# Analysing competing risks data using flexible parametric survival models: what tools are available in Stata, which ones to use and when?

Sarwar Islam Mozumder (sarwar.islam@le.ac.uk)

Biostatistics Research Group, Dept. of Health Sciences, University of Leicester

2018 London Stata Conference | 6 - 7 September 2018

## Overview

1. Introduction to survival analysis & competing risks
2. Fundamental relationships
3. Modelling on the cause-specific hazards scale
   - Cause-specific Cox PH model
   - Flexible parametric models (log-cumulative cause-specific hazards)
4. Modelling directly on the cause-specific cumulative incidence
   - Fine & Gray model
   - Flexible parametric models (log-cumulative subdistribution hazards)
5. Which scale is most appropriate?
6. Summary

# Survival analysis: the fundamentals

## Key components of a survival analysis

The study of time to a particular event of interest:

- Engineering e.g. time to failure of a component
- Economics e.g. duration of unemployment
- Medical e.g. time to death (survival time) of a cancer patient

## Key components of a survival analysis

The study of time to a particular event of interest:

- Engineering e.g. time to failure of a component
- Economics e.g. duration of unemployment
- Medical e.g. time to death (survival time) of a cancer patient

Censoring:

- Right censoring: survival time > follow-up time
  - Emmigration
  - Administrative (most common)
- Non-informative censoring: Loss to follow-up is not associated with factors related to the study

## Key components of a survival analysis

The study of time to a particular event of interest:

- Engineering e.g. time to failure of a component
- Economics e.g. duration of unemployment
- Medical e.g. time to death (survival time) of a cancer patient

Censoring:

- Right censoring: survival time > follow-up time
    - Emmigration *informative?*
    - Administrative (most common) *non-informative*
- Non-informative censoring: Loss to follow-up is not associated with factors related to the study

## Key components of a survival analysis

The study of time to a particular event of interest:

- Engineering e.g. time to failure of a component
- Economics e.g. duration of unemployment
- Medical e.g. time to death (survival time) of a cancer patient

Censoring:

- Right censoring: survival time > follow-up time
  - Emmigration
  - Administrative (most common)
- Non-informative censoring: Loss to follow-up is not associated with factors related to the study
- Independent and identically distributed (i.i.d) censoring: independence between survival time and censoring time (untestable)

## Some important notation

Let $T$ be a non-negative random variable that denotes observed survival time:

**(All-cause) Survival function**

$$S(t) = P(T \geq t)$$

## Some important notation

Let $T$ be a non-negative random variable that denotes observed survival time:

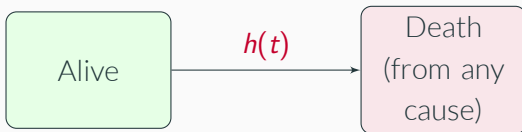**(All-cause) Survival function**

$$S(t) = P(T \geq t) = 1 - F(t)$$

**(All-cause) Cumulative incidence function (CIF)**

$$F(t) = P(T < t)$$
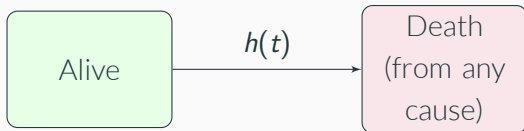
# A typical survival analysis: two-state model

# A typical survival analysis: two-state model



### (All-cause) Hazard rate, $h(t)$

Instantaneous mortality (failure) rate from any cause, given that the individual is still alive up to time $t$

## A typical survival analysis: two-state model

```
┌─────────┐            ┌──────────┐
│         │    h(t)    │  Death   │
│  Alive  │ ─────────▶ │ (from any│
│         │            │  cause)  │
└─────────┘            └──────────┘
```

**(All-cause) Hazard rate,** $h(t)$

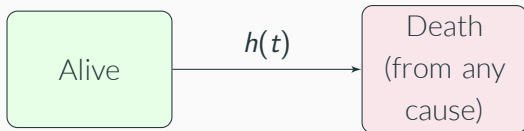$$h(t) = \lim_{\Delta t \to 0} \frac{P(t \leq T < t + \Delta t \mid T \geq t)}{\Delta t}$$

## A typical survival analysis: two-state model



**(All-cause) Hazard rate,** $h(t)$

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

**(All-cause) Survival function,** $S(t)$

$$S(t) = \exp\left(-\int_0^t h(u)\mathrm{d}u\right)$$

## Example dataset

Load public-use prostate cancer dataset:

```
. use "http://www.stata-journal.com/software/sj4-2/st0059/prostatecancer", clear
. tab status
```

|    status |  Freq. | Percent |    Cum. |
|----------:|-------:|--------:|--------:|
|    Censor |    150 |   29.64 |   29.64 |
|    Cancer |    155 |   30.63 |   60.28 |
|       CVD |    141 |   27.87 |   88.14 |
|     Other |     60 |   11.86 |  100.00 |
|     Total |    506 |  100.00 |         |

## The Kaplan-Meier estimator



```
. stset time, f(status==1,2,3) id(id) exit(time 60) scale(12)
. sts graph if agegrp == 1 & treatment == 1, ...
```

**What are competing risks?**

## Competing risks

Competing risks = when a patient dies from other causes that exclude the disease under study.

# Competing risks

Competing risks = when a patient dies from other causes that exclude the disease under study.

Non-informative censoring: Loss to follow-up is not associated with factors related to the study

## Competing risks

Competing risks = when a patient dies from other causes that exclude the disease under study.

~~Non-informative censoring: Loss to follow-up is not associated with factors related to the study~~

- Not valid under competing risks
- Death from ``competing'' causes may be due to adverse effects of treatment for disease

## Competing risks

Competing risks = when a patient dies from other causes that exclude the disease under study.

~~Non-informative censoring: Loss to follow-up is not associated with factors related to the study~~

- Not valid under competing risks
- Death from ``competing'' causes may be due to adverse effects of treatment for disease

Due to informative censoring - specialised competing risks methods are required to avoid biased estimation.

## No competing risks

## With competing risks



**Cause-specific hazard (CSH) rate,** $h_k^{cs}(t)$

Instantaneous mortality (failure) rate from cause $k$, given that the individual is still alive up to time $t$

## With competing risks



**Cause-specific hazard (CSH) rate,** $h_k^{cs}(t)$

$$h_k^{cs}(t) = \lim_{\Delta t \to 0} \frac{P(t < T \le t + \Delta t, D = k | T > t)}{\Delta t}$$

## The cause-specific CIF (transition probability)

Estimating the cause-specific CIF is of interest:

- Awkward interpretation on survival scale - what does it mean?

## The cause-specific CIF (transition probability)

Estimating the cause-specific CIF is of interest:

- Awkward interpretation on survival scale - what does it mean?
- The cause-specific survival function does not account for those who die from other competing causes before time $t$

## The cause-specific CIF (transition probability)

Estimating the cause-specific CIF is of interest:

- Awkward interpretation on survival scale - what does it mean?
- The cause-specific survival function does not account for those who die from other competing causes before time $t$
- Those who die from competing causes are removed from risk-set

## The cause-specific CIF (transition probability)

Estimating the cause-specific CIF is of interest:

- Awkward interpretation on survival scale - what does it mean?
- The cause-specific survival function does not account for those who die from other competing causes before time $t$
- Those who die from competing causes are removed from risk-set
- Better interpretation on mortality scale

### Cause-specific CIF, $F_k(t)$

Probability a patient will die from cause $D = k$ by time $t$ whilst also being at risk of dying from other competing causes of death

**Cause-specific CIF, $F_k(t)$**

Probability a patient will die from cause $D = k$ by time $t$ whilst also being at risk of dying from other competing causes of death

**Cause-specific CIF,** $F_k(t)$

$$F_k(t) = \int_0^t S(u) h_k^{cs}(u) \mathrm{d}u$$

## CSH relationship with cause-specific CIF

**Cause-specific CIF,** $F_k(t)$

$$F_k(t) = \int_0^t S(u) h_k^{cs}(u) \mathrm{d}u$$

$$S(t) = \prod_{k=1}^{K} S_k^{cs}(t) = \exp\left(-\sum_{k=1}^{K} \int_0^t h_k^{cs}(u) \mathrm{d}u\right)$$

