

## Competing Risks Analyses: Overview of Regression Models

Joseph C. Gardiner, Division of Biostatistics, Department of Epidemiology and Biostatistics,  
Michigan State University, East Lansing, MI 48824

### Abstract

Competing-risks analyses are methods for analyzing the time to a terminal event (such as death or failure) and its cause or type. The cumulative incidence function  $CIF(j, t)$  is the probability of death by time  $t$  from cause  $j$ . New options in the LIFETEST procedure provide for nonparametric estimation of the CIF from event times and their associated causes, allowing for right-censoring when the event and its cause are not observed. Cause-specific hazard functions that are derived from the CIFs are the analogs of the hazard function when only a single cause is present. Death by one cause precludes occurrence of death by any other cause, because an individual can die only once. Incorporating explanatory variables in hazard functions provides an approach to assessing their impact on overall survival and on the CIF. This semiparametric approach can be analyzed in the PHREG procedure. The Fine-Gray model defines a sub-distribution hazard function that has an expanded risk set, which consists of individuals at risk of the event by any cause at time  $t$ , together with those who experienced the event before  $t$  from any cause other than the cause of interest  $j$ . Finally, with additional assumptions a full parametric analysis is also feasible. We illustrate the application of these methods with empirical data sets.

### 1. Introduction

A typical survival analysis addresses estimation of the survival function from time to event data. The event is considered terminal (e.g., death or failure). Survival data in a sample of units are the observed death times or the last follow-up times in the situation of right censoring. The Kaplan-Meier estimator summarizes the survival experience by providing an estimator of the survival function. A competing risks situation occurs when both the death time  $T$  and its cause  $J$  are taken into consideration. The label  $J$  denotes the cause of failure. It is assumed that  $m$  mutually exclusive causes are operational. Death by one cause precludes occurrence of death by any other cause, because an individual can die only once. If censoring is present, then both  $T$  and  $J$  are not observed, but the last follow-up time is observed.

All definitions stem from the joint distribution of  $(T, J)$ . The *cumulative incidence function* (CIF) for cause  $j$  is  $F_j(t) = P[T \leq t, J = j]$  -- the probability of failure by time  $t$  from cause  $j$ . The overall *survival distribution function* is  $S(t) = P[T > t] = 1 - \sum_{j=1}^m F_j(t)$ . The *cause-specific hazard function* and *cumulative cause-specific hazard function* are  $\alpha_j(t) = \lim_{\Delta t \downarrow 0} P[T \leq t + \Delta t, J = j | T \geq t] / \Delta t$  and  $A_j(t) = \int_0^t \alpha_j(u) du$   $j=1, 2, \dots, m$ . The overall *hazard function* and *cumulative hazard function* are  $\alpha_0(t) = \sum_{j=1}^m \alpha_j(t)$  and  $A_0(t) = \int_0^t \alpha_0(u) du$ . The CIF and survival function can be expressed in terms of the hazards by

$$(1) \quad F_j(t) = \int_0^t S(u-) dA_j(u), \quad j=1, 2, \dots, m \quad \text{and} \quad S(t) = \exp(-A_0(t)).$$

Notice that the survival function and each CIF depend on all the cause-specific hazards.

### ***Multistate Model***

Competing risks may be formulated as a special case of a state-transition model (Andersen et al, 2002). An individual's event history is described by a stochastic process  $\{X(t): t \geq 0\}$  where  $X(t)$  is the state occupancy at time  $t$ . All individuals start in state '0'. During follow-up a transition  $0 \rightarrow j$  to the destination state  $j$  may occur. The time to exit state '0' is  $T = \inf\{t > 0: X(t) \neq 0\}$  -- the previously called failure time.

The *transition probabilities* are  $P_{0j}(0, t) = P[X(t) = j | X(0) = 0]$ ,  $j=0, \dots, m$  and the corresponding *transition intensities* are the same cause-specific hazards,  $\alpha_j(t) = \lim_{\Delta t \downarrow 0} P[X(t + \Delta t) = j | X(t-) = 0] / \Delta t$ ,  $j=1, 2, \dots, m$ .

Note that  $P_{00}(0, t) = S(t)$  is the survival function and  $P_{0j}(0, t) = F_j(t)$  since  $P[X(0) = 0] = 1$ . The analogy to a multistate model permits the application of the theory developed via counting processes.

## **2. Example**

Collett (2015) describes a study of 1761 patients who underwent a liver transplant between January 2000 and December 2010. Patients were followed to the end of 2012. The event of interest is the failure of the graft which if observed, was attributed for one of four causes: (1) organ rejection, (2) hepatic artery thrombosis, (3) recurrent disease, or (4) other reasons. Death with a functioning graft was regarded as a censoring event. Transplant success leads to heavy censoring: 260 (15%) patients have graft failure (see Table 1). The median age at transplant is 55 years, age range 18-74 years, and 61% are male. Patients in this study had one of three types of liver disease, primary biliary cirrhosis (PBC, 27%), primary sclerosing cholangitis (PSC, 20%), or alcoholic liver disease (ALD, 53%).

<b>Table 1: Causes of Graft Failure in N=1761 liver transplant patients</b>				
<b>Cause of Graft Failure</b>	<b>Frequency</b>	<b>Percent</b>	<b>Cumulative Frequency</b>	<b>Cumulative Percent</b>
<b>none†</b>	1501	85.24	1501	85.24
<b>reject</b>	29	1.65	1530	86.88
<b>thromb</b>	73	4.15	1603	91.03
<b>recdis</b>	42	2.39	1645	93.41
<b>other</b>	116	6.59	1761	100.00

