

Standard errors of observed and relative survival in `strs`

Paul Dickman

June 20, 2010

1 Overview

Prior to version 1.3.0 `strs` reported incorrect standard errors (for both all-cause and relative survival) when period analysis was performed.

By default, `strs` uses the actuarial approach (Section 2) to estimate observed survival. However, if late entry is detected (as in period analysis) then observed survival is estimated by transforming the estimated cumulative hazard (Section 3).

There was a bug in `strs` prior to version 1.3.0 that affected estimated standard errors and confidence intervals when period analysis was performed. The problem was that standard errors were calculated using Greenwood's method without appropriately adjusting the effective number at risk.

From version 1.3.0, standard errors are estimated using Greenwood's method when there is no late entry (i.e., no change) but when late entry is detected they are now estimated based on a transformation of the cumulative hazard (Section 3.1).

2 Estimating observed survival using the actuarial approach

The standard actuarial estimator for the interval-specific observed survival for interval i is

$$p_i = (1 - d_i/l'_i)$$

where d_i is the number of deaths in the interval and $l'_i = 1 - w_i/2$ is the 'effective number at risk' (w_i is the number censored during the interval).

2.1 Standard error – actuarial approach

Using the method described by Greenwood (1926) [1], the standard error of the cumulative observed survival proportion up until the end of interval i is given by

$$SE({}_1p_i) = {}_1p_i \left[\sum_{j=1}^i \frac{d_j}{l'_j(l'_j - d_j)} \right]^{\frac{1}{2}}, \quad (1)$$

where l'_i is the effective number at risk at the start of interval i and d_i the number of deaths during interval i . For a single interval, Equation 1 reduces to

$$SE(p_i) = p_i \left\{ \frac{d_i}{l'_i(l'_i - d_i)} \right\}^{\frac{1}{2}} = \sqrt{p_i(1 - p_i)/l'_i},$$

which is the familiar binomial formula for the standard error of the observed interval-specific survival proportion based on l'_i trials. It can also be shown that, in the absence of censoring, Equation 1 reduces to the binomial standard error.

The relevant Stata code (extracted from `strs.ado`) is

```
/* Standard errors using Greenwood's method */
gen se_p=sqrt(p*(1-p)/n_prime)
gen se_r= se_p/p_star

/* SE of the cumulative survival */
'byby' gen se_temp=sum( d/(n_prime*(n_prime-d)) )
gen se_cp=cp*sqrt(se_temp)
drop se_temp
```

3 Estimating survival by transforming the hazard

The standard actuarial estimator for the interval-specific observed survival for interval i is

$$p_i = (1 - d_i/l'_i)$$

where d_i is the number of deaths in the interval and $l'_i = 1 - w_i/2$ is the 'effective number at risk' (w_i is the number censored during the interval). In period analysis survival times can be left truncated in addition to being right censored so fewer subjects are at risk for the full interval. As such, w_i would need to represent the number of individuals whose survival time was left truncated or right censored. However, it is possible that some survival times would be both left truncated and right censored during the interval (and hence be, on average, at risk for only one quarter of the interval) so the estimator would need to be modified to account for this.

An alternative approach is to use the relationship that the survivor function is equal to the exponential of the negative of the cumulative hazard ($S = \exp(-\Lambda)$). We can estimate the average hazard for the interval as $\lambda_i = d_i/y_i$ where d_i is the number of deaths and y_i the person-time at risk in the interval. If the hazard is assumed to be constant at this value during the interval then the cumulative hazard for the interval is $\Lambda_i = k_i \times d_i/y_i$ where k_i is the width of the interval. Our estimate of the interval-specific observed survival is therefore

$$p_i = \exp(k_i \times -d_i/y_i)$$

Since this approach assumes the hazard is constant within the interval, it is sensitive to the choice of interval length, unlike the actuarial approach which gives the same estimates of cumulative observed survival independent of the choice of intervals.

3.1 Standard error – transformation approach

The variance of the cumulative hazard [2, equation 2.2] is

$$\text{var}(\Lambda) = \sum w^2 d/n^2$$

By the delta method, the variance of the survival proportion is given by

$$\begin{aligned}\text{var}(S) &= \text{var}(\exp(-\Lambda)) \\ &= \left[\frac{d}{d\Lambda} \exp(-\Lambda)\right]^2 \text{var}(\Lambda) \\ &= S^2 \text{var}(\Lambda)\end{aligned}$$

The relevant Stata code (extracted from `strs.ado`) is

```
/* SE of P */
    gen var_Lambda=(end-start)^2*d/y^2
    gen se_p=p*sqrt(var_Lambda)

/* SE of CP */
'byby' gen var_cLambda=sum( (end-start)^2*d/y^2 )
    gen se_cp=cp*sqrt(var_cLambda)
```

4 Standard error of the relative survival ratio

The variance of the expected survival proportion is very small in comparison to the variance of the observed survival proportion and, in practice, it is assumed that the expected survival proportion is a fixed constant. The variance of the relative survival ratio (both interval-specific and cumulative) is then given by

$$\begin{aligned}\text{var}(r) &= \text{var}(p/p^*) \\ &= \text{var}(p)/(p^*)^2.\end{aligned}\tag{2}$$

The standard error (SE) of the relative survival ratio is given by $\text{SE}(r) = \text{SE}(p)/p^*$.