## CSH relationship with cause-specific CIF

**Cause-specific CIF, $F_k(t)$**

$$F_k(t) = \int_0^t S(u) h_k^{cs}(u) \mathrm{d}u$$

$$S(t) = \prod_{k=1}^K S_k^{cs}(t) = \exp\left(-\sum_{k=1}^K \int_0^t h_k^{cs}(u)\mathrm{d}u\right)$$

**Note**

$$S_k^{cs}(t) = \exp\left(-\int_0^t h_k^{cs}(u)\mathrm{d}u\right) \neq 1 - F_k(t)$$

# Obtaining Aalen-Johansen (AJ) estimates of the cause-specific CIF

Non-parametric estimates of cause-specific CIFs obtained using
`stcompet`:

# Obtaining Aalen-Johansen (AJ) estimates of the cause-specific CIF

Non-parametric estimates of cause-specific CIFs obtained using `stcompet`:

```
. stset time, f(status==1) id(id) exit(time 60) scale(12)
. stcompet CIF1 = ci if agegrp == 0 & treatment == 1, compet1(2) compet2(3)
. stcompet CIF2 = ci if agegrp == 1 & treatment == 1, compet1(2) compet2(3)
```

Cancer-specific Survival

Patients aged under 75 years old and on treatment

```
. stset time, f(status==1) id(id) exit(time 60) scale(12)

. sts graph if agegrp == 0 & treatment == 1, failure ///
> addplot(line CIF1 _t if status == 1, sort connect(stepstair) ...
```

## Comparing AJ with 1 - KM estimates of the cancer-specific CIF



Cancer-specific Survival
Patients aged over 75 years old and on treatment

```
. stset time, f(status==1) id(id) exit(time 60) scale(12)

. sts graph if agegrp == 1 & treatment == 1, failure ///
> addplot(line CIF2 _t if status == 1, sort connect(stepstair)) ...
```

# Approaches for modelling (all) CSHs in Stata

## Standard approach: cause-specific Cox model

A common approach for modelling CSH function is by assuming proportional hazards (PH) using the Cox model.

**Cause-specific Cox PH model**

$$h_k^{cs}(t \mid \mathbf{x}_k) = h_{0k} \exp\left(\beta_k^{cs}\mathbf{x}_k\right)$$

$\beta_k^{cs}$: row vector of coefficients/log-CSH ratio for cause $k$

$\mathbf{x}_k$: column vector of covariates for cause $k$

$h_{0k}$: the baseline CSH function

## Standard approach: cause-specific Cox model

A common approach for modelling CSH function is by assuming proportional hazards (PH) using the Cox model.

**Cause-specific Cox PH model**

$$h_k^{cs}(t \mid \mathbf{x}_k) = h_{0k} \exp\left(\boldsymbol{\beta}_k^{cs} \mathbf{x}_k\right)$$

$\boldsymbol{\beta}_k^{cs}$: row vector of coefficients/log-CSH ratio for cause $k$

$\mathbf{x}_k$: column vector of covariates for cause $k$

$h_{0k}$: the baseline CSH function

CHR = association on the effect of a covariate on rate of dying from cause $k$

# stcox

```
. stset time, failure(status == 1) id(id) scale(12) exit(time 60)

. stcox treatment, nolog noshow

Cox regression -- Breslow method for ties

No. of subjects =          506                    Number of obs   =         506
No. of failures =          145
Time at risk    =   1457.966667
                                                  LR chi2(1)      =        6.14
Log likelihood  =   -834.85419                    Prob > chi2     =      0.0132

─────────────┬──────────────────────────────────────────────────────────────
          _t │  Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
─────────────┼──────────────────────────────────────────────────────────────
   treatment │    .6602897    .1116672    -2.45   0.014     .4740025    .9197894
─────────────┴──────────────────────────────────────────────────────────────

. predict h0_cancer, basehc

. gsort _t -_d

. by _t: replace h0_cancer = . if _n > 1

. gen h_cancer_trt0 = h0_cancer

. gen h_cancer_trt1 = h0_cancer*exp(_b[treatment])
```

## stcox

```
. stset time, failure(status == 2) id(id) scale(12) exit(time 60)

. stcox treatment, nolog noshow

Cox regression -- Breslow method for ties

No. of subjects =           506                Number of obs   =          506
No. of failures =           140
Time at risk    =   1457.966667
                                               LR chi2(1)      =         1.19
Log likelihood  =    -806.46297                Prob > chi2     =       0.2755
```

| _t | Haz. Ratio | Std. Err. | z | P>|z| | [95% Conf. Interval] |
|---|---|---|---|---|---|
| treatment | 1.20334 | .2048509 | 1.09 | 0.277 | .8619538    1.679937 |

```
. predict h0_cvd, basehc

. gsort _t -_d

. by _t: replace h0_cvd = . if _n > 1

. gen h_cvd_trt0 = h0_cvd

. gen h_cvd_trt1 = h0_cvd*exp(_b[treatment])
```

## stcox

```
. stset time, failure(status == 3) id(id) scale(12) exit(time 60)

. stcox treatment, nolog noshow

Cox regression -- Breslow method for ties

No. of subjects =          506                    Number of obs   =          506
No. of failures =           57
Time at risk    =   1457.966667
                                                  LR chi2(1)      =         2.67
Log likelihood  =    -324.95951                   Prob > chi2     =       0.1023
```

|          _t | Haz. Ratio | Std. Err. |     z | P>|z| | [95% Conf. Interval] |          |
|-------------|------------|-----------|-------|-------|----------------------|----------|
|   treatment |   .6460519 | .1745103  | -1.62 | 0.106 |             .3804893 | 1.096964 |

```
. predict h0_other, basehc

. gsort _t -_d

. by _t: replace h0_other = . if _n > 1

. gen h_other_trt0 = h0_other

. gen h_other_trt1 = h0_other*exp(_b[treatment])
```

## stcox

```
. drop if missing(h0_cancer) & missing(h0_other) & missing(h0_cvd)

. foreach i in cancer other cvd {
  2.        replace h0_`i´ = 0 if missing(h0_`i´)
  3.        replace h_`i´_trt0 = 0 if missing(h_`i´_trt0)
  4.        replace h_`i´_trt1 = 0 if missing(h_`i´_trt1)
  5. }

. sort _t

. gen S_1 = exp(sum(log(1- h_cancer_trt0 - h_other_trt0 - h_other_trt0)))

. gen S_2 = exp(sum(log(1- h_cancer_trt1 - h_other_trt1 - h_other_trt1)))

. foreach i in cancer other cvd {
  2.       gen cif_trt0_`i´ = sum(S_1[_n-1]*h_`i´_trt0)
  3.       gen cif_trt1_`i´ = sum(S_2[_n-1]*h_`i´_trt1)
  4. }

. foreach i in trt0 trt1 {
  2.       gen totcif2_`i´ = cif_`i´_cancer + cif_`i´_cvd
  3.       gen totcif3_`i´ = totcif2_`i´ + cif_`i´_other
  4. }
```

## stcox



```
. tw (rarea totcif3_trt1 totcif2_trt1 _t, sort connect(stepstair) ...) ///
> (rarea cif_trt1_cancer totcif2_trt1 _t, ...) ///
> (rarea zeros cif_trt1_cancer _t, ...), ...
```

## Remarks on Cox PH models for competing risks data

- Baseline hazard function is undefined - no risk in misspecification of underlying baseline distribution
- However, leads to difficulties in obtaining predictions to facilitate interpretation of model parameters:

## Remarks on Cox PH models for competing risks data

- Baseline hazard function is undefined - no risk in misspecification of underlying baseline distribution
- However, leads to difficulties in obtaining predictions to facilitate interpretation of model parameters:
    - Conditional and absolute measures

# Remarks on Cox PH models for competing risks data

- Baseline hazard function is undefined - no risk in misspecification of underlying baseline distribution
- However, leads to difficulties in obtaining predictions to facilitate interpretation of model parameters:
  - Conditional and absolute measures
  - Cause-specific CIF in presence of competing risks

## Remarks on Cox PH models for competing risks data

- Baseline hazard function is undefined - no risk in misspecification of underlying baseline distribution
- However, leads to difficulties in obtaining predictions to facilitate interpretation of model parameters:
  - Conditional and absolute measures
  - Cause-specific CIF in presence of competing risks
- To obtain such measures baseline hazard can be estimated non-parametrically as described by Breslow (1972)
- For a smooth function, further smoothing techniques must be applied