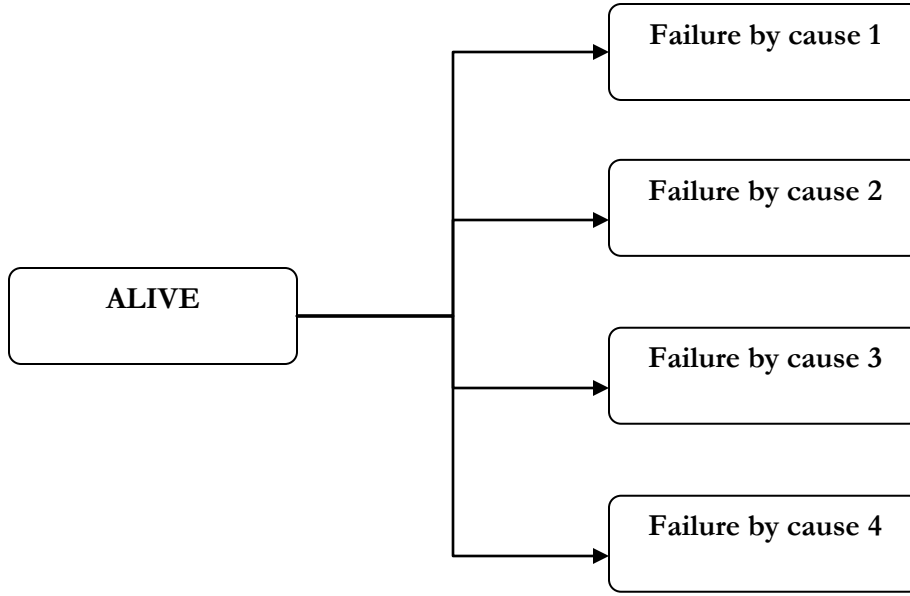
† Represents right censoring of graft failure time

The following formats are applied to the cause of graft failure (COF), liver disease (DISEASE), and patient (GENDER).

```
proc format;
value disease 1='PBC' 2='PSC' 3='ALD';
value gender 1='male' 2='female';
value cof 0='none' 1='reject' 2='thromb' 3='recdis' 4='other';
run;
```

Figure 1 is a representation of four competing risks as a multistate transition model.

**Figure 1: Competing risks as a multistate model**



CAUSE 1: Rejection

CAUSE 2: Thrombosis

CAUSE 3: Recurrent disease

CAUSE 4: Other causes

Only one path is observable in an individual  
 $0 \rightarrow 1$ ,  $0 \rightarrow 2$ ,  $0 \rightarrow 3$ , or  $0 \rightarrow 4$   
 Or, none (right censoring)

### 3. Nonparametric Analysis

From a random sample of observations  $\{(T_i, J_i) : 1 \leq i \leq n\}$  with potential random right censoring in the death times, estimation of the CIF follows from (1). The inputs are:

$Y(t) = \#$  at risk for any cause of failure at time just prior to  $t$

$N_j(t) =$  cumulative count of failures by cause  $j$  up to time  $t$

$N(t) =$  cumulative count of failures by *any* cause up to time  $t$ .

From (1) the CIF estimator is naturally given by

$$(2) \quad \hat{F}_j(t) = \sum_{u \leq t} \hat{S}(u-) \frac{\Delta N_j(u)}{Y(u)}$$

where  $\hat{S}(t) = \prod_{u \leq t} \left(1 - \frac{\Delta N(u)}{Y(u)}\right)$  denotes the Kaplan-Meier (KM) estimator of the survival function. If

there are no competing risks, that is only one cause of failure is operational,  $\hat{F}_j(t)$  reduces to  $1 - \hat{S}(t)$

because  $\hat{S}(t) = \left(1 - \frac{\Delta N(t)}{Y(t)}\right) \hat{S}(t-)$ . Also summation of (2) over all causes returns the KM estimator giving

reasons for using the KM estimator in (2) instead of the Breslow estimator  $\tilde{S}(t) = \exp\left(-\sum_{u \leq t} \frac{\Delta N(u)}{Y(u)}\right)$ .

PROC LIFETEST bears the brunt of computing  $\hat{F}_j(t)$ . The asymptotic normal distribution of  $\hat{F}_j(t)$  and an estimator of the asymptotic variance  $Var(\hat{F}_j(t))$  provides a basis for statistical inference. Analogous to the comparison of two or more survival functions from independent samples (e.g., log-rank test), we can test the equality of CIF's by Gray's test (Gray, 1988). LIFETEST offers two expressions for  $Var(\hat{F}_j(t))$ . The one based on the delta method (ERROR=delta) is

$$(3) \quad Var(\hat{F}_j(t)) = \sum_{u \leq t} \left\{ \hat{F}_j(t) - \hat{F}_j(u) \right\}^2 \frac{\Delta N(u)}{Y(u)(Y(u) - \Delta N(u))} + \sum_{u \leq t} \left\{ \hat{S}(u-) \right\}^2 \frac{\Delta N_j(u)(Y(u) - \Delta N_j(u))}{\{Y(u)\}^3} - 2 \sum_{u \leq t} \left\{ \hat{F}_j(t) - \hat{F}_j(u) \right\} \hat{S}(u-) \frac{\Delta N_j(u)}{\{Y(u)\}^2}.$$