5 Numerical example - cohort analysis

The following example illustrates a cohort analysis. `se_p` and `se_cp` are the standard errors using Greenwood's method whereas `se_p2` and `se_cp2` are the standard errors based on transforming the cumulative hazard. As can be seen, the two approaches are effectively identical.

```
. use melanoma if stage==1, clear
. stset surv_mm, fail(status==1 2) id(id) scale(12)
. strs using popmort, br(0(1)10) mergeby(_year sex _age) ///
> by(sex) list(n d w p se_p se_p2 cp se_cp se_cp2)
```

No late entry detected - p is estimated using the actuarial method

-> sex = Male

start	n	d	w	p	se_p	se_p2	cp	se_cp	se_cp2
0	2405	82	1	0.9659	0.0037	0.0037	0.9659	0.0037	0.0037
1	2322	181	143	0.9196	0.0057	0.0057	0.8882	0.0065	0.0065
2	1998	158	136	0.9181	0.0062	0.0063	0.8155	0.0081	0.0082
3	1704	104	125	0.9366	0.0060	0.0060	0.7638	0.0091	0.0091
4	1475	88	107	0.9381	0.0064	0.0064	0.7165	0.0098	0.0098
5	1280	70	110	0.9429	0.0066	0.0066	0.6756	0.0104	0.0104
6	1100	52	95	0.9506	0.0067	0.0067	0.6422	0.0109	0.0109
7	953	32	113	0.9643	0.0062	0.0062	0.6193	0.0112	0.0112
8	808	26	95	0.9658	0.0066	0.0066	0.5981	0.0116	0.0116
9	687	25	94	0.9609	0.0077	0.0077	0.5748	0.0120	0.0120

-> sex = Female

start	n	d	w	p	se_p	se_p2	cp	se_cp	se_cp2
0	2913	69	0	0.9763	0.0028	0.0028	0.9763	0.0028	0.0028
1	2844	148	156	0.9465	0.0043	0.0043	0.9241	0.0050	0.0049
2	2540	129	160	0.9476	0.0045	0.0045	0.8756	0.0063	0.0063
3	2251	107	146	0.9509	0.0046	0.0046	0.8326	0.0072	0.0072
4	1998	78	139	0.9596	0.0045	0.0045	0.7989	0.0079	0.0079
5	1781	68	130	0.9604	0.0047	0.0047	0.7673	0.0084	0.0084
6	1583	53	123	0.9652	0.0047	0.0047	0.7405	0.0089	0.0089
7	1407	43	140	0.9678	0.0048	0.0048	0.7167	0.0093	0.0093
8	1224	42	146	0.9635	0.0055	0.0055	0.6906	0.0098	0.0098
9	1036	25	115	0.9745	0.0050	0.0051	0.6729	0.0102	0.0102

6 Numerical example - period analysis

The following example illustrates a period analysis. `se_p` and `se_cp` are the standard errors using Greenwood's method with an inappropriate value of l' whereas `se_p2` and `se_cp2` are the standard errors based on transforming the cumulative hazard. As can be seen, the Greenwood estimates are biased downwards.

```
. stset exit, enter(time mdy(1,1,1994)) exit(time mdy(12,31,1995)) ///
>   origin(dx) f(status==1 2) id(id) scale(365.24)

. strs using popmort, br(0(1)10) mergeby(_year sex _age) ///
>   by(sex) list(n d w p se_p se_p2 cp se_cp se_cp2)
```

Late entry detected for at least one observation (probably because you are performing a period analysis). The conditional survival proportion (p) is estimated by transforming the estimated hazard; `n_prime` is not meaningful and is set to missing.

```
-> sex = Male
```

start	n	d	w	p	se_p	se_p2	cp	se_cp	se_cp2
0	311	13	0	0.9442	0.0130	0.0150	0.9442	0.0112	0.0150
1	443	20	143	0.9319	0.0131	0.0147	0.8799	0.0151	0.0197
2	407	18	136	0.9341	0.0135	0.0150	0.8220	0.0176	0.0227
3	380	21	125	0.9178	0.0154	0.0172	0.7544	0.0197	0.0251
4	339	12	107	0.9482	0.0131	0.0146	0.7154	0.0207	0.0262
5	340	15	110	0.9327	0.0148	0.0168	0.6672	0.0214	0.0273
6	322	9	95	0.9591	0.0120	0.0134	0.6399	0.0217	0.0276
7	320	8	113	0.9632	0.0116	0.0128	0.6163	0.0220	0.0278
8	274	8	95	0.9569	0.0135	0.0149	0.5898	0.0223	0.0282
9	234	8	154	0.9468	0.0179	0.0183	0.5584	0.0235	0.0288

```
-> sex = Female
```

start	n	d	w	p	se_p	se_p2	cp	se_cp	se_cp2
0	337	7	0	0.9713	0.0091	0.0107	0.9713	0.0077	0.0107
1	489	14	154	0.9592	0.0098	0.0107	0.9316	0.0113	0.0146
2	483	15	160	0.9524	0.0106	0.0120	0.8873	0.0139	0.0178
3	449	23	146	0.9229	0.0138	0.0154	0.8189	0.0167	0.0214
4	412	12	139	0.9565	0.0110	0.0123	0.7833	0.0179	0.0228
5	410	8	129	0.9708	0.0091	0.0102	0.7604	0.0185	0.0235
6	423	13	122	0.9543	0.0110	0.0124	0.7257	0.0191	0.0244
7	404	2	140	0.9929	0.0046	0.0050	0.7205	0.0192	0.0245
8	354	3	146	0.9875	0.0066	0.0072	0.7115	0.0195	0.0247
9	312	3	219	0.9846	0.0086	0.0088	0.7005	0.0201	0.0251

References

- [1] Greenwood M. *The Errors of Sampling of the Survivorship Table*, vol. 33 of *Reports on Public Health and Medical Subjects*. London: Her Majesty's Stationery Office, 1926.
- [2] Breslow NE, Day NE. *Statistical Methods in Cancer Research: Volume II - The Design and Analysis of Cohort Studies*. IARC Scientific Publications No. 82. Lyon: IARC, 1987.