## Remarks on Cox PH models for competing risks data

- Baseline hazard function is undefined - no risk in misspecification of underlying baseline distribution
- However, leads to difficulties in obtaining predictions to facilitate interpretation of model parameters:
  - Conditional and absolute measures
  - Cause-specific CIF in presence of competing risks
- To obtain such measures baseline hazard can be estimated non-parametrically as described by Breslow (1972)
- For a smooth function, further smoothing techniques must be applied
- Computationally intensive methods such as bootstrapping is required for SEs/CIs

# Flexible parametric survival models (FPMs) [Royston and Parmar, 2002]

- Models and more accurately captures complex shapes of the (log-cumulative) baseline hazard function
- A generalisation of the Weibull distribution is used with restricted cubic splines (RCS) that allows for more flexibility

## Flexible parametric survival models (FPMs) [Royston and Parmar, 2002]

- Models and more accurately captures complex shapes of the (log-cumulative) baseline hazard function
- A generalisation of the Weibull distribution is used with restricted cubic splines (RCS) that allows for more flexibility

**Cause-specific log-cumulative PH FPM**

$$\ln\left(H_k^{cs}(t \mid \mathbf{x}_k)\right) = s_k(\ln t; \boldsymbol{\gamma}_k, \mathbf{m}_{0k}) + \boldsymbol{\beta}_k^{cs}\mathbf{x}_k$$

$s_k(\ln t; \boldsymbol{\gamma}_k, \mathbf{m}_{0k})$: baseline restricted cubic spline function on log-time

# Flexible parametric survival models (FPMs) [Royston and Parmar, 2002]

- Models and more accurately captures complex shapes of the (log-cumulative) baseline hazard function
- A generalisation of the Weibull distribution is used with restricted cubic splines (RCS) that allows for more flexibility
- Can also easily include time-dependent effects (TDE)

**Cause-specific log-cumulative non-PH FPM**

$$\ln\left(H_k^{cs}(t \mid \mathbf{x}_k)\right) = s_k(\ln t; \boldsymbol{\gamma}_k, \mathbf{m}_{0k}) + \boldsymbol{\beta}_k^{cs}\mathbf{x}_k + \sum_{l=1}^{E} s_k(\ln t; \boldsymbol{\alpha}_{lk}, \mathbf{m}_{lk})\mathbf{x}_{lk}$$

$s_k(\ln t; \boldsymbol{\alpha}_{lk}, \mathbf{m}_{lk})\mathbf{x}_{lk}$: interaction between spline variables and covariates for TDEs

## stpm2 [Lambert and Royston, 2009]

```
. stset time, failure(status == 1) id(id) scale(12) exit(time 60)

. stpm2 treatment, scale(hazard) df(4) eform nolog

Log likelihood =  -440.316                    Number of obs    =      506
```

| | exp(b) | Std. Err. | z | P>\|z\| | [95% Conf. Interval] |
|---|---|---|---|---|---|
| xb | | | | | | |
| treatment | .6594084 | .111509 | -2.46 | 0.014 | .4733827 | .9185368 |
| _rcs1 | 3.389716 | .4258797 | 9.72 | 0.000 | 2.649838 | 4.336179 |
| _rcs2 | .8879662 | .0724157 | -1.46 | 0.145 | .7567963 | 1.041871 |
| _rcs3 | 1.06315 | .0411503 | 1.58 | 0.114 | .9854806 | 1.146942 |
| _rcs4 | 1.016818 | .0199075 | 0.85 | 0.394 | .9785387 | 1.056594 |
| _cons | .229559 | .0272468 | -12.40 | 0.000 | .1819129 | .2896844 |

```
Note: Estimates are transformed only in the first equation.

. stcox treatment, nolog noshow

Cox regression -- Breslow method for ties
```

| _t | Haz. Ratio | Std. Err. | z | P>\|z\| | [95% Conf. Interval] |
|---|---|---|---|---|---|
| treatment | .6602897 | .1116672 | -2.45 | 0.014 | .4740025 | .9197894 |

## stpm2 [Lambert and Royston, 2009]

```
. stset time, failure(status == 2) id(id) scale(12) exit(time 60)

. stpm2 treatment, scale(hazard) df(4) eform nolog

Log likelihood = -448.73758                        Number of obs    =      506
```

|             | exp(b)   | Std. Err. | z      | P>|z| | [95% Conf. Interval] |          |
|-------------|----------|-----------|--------|-------|----------------------|----------|
| xb          |          |           |        |       |                      |          |
| treatment   | 1.202808 | .2047249  | 1.08   | 0.278 | .8616223             | 1.679097 |
| _rcs1       | 2.82908  | .2642265  | 11.13  | 0.000 | 2.355841             | 3.397384 |
| _rcs2       | .8685486 | .0544436  | -2.25  | 0.025 | .7681357             | .9820878 |
| _rcs3       | .9529595 | .0319403  | -1.44  | 0.151 | .8923696             | 1.017663 |
| _rcs4       | 1.027927 | .0213538  | 1.33   | 0.185 | .986915              | 1.070644 |
| _cons       | .17767   | .0237024  | -12.95 | 0.000 | .1367912             | .2307651 |

```
Note: Estimates are transformed only in the first equation.

. stcox treatment, nolog noshow
```

| _t        | Haz. Ratio | Std. Err. | z    | P>|z| | [95% Conf. Interval] |          |
|-----------|------------|-----------|------|-------|----------------------|----------|
| treatment | 1.20334    | .2048509  | 1.09 | 0.277 | .8619538             | 1.679937 |

## stpm2 [Lambert and Royston, 2009]

```
. stset time, failure(status == 3) id(id) scale(12) exit(time 60)

. stpm2 treatment, scale(hazard) df(4) eform nolog

Log likelihood = -231.45608                    Number of obs    =       506
```

|            | exp(b)    | Std. Err. | z      | P>\|z\| | [95% Conf. Interval] |          |
|------------|-----------|-----------|--------|---------|----------------------|----------|
| xb         |           |           |        |         |                      |          |
| treatment  | .6432149  | .1737196  | -1.63  | 0.102   | .3788467             | 1.092066 |
| _rcs1      | 2.638735  | .3351586  | 7.64   | 0.000   | 2.057219             | 3.384628 |
| _rcs2      | .7913665  | .0590788  | -3.13  | 0.002   | .683647              | .9160589 |
| _rcs3      | .9369818  | .0467358  | -1.30  | 0.192   | .8497164             | 1.033209 |
| _rcs4      | 1.029843  | .031817   | 0.95   | 0.341   | .9693337             | 1.09413  |
| _cons      | .097687   | .0179093  | -12.69 | 0.000   | .0681998             | .1399235 |

```
Note: Estimates are transformed only in the first equation.

. stcox treatment, nolog noshow
```

| _t         | Haz. Ratio | Std. Err. | z     | P>\|z\| | [95% Conf. Interval] |          |
|------------|------------|-----------|-------|---------|----------------------|----------|
| treatment  | .6460519   | .1745103  | -1.62 | 0.106   | .3804893             | 1.096964 |

## stpm2 [Lambert and Royston, 2009]

```
. stset time, failure(status == 3) id(id) scale(12) exit(time 60)

. stpm2 treatment, scale(hazard) df(4) tvc(treatment) dftvc(2) eform nolog
Log likelihood = -230.90611                    Number of obs    =        506
```

|                | exp(b)   | Std. Err. | z      | P>\|z\| | [95% Conf. Interval] |          |
|----------------|----------|-----------|--------|---------|----------------------|----------|
| xb             |          |           |        |         |                      |          |
| treatment      | .711078  | .2158501  | -1.12  | 0.261   | .3922222             | 1.289147 |
| _rcs1          | 2.805675 | .4977957  | 5.81   | 0.000   | 1.981588             | 3.972477 |
| _rcs2          | .7487466 | .0683538  | -3.17  | 0.002   | .6260772             | .895451  |
| _rcs3          | .9426525 | .0484762  | -1.15  | 0.251   | .8522722             | 1.042617 |
| _rcs4          | 1.032005 | .0318598  | 1.02   | 0.308   | .9714123             | 1.096377 |
| _rcs_treatment1 | .9101003 | .2468771  | -0.35  | 0.728   | .5347974             | 1.548778 |
| _rcs_treatment2 | 1.161785 | .1848084  | 0.94   | 0.346   | .8505949             | 1.586824 |
| _cons          | .0931347 | .0183948  | -12.02 | 0.000   | .0632401             | .1371608 |

```
Note: Estimates are transformed only in the first equation.
```