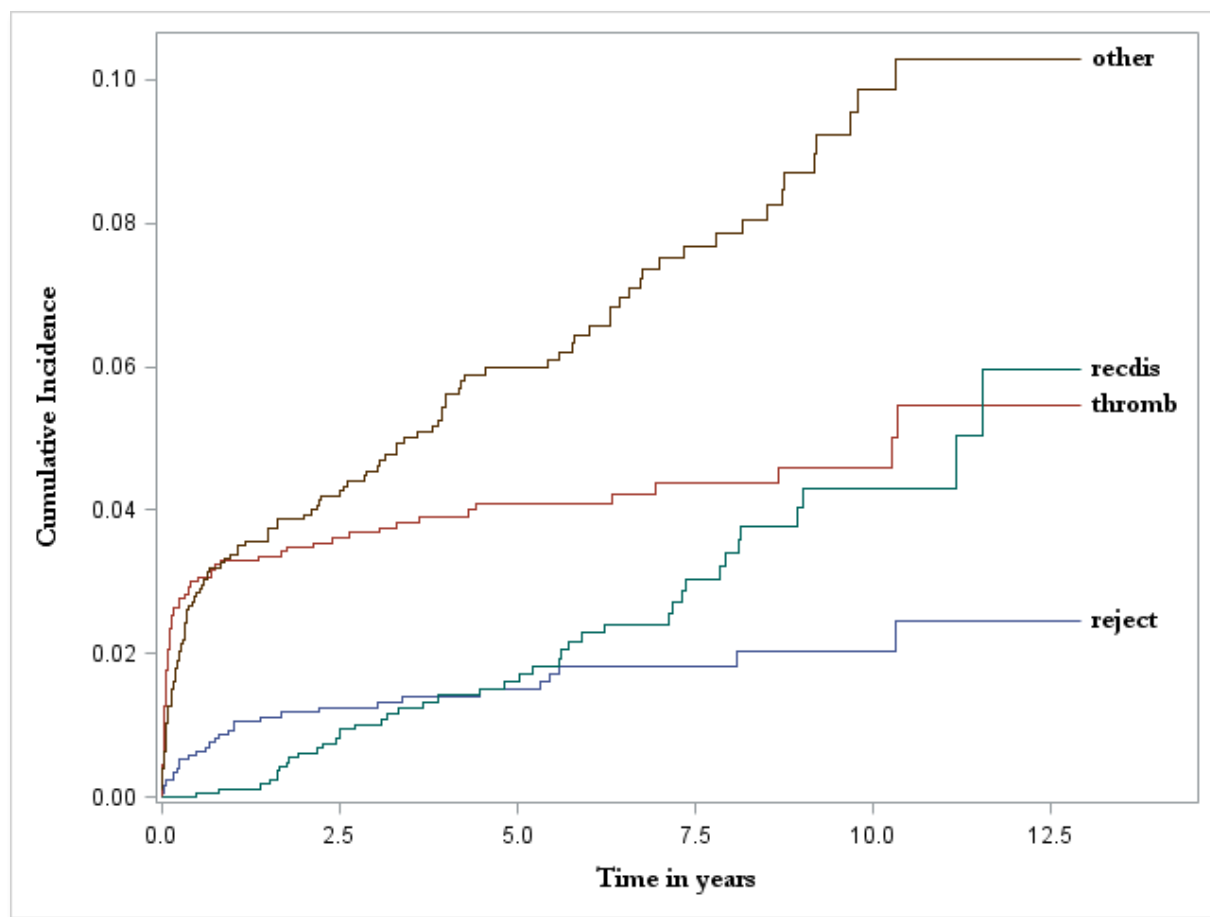
Another formula (ERROR=Aalen) is derived from application of the theory of counting processes in derivation of the distribution of  $\hat{F}_j(t)$ . Furthermore, different methods for adjustment for ties in failure times would lead to slightly different expressions for the variance. Formula (3) is a reasonable compromise and is in the spirit of Greenwood's formula for the asymptotic variance of the KM estimator. For further elaboration see Allignol et al (2010), Aalen et al (2008) and Andersen et al (1993).

Beginning in SAS® Analytics 14.1 use PROC LIFETEST to obtain CIF estimates, standards errors, confidence intervals (pointwise) and plotting of the CIFs. A macro %CIF is available in the AUTOCALL library in version 13.2.

```
proc lifetest data=liver plots=cif(CL) conftype=loglog
    error=delta outcif=cif_est;
time time*cof(0)/failcode=1 2 3 4;
format cof cof.;
run;
```

The TIME statement requests estimation of the CIFs for each type of graft failure, COF=1, 2, 3, 4. The value COF=0 denotes a right censored failure time. To restrict the analysis to a single failure type, for example graft failure due to rejection, specify **failcode=1**. Other options in the LIFETEST statement request that estimates be saved in the data set **cif\_est**. It also contains standard errors computed from expression (3) (**error=delta**) and pointwise 95% confidence intervals obtained after log(-log) transformation (**conftype=loglog**). The PLOTS option requests a plot of the CIF estimates and confidence intervals. Figure 2 overlays the four CIFs produced by PROC SGPLOT on the output data set **cif\_est**.

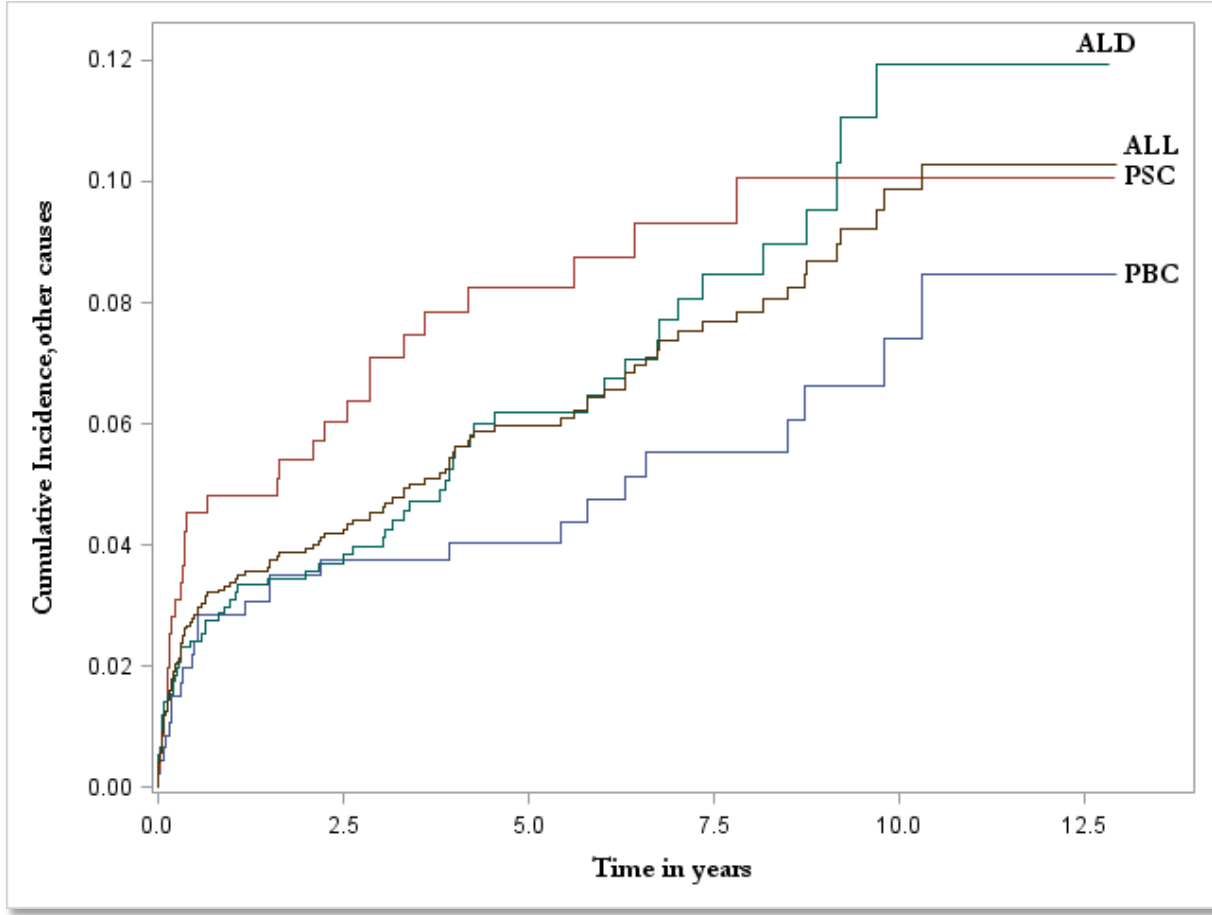
Figure 2: Cumulative incidence functions for four causes of graft failure in liver transplantation



