p_star));
invlink ilink = exp(-exp(_xbeta_))*p_star;
class fu sex age yydx;
model ns/l_prime = fu sex yydx age / error=bin type3;
format fu fu. age age. yydx yydx.;
run;
```

Files available on the website

- The code in survival.sas produces life table estimates of relative survival stratified by sex, age, and calendar period of diagnosis. In addition, two output data sets are created (one containing grouped data and one containing individual patient data) which are used as input data sets for modelling.
- The code in models.sas estimates the excess mortality model using several different approaches.
- The code in survival_period.sas estimates survival using a period approach.
- README.PDF contains further details.

Estimating and modelling relative survival using SAS

33