## Estimating cause-specific CIFs after fitting FPMs

**Cause-specific CIF,** $F_k(t)$

$$F_k(t) = \int_0^t \exp\left(-\sum_{k=1}^K \int_0^t h_k^{cs}(u)\mathrm{d}u\right) h_k^{cs}(u)\mathrm{d}u$$

**Cause-specific CIF, $F_k(t)$**

$$F_k(t) = \int_0^t \exp\left(-\sum_{k=1}^{K} \int_0^t h_k^{cs}(u)\mathrm{d}u\right) h_k^{cs}(u)\mathrm{d}u$$

Must be obtained by numerical approximation:

- Trapezoid method - `stpm2cif` [Hinchliffe and Lambert, 2013]
- Gauss-Legendre quadrature - `stpm2cr` [Mozumder et al., 2017]

## stpm2cif: Data setup

```
. local knotstvc_opt
. local bknotstvc_opt
. local k = 1
. foreach cause in _cancer _cvd _other {
  2.        stset time, failure(status == `k´) exit(time 60) scale(12)
  3.        cap stpm2 treatment, df(4) scale(h) eform nolog
  4.        estimates store stpm2`cause´
  5.        local bhknots`cause´ `e(bhknots)´
  6.        local boundknots`cause´ `e(boundary_knots)´
  7.        local knotstvc_opt `knotstvc_opt´ `cause´ `bhknots`cause´´
  8.        local bknotstvc_opt `bknotstvc_opt´ `cause´ `boundknots`cause´´
  9.        local k = `k´ + 1
 10. }
```

## stpm2cif: Data setup

```
. expand 3 // augment data k = 3 times
. bysort id: gen _cause=_n
. //create dummy variables for each cause of death
. gen _cvd=_cause==2
. gen _other=_cause==3
. gen _cancer=_cause==1
. //create cause of death event indicator variable
. gen _event=(_cause==status)
. label values _cause status
. foreach cause in _cancer _cvd _other {
  2.        gen treatment`cause´ = treatment*`cause´
  3. }
```

## stpm2cif: Data setup

```
. list id status time treatment _cause _event in 1/9, sep(9)
```
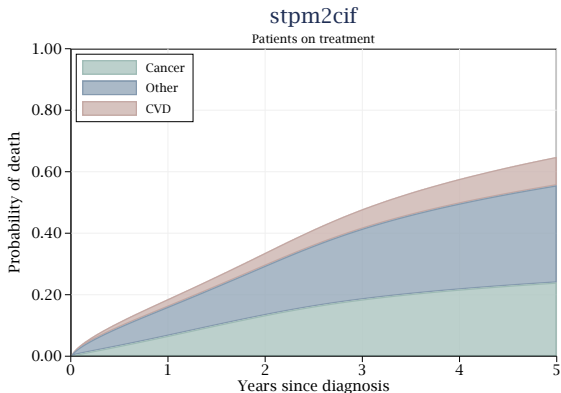
|     | id | status | time | treatm_t | _cause | _event |
|-----|----|--------|------|----------|--------|--------|
| 1.  | 1  | Censor | 72   | 0        | 1      | 0      |
| 2.  | 1  | Censor | 72   | 0        | 2      | 0      |
| 3.  | 1  | Censor | 72   | 0        | 3      | 0      |
| 4.  | 2  | Cancer | 1    | 0        | 1      | 1      |
| 5.  | 2  | Cancer | 1    | 0        | 2      | 0      |
| 6.  | 2  | Cancer | 1    | 0        | 3      | 0      |
| 7.  | 3  | CVD    | 40   | 1        | 1      | 0      |
| 8.  | 3  | CVD    | 40   | 1        | 2      | 1      |
| 9.  | 3  | CVD    | 40   | 1        | 3      | 0      |

## stpm2cif: Fitting the model

```
. stset time, failure(_event == 1) exit(time 60) scale(12)

. stpm2 treatment_cancer _cancer treatment_cvd _cvd treatment_other _other ///
> , scale(h) knotstvc(`knotstvc_opt´) bknotstvc(`bknotstvc_opt´) ///
> tvc(_cancer _cvd _other) rcsbaseoff nocons eform nolog

Log likelihood = -1120.5192                    Number of obs    =    1,518
```

|  | exp(b) | Std. Err. | z | P>|z| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| xb |  |  |  |  |  |  |
| treatment_cancer | .6593781 | .111504 | -2.46 | 0.014 | .4733607 | .9184951 |
| _cancer | .2295677 | .0272475 | -12.40 | 0.000 | .1819204 | .2896945 |
| treatment_cvd | 1.202808 | .2047249 | 1.08 | 0.278 | .8616223 | 1.679097 |
| _cvd | .17767 | .0237024 | -12.95 | 0.000 | .1367912 | .2307651 |
| treatment_other | .6432149 | .1737196 | -1.63 | 0.102 | .3788467 | 1.092066 |
| _other | .097687 | .0179093 | -12.69 | 0.000 | .0681998 | .1399235 |
| (*output omitted*) |  |  |  |  |  |  |

```
Note: Estimates are transformed only in the first equation.
```

## stpm2cif: Post-estimation

```
. stpm2cif cancer cvd other, cause1(treatment_cancer 1 _cancer 1) ///
> cause2(treatment_cvd 1 _cvd 1) cause3(treatment_other 1 _other 1) ci

. gen _totcif2_trt1 = CIF_cancer + CIF_cvd

. gen _totcif3_trt1 = _totcif2_trt1 + CIF_other
```
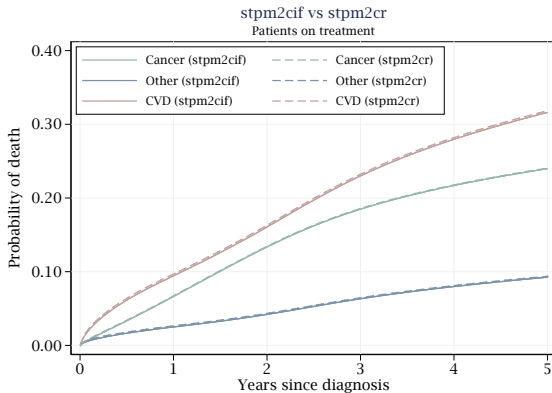
## stpm2cif: Post-estimation



```
. gen zeros = 0

. tw (rarea _totcif3_trt1 _totcif2_trt1  _newt, sort color(erose%80)) ///
> (rarea CIF_cancer _totcif2_trt1  _newt, sort color(emidblue%80)) ///
> (rarea zeros CIF_cancer _newt, sort color(eltgreen%80)), ...
```

```
. stset time, failure(status == 1,2,3) exit(time 60) scale(12)

. stpm2cr [cancer: treatment, scale(hazard) df(4)] ///
> [cvd: treatment, scale(hazard) df(4)] ///
> [other: treatment, scale(hazard) df(4)], ///
> events(status) cause(1 2 3) cens(0) eform model(csh)
```
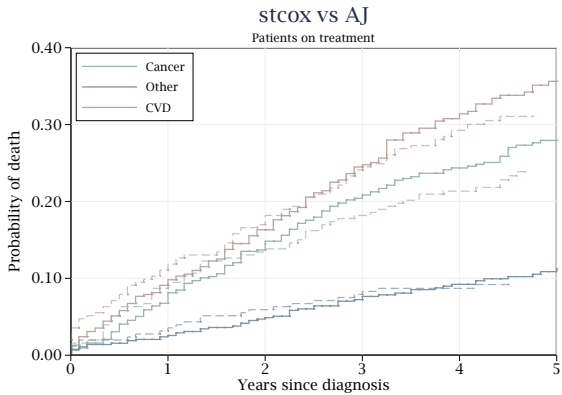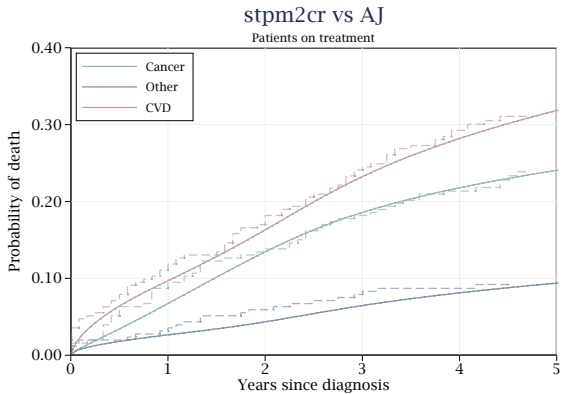
## stpm2cr: Post-estimation



```
. range newt 0 5 100

. predict cifgq_trt1, cif at(treatment 1) timevar(newt) ci
```
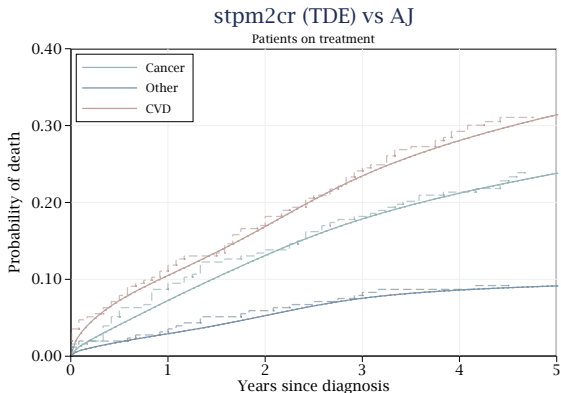
# Comparison with AJ estimates