A STRATA statement is needed to produce Gray's test of comparison of two or more CIFs (for the same failure code) in independent groups. For example, we might compare by disease groups (PBC, PSC, ALD) the cumulative incidence of graft failure due to other causes (**failcode=4**). The 2 DF chi-square test is not significant ( $p=0.184$ ). Figure 3 overlays the CIFs for the three disease groups and the overall CIF --- labelled 'ALL'. This is the same curve labelled 'other' in figure 2. Note that each curve is extended from its last failure time to the maximum observed time in the data set.

In general a  $K$ -sample test compares the CIFs across  $K$ -groups resulting in a  $K-1$  DF chi-square test of the null hypothesis  $H_0 : F(t,1) = F(t,2) = \dots = F(t,K)$  for all  $t$ . For the specified cause of failure  $F(t,k)$  denotes the CIF in group  $k$ . If our competing risks data are stratified by say, **center** a stratified test of  $H_0$  can be carried out. The modified syntax is

```
strata center/group=disease;
time time*cof(0)/failcode=4;
```

**Figure 3: Cumulative incidence functions by disease group for graft failure due to other causes**

#### 4. Regression models

The comparison of CIFs in independent groups is a special case of more general regression models that allow for assessing the influence of covariates on the CIF. For regression analyses with covariates  $\mathbf{z}$  posit a model for the cause-specific hazards  $\alpha_j(t)$ ,  $j = 1, \dots, m$  akin to a Cox model:

$$(4) \quad \alpha_j(t | \mathbf{z}) = \alpha_{j,0}(t) \exp(\mathbf{z}' \boldsymbol{\beta}_j)$$

where the baseline cause-specific hazard  $\alpha_{j,0}(t)$  and the regression coefficients  $\boldsymbol{\beta}_j$  are also made cause-specific. However this need not be the case. We could use a model that has the same  $\boldsymbol{\beta}$ -coefficient for some covariates (e.g., age, gender) across the different causes, while having cause-specific  $\boldsymbol{\beta}$ -coefficients for some other covariates (e.g., disease-type). Although time-varying covariates could be included in (4) just as one could do in a multistate model, the same issues of interpretation of the estimate of the survival function  $S(t | \mathbf{z})$  persists for estimates of the CIFs  $F_j(t | \mathbf{z})$ . To avoid this morass, our discussion here is restricted to time-invariant covariates—the baseline variables age, gender and disease in our example on graft failure in liver transplantation.

Estimation the parameters in model (4) is carried out in PROC PHREG by maximization of the partial

$$\text{likelihood } L(\beta) = \prod_{i \geq 0} \prod_{i=1}^n \prod_{j=1}^m \left( \frac{Y_i(t) \exp(\mathbf{z}'_i \beta_j)}{\sum_{k=1}^n Y_k(t) \exp(\mathbf{z}'_k \beta_j)} \right)^{\Delta N_{ji}(t)} \quad \text{where } Y_i(t) \text{ is the at-risk indicator at time } t$$

for subject  $i$  for any cause and  $N_{ji}(t)$  is the indicator for failure due to cause  $j$  by time  $t$ . Since  $L(\beta)$  is the product of  $m$  partial likelihoods, one for each cause of failure, we can estimate  $\beta = \{\beta_j, j = 1, \dots, m\}$  by separate calls to PHREG. Instead, in our application consider a single invocation on the expanded data set LIVER2.

```
data LIVER2;
set liver;
do i=1 to 4;
stratum=compress('0' || vvalue(i));
fstatus=(cof=i);
output;
end;
run;
```

In LIVER2 the original data set is quadrupled to have 4 blocks of 1761 patients. Each block represents those at risk for the cause of failure identified by STRATUM. For example in Table 2 STRATUM=01 denotes the cohort at risk for the 0→1(reject) cause of failure. The final status indicator of a patient is FSTATUS—with value 1 for event, and value 0 for censoring and the competing causes 2, 3 and 4.

stratum	cof	fstatus	Frequency	Percent	Cumulative Frequency	Cumulative Percent
01	0	0	1501	21.31	1501	21.31
01	1	1	29	0.41	1530	21.72
01	2	0	73	1.04	1603	22.76
01	3	0	42	0.60	1645	23.35
01	4	0	116	1.65	1761	25.00
02	0	0	1501	21.31	3262	46.31
02	1	0	29	0.41	3291	46.72
02	2	1	73	1.04	3364	47.76
02	3	0	42	0.60	3406	48.35
02	4	0	116	1.65	3522	50.00
03	0	0	1501	21.31	5023	71.31
03	1	0	29	0.41	5052	71.72
03	2	0	73	1.04	5125	72.76
03	3	1	42	0.60	5167	73.35
03	4	0	116	1.65	5283	75.00
04	0	0	1501	21.31	6784	96.31
04	1	0	29	0.41	6813	96.72
04	2	0	73	1.04	6886	97.76
04	3	0	42	0.60	6928	98.35
04	4	1	116	1.65	7044	100.00

The four causes of failure are analyzed simultaneously with cause-specific hazard ratio estimation by

```
ods output hazardratios=hrs type3=type3;
proc phreg data=liver2;
strata stratum;
class stratum;
class gender(ref='female') disease(ref='ALD')/param=glm;
model time*fstatus(0)=age*stratum gender*Stratum disease*stratum;
format disease disease. gender gender. cof cof.;
hazardratio age/at(stratum=all) units=5 cl=wald;
hazardratio gender/at(stratum=all) diff=all cl=wald;
hazardratio disease/at(stratum=all) diff=all cl=wald;
run;
```

Significance tests of the covariate effects overall are in the data set **type3**. Table 3 shows that there are significant effects for age and for disease, but not for gender.

Table 3: Covariate effects (Type 3 Tests)			
Effect	DF	Wald Chi-Square	Pr > ChiSq
age*stratum	4	15.1941	0.0043
stratum*gender	4	1.7737	0.7773
stratum*disease	8	19.5035	0.0124

An examination for the cause-specific hazard ratios in the data set **hrs** indicates that age is significant only for COF=4. For a 5 year increase in age, the cause-specific hazard ratio is 1.14 (95% CI: 1.03, 1.27). With respect to disease we see that PBC and PSC differ, for failure due to recurrent disease (COF=3) and other causes (COF=4). The default DIFF=ALL creates redundant comparisons. The data set has three comparisons per stratum instead of two with one as referent. To display the cause-specific hazard ratios for disease by failure type use

```
%let opts=(Color=black Family=Garamond Size=12 Style=normal
Weight=Bold);
proc sgplot data=HRS(where=(substr(description,1,1)='d'));
scatter y=description x=hazardratio/xerrorlower=Waldlower
xerrorupper=Waldupper markerattrs=(symbol=DiamondFilled size=8);
refline 1 / axis=x;
xaxis label="Cause-specific Hazard Ratio" min=0 LABELATTRS=&opts
VALUEATTRS=&opts;
yaxis label="Comparisons" LABELATTRS=&opts VALUEATTRS=&opts;
run;
```

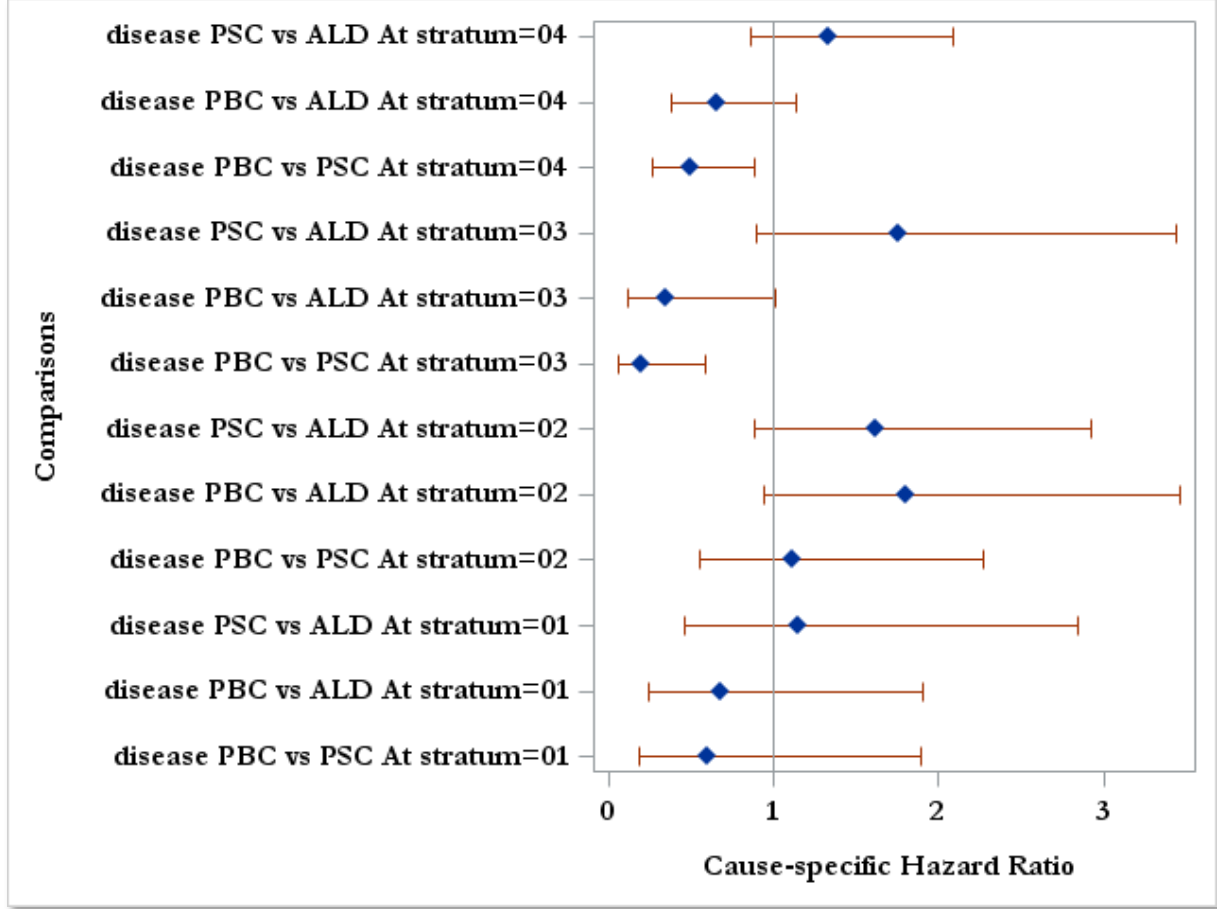
Figure 4 displays 95% CIs are for all 12 pairwise comparisons. We could select a subset of linearly independent comparisons, the maximum is 8. Use an ESTIMATE statement to get multiplicity-adjusted CIs and p-values. The syntax makes 4 comparisons between disease groups for graft failure due to causes 3 and 4 only. The non-positional coefficients are implied by the CLASS options and parameterization **param=glm**. Output is not shown.

```
estimate 'PBC vs ALD by cause 3' stratum*disease [1, 3 1] [-1, 3 3],
        'PBC vs PSC by cause 3' stratum*disease [1, 3 1] [-1, 3 2],
        'PBC vs ALD by cause 4' stratum*disease [1, 4 1] [-1, 4 3],
```



```
'PBC vs PSC by cause 4' stratum*disease [1, 4 1] [-1, 4 2]/joint
exp cl adjust=sidak;
```