# Comparison with AJ estimates

## Comparison with AJ estimates



stpm2cr (TDE) vs AJ
Patients on treatment

```
. stpm2cr [cancer: treatment, scale(hazard) df(4) tvc(treatment) dftvc(3)] ///
> [cvd: treatment, scale(hazard) df(4) tvc(treatment) dftvc(3)] ///
> [other: treatment, scale(hazard) df(4) tvc(treatment) dftvc(3)], ///
> events(status) cause(1 2 3) cens(0) eform model(csh)
```

## Note on computational time

```
. expand 500 //now 253,000 observations

. replace time = time + runiform()*0.0001

. replace id = _n
variable id was int now long
```

|  | Time (secs) |
|---|---|
| **stpm2cr model** | 52.60 |
| **stpm2 (stacked data)** | 76.59 |
| **stpm2cr predict (w/ CIs)** | 2.56 |
| **stpm2cif (w/ CIs)** | 11.10 |

- Restricted mean lifetime (RML) [Royston and Parmar, 2013; Andersen, 2013]
- Absolute & relative CIF measures
- Subdistribution hazard [Beyersmann et al., 2009]
- *Standardisation (to come)*

- Restricted mean lifetime (RML) [Royston and Parmar, 2013; Andersen, 2013] - double integration
- Absolute & relative CIF measures
- Subdistribution hazard [Beyersmann et al., 2009]
- *Standardisation (to come)*- predict for and average over every individual in study population

# Using the multistate package

## multistate [Crowther and Lambert, 2017]

- Written mainly by Michael (& Paul) for more complex multi-state models e.g. illness-death models
- Competing risks is a special case of multi-state models
- Can use `multistate` package to obtain equivalent non-parametric estimates and fit parametric models in presence of competing risks
- Uses a simulation approach for calculating transition probabilities i.e. cause-specific CIFs

## msset

```
. tab status, gen(cause)
. rename cause2 _cancer
. rename cause3 _cvd
. rename cause4 _other
. msset, id(id) states(_cancer _cvd _other) times(time time time) cr
. li id treatment status time _from _to _trans _start _stop _status _flag in 1/9, sep(9) noobs
```

| id | treatm_t | status | time | _from | _to | _trans | _start | _stop | _status | _flag |
|----|----------|--------|---------|-------|-----|--------|--------|-----------|---------|-------|
| 1  | 0        | Censor | 72.0024 | 1     | 2   | 1      | 0      | 72.002434 | 0       | 0     |
| 1  | 0        | Censor | 72.0024 | 1     | 3   | 2      | 0      | 72.002434 | 0       | 0     |
| 1  | 0        | Censor | 72.0024 | 1     | 4   | 3      | 0      | 72.002434 | 0       | 0     |
| 2  | 0        | Cancer | 1.00301 | 1     | 2   | 1      | 0      | 1.0030106 | 1       | 0     |
| 2  | 0        | Cancer | 1.00301 | 1     | 3   | 2      | 0      | 1.0030106 | 0       | 0     |
| 2  | 0        | Cancer | 1.00301 | 1     | 4   | 3      | 0      | 1.0030106 | 0       | 0     |
| 3  | 1        | CVD    | 40.008  | 1     | 2   | 1      | 0      | 40.007992 | 0       | 0     |
| 3  | 1        | CVD    | 40.008  | 1     | 3   | 2      | 0      | 40.007992 | 1       | 0     |
| 3  | 1        | CVD    | 40.008  | 1     | 4   | 3      | 0      | 40.007992 | 0       | 0     |

## msaj

```
. stset _stop, failure(_status == 1) scale(12) exit(time 60)

. msaj if treatment == 1, cr //ci

. sort _t

. li id status _trans _d _t P_AJ_? if P_AJ_1 != . in 1/45, noobs
```
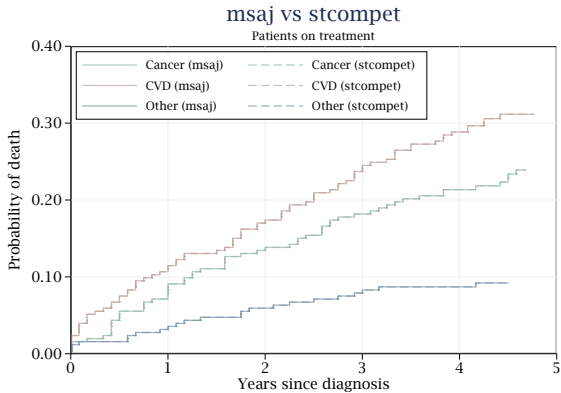
| id | status | _trans | _d | _t | P_AJ_1 | P_AJ_2 | P_AJ_3 | P_AJ_4 |
|----|--------|--------|----|----|--------|--------|--------|--------|
| 202 | Cancer | 1 | 1 | .00841895 | .99604743 | .00395257 | 0 | 0 |
| 105 | CVD | 2 | 1 | .00854265 | .99209486 | .00395257 | .00395257 | 0 |
| 151 | Other | 3 | 1 | .00855531 | .98814229 | .00395257 | .00395257 | .00395257 |
| 382 | CVD | 2 | 1 | .00866204 | .98418972 | .00395257 | .00790514 | .00395257 |
| 437 | CVD | 2 | 1 | .00869011 | .98023715 | .00395257 | .01185771 | .00395257 |
| 120 | Cancer | 1 | 1 | .00869888 | .97628458 | .00790514 | .01185771 | .00395257 |
| 502 | Cancer | 1 | 1 | .00881231 | .97233202 | .01185771 | .01185771 | .00395257 |
| 464 | CVD | 2 | 1 | .00886007 | .96837945 | .01185771 | .01581028 | .00395257 |
| 93 | Other | 3 | 1 | .00898155 | .96442688 | .01185771 | .01581028 | .00790514 |
| 492 | CVD | 2 | 1 | .00904977 | .96047431 | .01185771 | .01976285 | .00790514 |

```
. bysort P_AJ_2 (_t): gen first1 = _n==1

. bysort P_AJ_3 (_t): gen first2 = _n==1

. bysort P_AJ_4 (_t): gen first3 = _n==1
```

```
. stpm2 treatment if _trans==1, df(4) scale(h) eform nolog
. estimates store m1
. stpm2 treatment if _trans==2, df(4) scale(h) eform nolog
. estimates store m2
. stpm2 treatment if _trans==3, df(4) scale(h) eform nolog
. estimates store m3
. range tempt 0 5 100
. predictms , cr timevar(tempt) models(m1 m2 m3) at1(treatment 1)
```
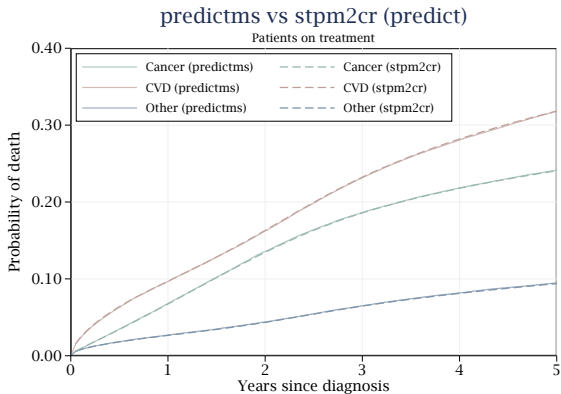
## predictms - without msset

```
. forvalues k = 1/3 {
2.         stset time, failure(status == `k´) id(id) scale(12) exit(time 60)
3.         stpm2 treatment, df(4) scale(h) eform nolog
4.         estimates store m`k´
5. }
. range tempt 0 5 100
. predictms , cr timevar(tempt) models(m1 m2 m3) at1(treatment 1)
```