Figure 4: Cause-specific hazard ratios for disease and 95% confidence intervals



### Estimation of Cumulative Incidence Functions

After estimation of the cause-specific hazard model (4) the next step is the estimation of the CIF  $F_j(t | \mathbf{z}_0)$  at a specified covariate profile  $\mathbf{z}_0$ . The preceding discussion suggests a reduced model without gender, with disease effects for causes 3 and 4, and an age effect for cause 4. Therefore we consider

$$\text{CAUSE}=1 \text{ (rejection)} \quad \alpha_1(t | \mathbf{z}) = \alpha_{1,0}(t)$$

$$\text{CAUSE}=2 \text{ (thrombosis)} \quad \alpha_2(t | \mathbf{z}) = \alpha_{2,0}(t)$$

$$\text{CAUSE}=3 \text{ (recurrent disease)} \quad \alpha_3(t | \mathbf{z}) = \alpha_{3,0}(t) \exp(\beta_{31}[\text{DISEASE} = \text{PBC}] + \beta_{32}[\text{DISEASE} = \text{PSC}])$$

$$\text{CAUSE}=4 \text{ (other)} \quad \alpha_4(t | \mathbf{z}) = \alpha_{4,0}(t) \exp(\beta_{41}[\text{DISEASE} = \text{PBC}] + \beta_{42}[\text{DISEASE} = \text{PSC}] + \beta_{43} \text{AGE})$$

This model does not share covariate effects.

To estimate this model we must use explicit ‘interaction’ terms with **I3**, **I4**.

```
proc phreg data=liver2;
strata stratum;
class disease(ref='ALD')/param=ref;
model time*fstatus(0)=disease*I3 disease*I4 age*I4;
I3=(stratum='03');I4=(stratum='04');
format disease disease.;
baseline out=stats_ph(where=(stratum=stratum_s)) covariates=covar
cumhaz=cumhaz survival=survival /method=ch;
run;
```

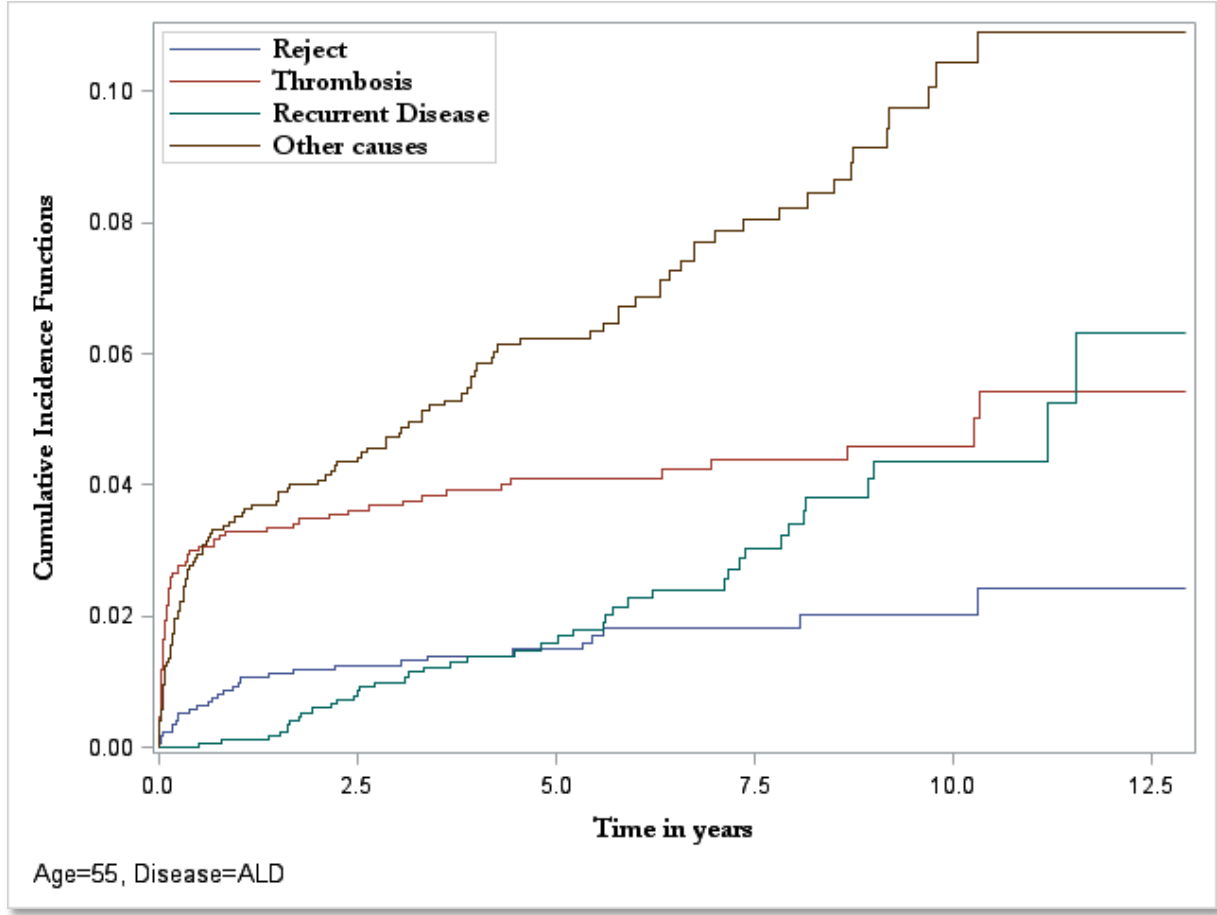
The COVAR data set (Table 4) has 8 profiles  $\mathbf{z}_0$  for which estimates of the cumulative cause-specific hazard  $\mathcal{A}_j(t|\mathbf{z}_0)$  are requested. For computing the CIF estimator  $\hat{F}_j(t|\mathbf{z}_0) = \int_0^t \hat{S}(u|\mathbf{z}_0) d\hat{\mathcal{A}}_j(u|\mathbf{z}_0)$  we will use  $\hat{S}(t|\mathbf{z}_0) = \prod_{j=1}^4 \exp(-\hat{\mathcal{A}}_j(t|\mathbf{z}_0))$ . The data set **stats\_ph** is vertically stacked with variables CUMHAZ for  $\hat{\mathcal{A}}_j(t|\mathbf{z}_0)$  and SURVIVAL for  $\exp(-\hat{\mathcal{A}}_j(t|\mathbf{z}_0))$  at the failure times by cause in each stratum.

Table 4: Covariate profiles data set COVAR					
Obs	disease	Stratum_s	age	I3	I4
1	ALD	01	55	0	0
2	ALD	02	55	0	0
3	PBC	03	55	1	0
4	PSC	03	55	1	0
5	ALD	03	55	1	0
6	PBC	04	55	0	1
7	PSC	04	55	0	1
8	ALD	04	55	0	1

Focus on the CIFs for a patient, age=55 and disease=3 (ALD). We will need information from all four strata with all relevant data on the same record line for calculations. This could be achieved by creating a data set after construction of the data set JOIN. See Appendix. A SAS macro %Cuminc by Rosthøj et al (2004) is available with the initial effort placed on construction of the input data set to run the PHREG procedure. Figure 5 shows the four CIFs. They appear similar to the non-parametric estimates in figure 2. It is a consequence of our model wherein covariates affects estimates  $\hat{\mathcal{A}}_j(t|\mathbf{z}_0)$  for  $j=3, 4$  and  $\hat{S}(t|\mathbf{z}_0)$ .

```
data join;
merge stats_ph(where=(stratum='01') rename=(survival=surv01
cumhaz=cumhaz01))
stats_ph(where=(stratum='02') rename=(survival=surv02 cumhaz=cumhaz02))
stats_ph(where=(stratum='03' and disease=3) rename=(survival=surv03
cumhaz=cumhaz03))
stats_ph(where=(stratum='04' and disease=3) rename=(survival=surv04
cumhaz=cumhaz04));
by time;
keep time surv: cumhaz:;
run;
```

Figure 5: Cumulative incidence functions from multistate model for four types of graft failure



### Variance Estimation of the CIF

The expression for the asymptotic variance  $Var(\hat{F}_j(t | \mathbf{z}_0)) = Var_1(\hat{F}_j(t | \mathbf{z}_0)) + Var_2(\hat{F}_j(t | \mathbf{z}_0))$  is derived in Anderson et al (1993) as a special case of the discussion on transition probabilities. The general formulae simplify considerably in the situation of competing risks.

$Var_1(\hat{F}_j(t | \mathbf{z}_0)) = [vec_{1j}(t)]' [Var(\hat{\beta})] [vec_{1j}(t)]$ , where  $Var(\hat{\beta})$  is the estimated covariance matrix of  $\hat{\beta}$ ,

$$Var_2(\hat{F}_j(t | \mathbf{z}_0)) = \sum_{l=1}^m \int_0^t \left\{ \hat{S}(u- | \mathbf{z}_0) \right\}^2 \left\{ \delta_{lj} - \left( \hat{F}_j(t | \mathbf{z}_0) - \hat{F}_j(u | \mathbf{z}_0) \right) \right\}^2 \frac{\exp(2\mathbf{z}'_{0l}\hat{\beta})}{\{S_l^{(0)}(u, \hat{\beta})\}^2} dN_l(u).$$

$$\text{where } vec_{1j}(t) = \int_0^t \left\{ \hat{S}(u- | \mathbf{z}_0) \left( \hat{F}_j(t | \mathbf{z}_0) - \hat{F}_j(u | \mathbf{z}_0) \right) \right\} d\mathbf{W}_{00}(u) + \int_0^t \left\{ \hat{S}(u- | \mathbf{z}_0) \right\} d\mathbf{W}_{0j}(u)$$

$$\mathbf{W}_{0l}(t) = \exp(\mathbf{z}'_0 \hat{\beta}) \int_0^t (\mathbf{z}_{0l} - \bar{\mathbf{z}}_l(u)) \frac{dN_l(u)}{\{S_l^{(0)}(u, \hat{\beta})\}}, \quad l \neq 0, \quad \mathbf{W}_{00}(t) = -\sum_{l=1}^m \mathbf{W}_{0l}(t), \quad \bar{\mathbf{z}}_l(t) = \frac{S_l^{(1)}(t, \hat{\beta})}{S_l^{(0)}(t, \hat{\beta})} \text{ and}$$

$$S_l^{(0)}(t, \hat{\beta}) = \sum_{i=1}^n Y_i(u) \exp(\mathbf{z}'_i \hat{\beta}) \quad \text{and} \quad S_l^{(1)}(t, \hat{\beta}) = \sum_{i=1}^n Y_i(u) \exp(\mathbf{z}'_i \hat{\beta}) \mathbf{z}_{il}.$$

The data from cause of failure (stratum)  $l$  are the individual covariates  $\mathbf{z}_{il}$ , the at-risk indicator  $Y_{il}(t)$  and total count of failures  $N_{il}(t)$ . The fixed covariate profile  $\mathbf{z}_0$  at which computations are made could differ by cause, as in our illustration.