# predictms

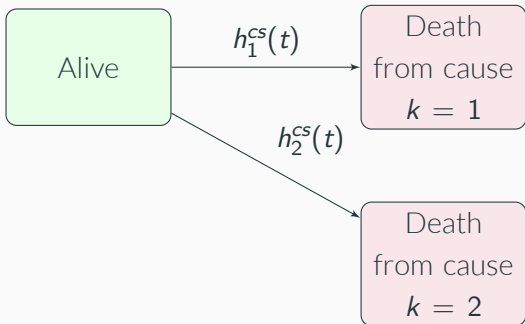## Summary of FPM tools for estimating cause-specific CIFs using CSHs

- Post-estimation command, `stpm2cif`
  - Requires augmenting data before `stpm2`
  - Fitting a single model means interpretation is difficult and more room for errors
  - Uses a basic numerical integration method - slow for larger datasets

## Summary of FPM tools for estimating cause-specific CIFs using CSHs

- Post-estimation command, `stpm2cif`
  - Requires augmenting data before `stpm2`
  - Fitting a single model means interpretation is difficult and more room for errors
  - Uses a basic numerical integration method - slow for larger datasets
- Using `stpm2cr` as a wrapper followed by `predict`
  - Fits separate `stpm2` models for each cause of death without data augmentation
  - Uses quicker numerical integration method
  - Can obtain other useful predictions e.g. restricted mean lifetime/comparative predictions

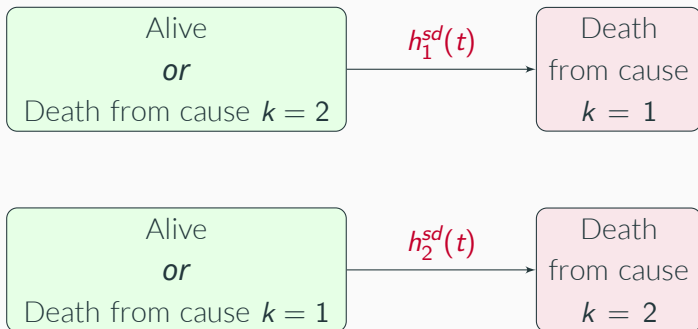## Summary of FPM tools for estimating cause-specific CIFs using CSHs

- Via the `predictms` command provided as a part of the `multistate` package
    - Uses a simulation approach. Can alternatively use AJ estimator to save on computational time
    - Can also be used without requiring `msset`
    - Extremely versatile - has some very useful features and post-estimation options

**What about modelling covariate effects on the risk of dying from a particular cause?**
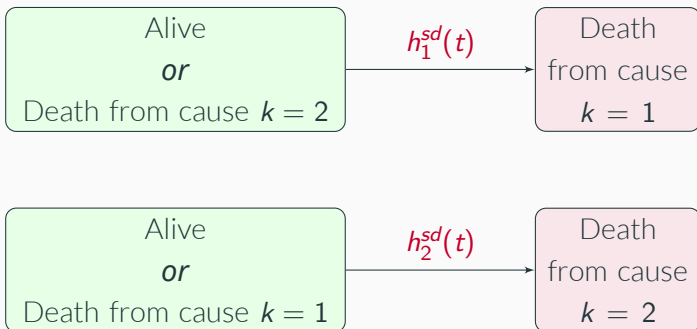
## Cause-specific hazards

## Subdistribution hazards



**Subdistribution hazard (SDH) rate,** $h_k^{sd}(t)$

The instantaneous rate of failure at time $t$ from cause $D = k$
amongst those who have not died, or have died from any of the
other causes, where $D \neq k$

## Subdistribution hazards



**Subdsitribution hazard (SDH) rate,** $h_k^{sd}(t)$

$$h_k^{sd}(t) = \lim_{\Delta t \to 0} \frac{P(t < T \le t + \Delta t, D = k | T > t \cup (T \le t \cap D \ne k)}{\Delta t}$$

# SDH relationship with cause-specific CIF

**Cause-specific CIF,** $F_k(t)$

$$F_k(t) = 1 - \exp\left[-\int_0^t h_k^{sd}(u)\mathrm{d}u\right]$$

# SDH relationship with cause-specific CIF

**Cause-specific CIF,** $F_k(t)$

$$F_k(t) = 1 - \exp\left[-\int_0^t h_k^{sd}(u)\mathrm{d}u\right]$$

**Note**

$$1 - F_k(t) = P(D \neq k) + S_k^{sd}(t)$$

## Standard approach: Fine & Gray model

Derived in a similar way to cause-specific Cox PH model as described by Fine and Gray [1999].

**SDH Regression Model (Fine & Gray Model)**

$$h_k^{sd}(t \mid \mathbf{x}_k) = h_{0k} \exp\left(\beta_k^{sd}\mathbf{x}_k\right)$$

$\beta_k^{sd}$: row vector of coefficients/log-SDH ratio for cause $k$

$\mathbf{x}_k$: column vector of covariates for cause $k$

$h_{0k}$: the baseline SDH function

## Standard approach: Fine & Gray model

Derived in a similar way to cause-specific Cox PH model as described by Fine and Gray [1999].

**SDH Regression Model (Fine & Gray Model)**

$$h_k^{sd}(t \mid \mathbf{x}_k) = h_{0k} \exp\left(\beta_k^{sd} \mathbf{x}_k\right)$$

$\beta_k^{sd}$: row vector of coefficients/log-SDH ratio for cause $k$

$\mathbf{x}_k$: column vector of covariates for cause $k$

$h_{0k}$: the baseline SDH function

SHR = association on the effect of a covariate on risk of dying from cause $k$

# Time-dependent censoring weights

- Need to consider those who have already died from other competing causes of death in risk-set
- Calculate missing censoring times for those that died from other causes by applying time-dependent weights to partial likelihood
- Influence of weights decreases over-time as the probability of being censored increases
- Further details given by Lambert et al. [2017] and Geskus [2011]

## stcrreg

```
. *Cancer
. stset time, failure(status == 1) exit(time 60) scale(12)

. stcrreg treatment, compete(status == 2, 3)

        failure _d:  status == 1
  analysis time _t:  time/12
  exit on or before:  time 60

Iteration 0:   log pseudolikelihood = -875.12133
Iteration 1:   log pseudolikelihood =  -875.1123
Iteration 2:   log pseudolikelihood =  -875.1123

Competing-risks regression                    No. of obs       =        506
                                               No. of subjects  =        506
Failure event   : status == 1                  No. failed       =        145
Competing events: status == 2 3                No. competing    =        197
                                               No. censored     =        164

                                               Wald chi2(1)     =       6.74
Log pseudolikelihood = -875.1123               Prob > chi2      =     0.0094
```

|           |          | Robust    |       |       |                      |
| --------- | -------- | --------- | ----- | ----- | -------------------- |
| _t        | SHR      | Std. Err. | z     | P>\|z\| | [95% Conf. Interval] |
| treatment | .6454653 | .1088223  | -2.60 | 0.009 | .463836    .8982171  |

```
. stcurve, cif at(treatment=1) outfile(cancer1, replace) range(0 5)
```

## stcrreg

```
. *CVD
. stset time, failure(status == 2) exit(time 60) scale(12)

. stcrreg treatment, compete(status == 1, 3)

         failure _d:  status == 2
   analysis time _t:  time/12
  exit on or before:  time 60

Iteration 0:   log pseudolikelihood = -848.00112
Iteration 1:   log pseudolikelihood = -847.83627
Iteration 2:   log pseudolikelihood = -847.83627

Competing-risks regression                      No. of obs       =        506
                                                No. of subjects  =        506
Failure event   : status == 2                   No. failed       =        140
Competing events: status == 1 3                 No. competing    =        202
                                                No. censored     =        164

                                                Wald chi2(1)     =       2.79
Log pseudolikelihood = -847.83627               Prob > chi2      =     0.0949

─────────────────────────────────────────────────────────────────────────────
                          Robust
         _t        SHR   Std. Err.      z    P>|z|     [95% Conf. Interval]
─────────────────────────────────────────────────────────────────────────────
  treatment   1.326649   .2245377    1.67   0.095      .9521137    1.848517

. stcurve, cif at(treatment=1) outfile(cvd1, replace) range(0 5)
```

## stcrreg

```
. *Other causes
. stset time, failure(status == 3) exit(time 60) scale(12)

. stcrreg treatment, compete(status == 1, 2)

        failure _d:  status == 3
  analysis time _t:  time/12
  exit on or before:  time 60

Iteration 0:   log pseudolikelihood = -349.42345
Iteration 1:   log pseudolikelihood = -349.41144
Iteration 2:   log pseudolikelihood = -349.41144

Competing-risks regression                    No. of obs      =        506
                                              No. of subjects =        506
Failure event   : status == 3                 No. failed      =         57
Competing events: status == 1 2               No. competing   =        285
                                              No. censored    =        164

                                              Wald chi2(1)    =       2.14
Log pseudolikelihood = -349.41144             Prob > chi2     =     0.1432
```

|  | | Robust | | | | |
| _t | SHR | Std. Err. | z | P>\|z\| | [95% Conf. Interval] |
|---|---|---|---|---|---|---|
| treatment | .6736976 | .1817566 | -1.46 | 0.143 | .3970267 | 1.143169 |

```
. stcurve, cif at(treatment=1) outfile(other1, replace) range(0 5)
```

36/47

**Log-cumulative SDH FPM**

$$\ln\left(H_k^{sd}(t \mid \mathbf{x}_k)\right) = s_k(\ln t; \boldsymbol{\gamma}_k, \mathbf{m}_{0k}) + \boldsymbol{\beta}_k^{sd}\mathbf{x}_k$$

**Log-cumulative non-proportional SDH FPM**

$$\ln\left(H_k^{sd}(t \mid \mathbf{x}_k)\right) = s_k(\ln t; \boldsymbol{\gamma}_k, \mathbf{m}_{0k}) + \beta_k^{sd}\mathbf{x}_k + \sum_{l=1}^{E} s_k(\ln t; \boldsymbol{\alpha}_{lk}, \mathbf{m}_{lk})\mathbf{x}_{lk}$$

## FPMs on (log-cumulative) SDH scale

**Log-cumulative non-proportional SDH FPM**

$$\ln\left(H_k^{sd}(t \mid \mathbf{x}_k)\right) = s_k(\ln t; \boldsymbol{\gamma}_k, \mathbf{m}_{0k}) + \beta_k^{sd}\mathbf{x}_k + \sum_{l=1}^{E} s_k(\ln t; \boldsymbol{\alpha}_{lk}, \mathbf{m}_{lk})\mathbf{x}_{lk}$$

1. Apply time-dependent censoring weights to the likelihood function for each cause $k$ (`stcrprep`) [Lambert et al., 2017]
2. Model all $k$ causes of death simultaneously directly using the full likelihood function (`stpm2cr`) [Mozumder et al., 2017; Jeong and Fine, 2007]

## stcrprep

```
. stset time, failure(status == 1,2,3) exit(time 60) scale(12) id(id)
. gen cod2 = cond(_d==0,0,status)
. stcrprep, events(cod2) keep(treatment ) trans(1 2 3) wtstpm2 censcov(treatment) every(1)
. gen event = cod2 == failcode
. stset tstop [iw=weight_c], failure(event) enter(tstart) noshow
  (output omitted)
```
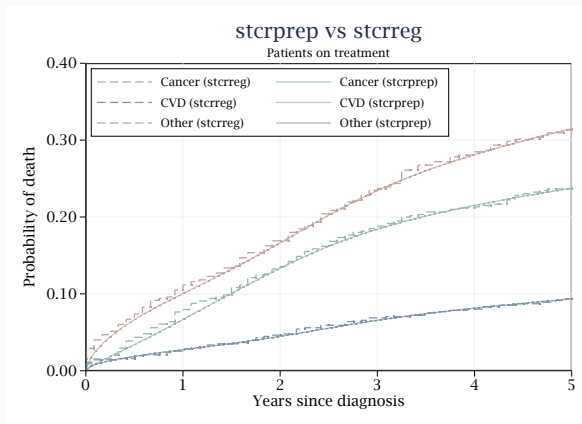
# stcrprep

```
. stpm2 treatment_cancer _cancer treatment_cvd _cvd treatment_other _other ///
> , scale(h) knotstvc(`knotstvc_opt´) bknotstvc(`bknotstvc_opt´) ///
> tvc(_cancer _cvd _other) rcsbaseoff nocons eform nolog
note: delayed entry models are being fitted

Log likelihood = -1228.025                      Number of obs    =      3,688
```

|                  | exp(b)    | Std. Err. | z      | P>|z| | [95% Conf. Interval] |           |
|------------------|-----------|-----------|--------|-------|----------------------|-----------|
| xb               |           |           |        |       |                      |           |
| treatment_cancer | .6408643  | .1083623  | -2.63  | 0.009 | .4600852             | .8926761  |
| _cancer          | .3060732  | .0335208  | -10.81 | 0.000 | .2469463             | .3793569  |
| treatment_cvd    | 1.329932  | .2263497  | 1.68   | 0.094 | .9527038             | 1.856525  |
| _cvd             | .2029639  | .0262824  | -12.32 | 0.000 | .1574686             | .2616034  |
| treatment_other  | .6740861  | .1819979  | -1.46  | 0.144 | .3970979             | 1.144282  |
| _other           | .1034306  | .0183681  | -12.78 | 0.000 | .0730273             | .1464916  |
| *(output omitted)* |         |           |        |       |                      |           |

```
Note: Estimates are transformed only in the first equation.

. predict cif_stcrprep_cancer, at(treatment_cancer 1 _cancer 1) zeros failure timevar(tempt)

. predict cif_stcrprep_cvd, at(treatment_cvd 1 _cvd 1) zeros failure timevar(tempt)

. predict cif_stcrprep_other, at(treatment_other 1 _other 1) zeros failure timevar(tempt)
```

# stcrprep

## stpm2cr

```
. stset time, failure(status == 1,2,3) exit(time 60) scale(12)

. stpm2cr [cancer: treatment, scale(hazard) df(4)] ///
> [cvd: treatment, scale(hazard) df(4)] ///
> [other: treatment, scale(hazard) df(4)], ///
> events(status) cause(1 2 3) cens(0) eform
  (output omitted)

. predict cifgq_trt1, cif at(treatment 1) timevar(tempt)
Calculating predictions for the following causes: 1 2 3
```

```
. stset time, failure(status == 1,2,3) exit(time 60) scale(12)

. stpm2cr [cancer: treatment, scale(hazard) df(4)] ///
> [cvd: treatment, scale(hazard) df(4)] ///
> [other: treatment, scale(hazard) df(4)], ///
> events(status) cause(1 2 3) cens(0) eform
  (output omitted)

. predict cifgq_trt1, cif at(treatment 1) timevar(tempt)
Calculating predictions for the following causes: 1 2 3
```

Above is not comparable with time-dependent censoring weights
approach as we assume proportionality for the competing causes
of death.

## stpm2cr

```
. stpm2cr [cancer: treatment, scale(hazard) df(4)] ///
> [cvd: treatment, scale(hazard) df(4) tvc(treatment) dftvc(3)] ///
> [other: treatment, scale(hazard) df(4) tvc(treatment) dftvc(3)], ///
> events(status) cause(1 2 3) cens(0) eform
  (output omitted)
Log likelihood = -1117.3418                 Number of obs    =       506
```

| | exp(b) | Std. Err. | z | P>|z| | [95% Conf. Interval] |
|---|---|---|---|---|---|
| **cancer** | | | | | | |
| treatment | .647454 | .1094638 | -2.57 | 0.010 | .464834 | .9018201 |
| (output omitted) | | | | | | |
| _cons | .1889881 | .0229604 | -13.71 | 0.000 | .1489433 | .2397993 |
| (output omitted) | | | | | | |

## stpm2cr

```
. stpm2cr [cancer: treatment, scale(hazard) df(4) tvc(treatment) dftvc(3)] ///
> [cvd: treatment, scale(hazard) df(4)] ///
> [other: treatment, scale(hazard) df(4) tvc(treatment) dftvc(3)], ///
> events(status) cause(1 2 3) cens(0) eform
  (output omitted)
```
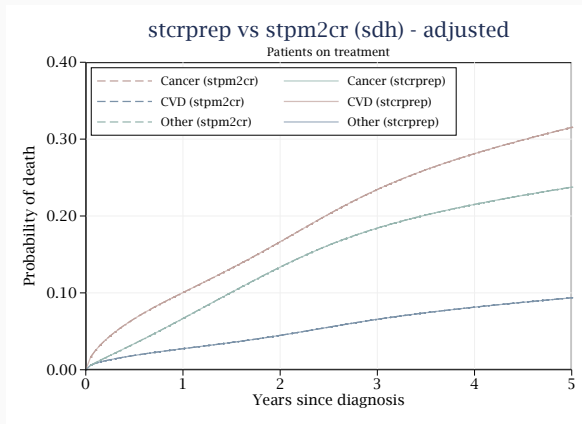
|                   | exp(b)    | Std. Err. | z      | P>|z| | [95% Conf. Interval] |          |
|-------------------|-----------|-----------|--------|-------|----------------------|----------|
| (output omitted)  |           |           |        |       |                      |          |
| cvd               |           |           |        |       |                      |          |
| treatment         | 1.336129  | .2273682  | 1.70   | 0.089 | .9571939             | 1.865077 |
| (output omitted)  |           |           |        |       |                      |          |
| _cons             | .1366028  | .0187788  | -14.48 | 0.000 | .1043385             | .178844  |
| (output omitted)  |           |           |        |       |                      |          |

## stpm2cr

```
. stpm2cr [cancer: treatment, scale(hazard) df(4) tvc(treatment) dftvc(3)] ///
> [cvd: treatment, scale(hazard) df(4) tvc(treatment) dftvc(3)] ///
> [other: treatment, scale(hazard) df(4)], ///
> events(status) cause(1 2 3) cens(0) eform
  (output omitted)
```

|  | exp(b) | Std. Err. | z | P>|z| | [95% Conf. Interval] |  |
|---|---|---|---|---|---|---|
| (output omitted) |  |  |  |  |  |  |
| other |  |  |  |  |  |  |
| treatment | .6771057 | .1827954 | -1.44 | 0.149 | .3988974 | 1.149349 |
| (output omitted) |  |  |  |  |  |  |
| _cons | .0720086 | .0138407 | -13.69 | 0.000 | .0494056 | .1049525 |

# Comparing stcrprep and stpm2cr



stcrprep vs stpm2cr (sdh) - adjusted

```
. expand 100 //now 50,060 observations

. replace time = time + runiform()*0.0001

. replace id = _n
variable id was int now long
```

|  | Time |
|---|---|
| **stcrreg (total)** | 53 mins |
| **stcrprep (total)** | 1 min |
| **stpm2cr** | 17 secs |

## On which scale should we model?

Cause-specific hazards

- Risk-set is defined in usual way - easy to understand

Subdistribution hazards

- Maintains direct relationship with cause-specific CIF

## On which scale should we model?

**Cause-specific hazards**

- Risk-set is defined in usual way - easy to understand
- Infer covariate effects on the rate of dying from a cause
  - For research questions on aetiology and causal effects

**Subdistribution hazards**

- Maintains direct relationship with cause-specific CIF
- Infer covariate effects on the risk of dying from a cause
  - For research questions on prognosis

## On which scale should we model?

**Cause-specific hazards**

- Risk-set is defined in usual way - easy to understand
- Infer covariate effects on the rate of dying from a cause
  - For research questions on aetiology and causal effects

**Subdistribution hazards**

- Maintains direct relationship with cause-specific CIF
- Infer covariate effects on the risk of dying from a cause
  - For research questions on prognosis

Many recommend inferences on all CSHs and cause-specific CIFs for a better understanding on the overall impact of cancer [Lambert et al., 2017; Latouche et al., 2013; Beyersmann et al., 2007]

## What next?

- Standardisation post-estimation for FPMs on cause-specific log-cumulative hazard scale
- Standardisation post-estimation after `stpm2cr`
- Restricted mean survival time [Royston and Parmar, 2011] for `stpm2cr` and `stcrprep`
- Expected number of life-years lost decomposed by cause of death [Andersen, 2013]

P. K. Andersen. Decomposition of number of life years lost according to causes of death. *Statistics in Medicine*, 32:5278--85, Jul 2013.

J. Beyersmann, M. Dettenkofer, H. Bertz, and M. Schumacher. A competing risks analysis of bloodstream infection after stem-cell transplantation using subdistribution hazards and cause-specific hazards. *Statistics in Medicine*, 26 (30):5360--5369, Dec. 2007.

J. Beyersmann, A. Latouche, A. Buchholz, and M. Schumacher. Simulating competing risks data in survival analysis. *Stat Med*, 28(6):956--971, 2009.

M. J. Crowther and P. C. Lambert. Parametric multistate survival models: Flexible modelling allowing transition-specific distributions with application to estimating clinically useful measures of effect differences. *Statistics in medicine*, 36(29):4719--4742, 2017.

## References ii

J. P. Fine and R. J. Gray. A proportional hazards model for the subdistribution of a competing risk. *Journal of the American Statistical Association*, 446: 496--509., 1999.

R. B. Geskus. Cause-specific cumulative incidence estimation and the Fine and Gray model under both left truncation and right censoring. *Biometrics*, 67(1): 39--49, Mar 2011.

S. R. Hinchliffe and P. C. Lambert. Extending the flexible parametric survival model for competing risks. *The Stata Journal*, 13:344--355, 2013.

J.-H. Jeong and J. P. Fine. Parametric regression on cumulative incidence function. *Biostatistics*, 8(2):184--196, Apr 2007.

P. C. Lambert and P. Royston. Further development of flexible parametric models for survival analysis. *The Stata Journal*, 9:265--290, 2009.

P. C. Lambert, S. R. Wilkes, and M. J. Crowther. Flexible parametric modelling of the cause-specific cumulative incidence function. *Statistics in medicine*, 36(9): 1429--1446, 2017.

## References  iii

A. Latouche, A. Allignol, J. Beyersmann, M. Labopin, and J. P. Fine. A competing risks analysis should report results on all cause-specific hazards and cumulative incidence functions. *J Clin Epidemiol*, 66(6):648--653, Jun 2013.

S. I. Mozumder, M. J. Rutherford, P. C. Lambert, et al. A flexible parametric competing-risks model using a direct likelihood approach for the cause-specific cumulative incidence function. *Stata Journal*, 17(2):462--489, 2017.

P. Royston and M. K. B. Parmar. Flexible parametric proportional-hazards and proportional-odds models for censored survival data, with application to prognostic modelling and estimation of treatment effects. *Statistics in Medicine*, 21(15):2175--2197, Aug 2002.

P. Royston and M. K. B. Parmar. The use of restricted mean survival time to estimate the treatment effect in randomized clinical trials when the proportional hazards assumption is in doubt. *Stat Med*, 30(19):2409--2421, Aug 2011.

P. Royston and M. K. B. Parmar. Restricted mean survival time: an alternative to the hazard ratio for the design and analysis of randomized trials with a time-to-event outcome. *BMC medical research methodology*, 13:152, 2013. ISSN 1471-2288.