The ingredients for  $Var_2(\hat{F}_j(t | \mathbf{z}_0))$  can be obtained from the OUTPUT data set from PHREG. Minimally, we get  $\mathbf{z}'_{il}\hat{\boldsymbol{\beta}}$  (XBETA) and  $\hat{A}_l(t | \mathbf{z}_{il}) = \exp(\mathbf{z}'_{il}\hat{\boldsymbol{\beta}}) \int_0^t \{S_l^{(0)}(u, \hat{\boldsymbol{\beta}})\}^{-1} dN_{il}(u)$  (-LOGSURV) at the observed time  $t$ , censored or not. We must get  $S_l^{(0)}(t, \hat{\boldsymbol{\beta}})$  and  $S_l^{(1)}(t, \hat{\boldsymbol{\beta}})$  to construct  $vec_{1j}(t)$  and via PROC IML to arrive at  $Var_1(\hat{F}_j(t | \mathbf{z}_0))$ . Although the programming is pedestrian it is not for the faint-hearted. Fortunately an elegant macro %cumincV by Rosthøj et al (2004) is available for the variance calculations. When covariates are absent  $Var(\hat{F}_j(t))$  reduces to

$$Var(\hat{F}_j(t)) = \sum_{u \leq t} \left\{ \hat{F}_j(t) - \hat{F}_j(u) \right\}^2 \left\{ \hat{S}(u-) \right\}^2 \frac{\Delta N(u)}{\{Y(u)\}^2} + \sum_{u \leq t} \left\{ \hat{S}(u-) \right\}^2 \frac{\Delta N_j(u)}{\{Y(u)\}^2} - 2 \sum_{u \leq t} \left\{ \hat{F}_j(t) - \hat{F}_j(u) \right\} \left\{ \hat{S}(u-) \right\}^2 \frac{\Delta N_j(u)}{\{Y(u)\}^2}$$

This formula is given in Anderson et al (1993), page 299 and also in Aalen et al (2008), page 126. Compare with the expression (3) in section 3 used by LIFETEST.

### ***Fine-Gray model***

In the multistate formulation the CIF  $F_j(t)$  for cause  $j$ , depends on all competing risks. Fine and Gray

(1999) introduced the *sub-distribution* or *sub-hazard function*,  $\bar{h}_j(t) = -\frac{d}{dt}(\log(1 - F_j(t)))$ , or formally

$\bar{h}_j(t) = \lim_{\Delta t \downarrow 0} P[T \leq t + \Delta t, J = j | T \geq t, \text{or } (T < t, J \neq j)] / \Delta t$ . Then  $F_j(t) = 1 - \exp\left(-\int_0^t \bar{h}_j(u) du\right)$  which at first glance might be simpler than equation (1) with cause-specific hazard  $\alpha_j(t)$ . The mathematical

relationship is  $\alpha_j(t) = \left(\frac{1 - F_j(t-)}{S(t-)}\right) \bar{h}_j(t) = \left(1 + \frac{\sum_{k \neq j} F_k(t-)}{S(t-)}\right) \bar{h}_j(t)$  but beyond that  $\alpha_j(t)$  and  $\bar{h}_j(t)$  have

fundamentally different interpretations. For example, the conditioning set in the latter keeps a subject at risk at time  $t$  who has failed before  $t$  by a cause other than  $j$ . The notion of keeping subjects alive after they have died could be jarring to many (Andersen and Keiding, 2012). Nevertheless, simplicity lies in the direct relationship between  $\bar{h}_j(t)$  and  $F_j(t)$  with the feasibility of estimating proportional sub-hazard models for only the risk  $j$  of interest, given by  $\bar{h}_j(t | \mathbf{z}) = \bar{h}_{j,0}(t) \exp(\mathbf{z}'\bar{\boldsymbol{\beta}}_j)$  where  $\bar{h}_{j,0}(t)$  is an unspecified baseline sub-hazard. This model differs fundamentally from model (4)  $\alpha_j(t | \mathbf{z}) = \alpha_{j,0}(t) \exp(\mathbf{z}'\boldsymbol{\beta}_j)$ ,  $j = 1, \dots, m$ , that posits proportionality for all cause-specific hazards. The regression coefficients have very different interpretations. Proportionality for cause-specific hazards does not imply proportionality for sub-hazards, and vice-versa. Comparing coefficients from one model with the other is not recommended.

PROC PHREG offers options for analyzing the Fine-Gray model. Focus on a specific cause say  $j=1$ .

Estimation of  $\bar{\beta}_1$  is via a pseudo partial likelihood  $L(\bar{\beta}_1) = \prod_t \prod_{i=1}^n \left( \frac{\exp(\mathbf{z}'_i \bar{\beta}_1)}{\sum_{k=1}^n Y_k(t) w_k(t) \exp(\mathbf{z}'_k \bar{\beta}_1)} \right)^{\Delta N_i(t)}$

where  $N_i(t) = [T_i \leq t, J_i = 1]$ ,  $Y_i(t) = 1 - N_i(t-)$  and the weights  $w_i(t) = \frac{r_i(t) \hat{G}(t)}{\hat{G}(T_i \wedge t)}$  are constructed from

$r_i(t) = [U_i \geq T_i \wedge t]$  and the KM estimator  $\hat{G}$  of the distribution of the censoring time  $U_i$ . If  $T_i$  was censored before  $t$  then  $r_i(t) = 0$ . The subscript 1 on  $\bar{\beta}_1$  emphasizes that the model of interest is only for cause 1. The product in  $L(\bar{\beta}_1)$  is actually only over individuals who have the event by cause 1. At a failure time  $t$  by cause 1 the at-risk indicators include those who experienced an event before  $t$ , but from a cause other than cause 1. They will get a weight  $< 1$ .

For illustration consider the liver transplant graft failure data set of  $n=1761$  patients and focus on COF=4, failure due to other causes. Our sub-hazard model  $\bar{h}_4(t | \mathbf{z}) = \bar{h}_{4,0}(t) \exp(\mathbf{z}' \bar{\beta}_4)$  has the covariates **age** **gender** **Disease**.

```
proc phreg data=liver;
  class gender(ref='female') disease(ref='ALD')/param=ref;
  model time_yr*cof(0)=age gender Disease /rl eventcode=4;
format disease disease. gender gender. cof cof.;
run;
```

Optimization of the pseudo partial likelihood produces estimates and standard errors for  $\bar{\beta}_4$ . The **rl** option computes 95% CIs for the sub-hazard ratios  $\exp(\bar{\beta}_4)$  based on asymptotic normality. Results are summarized in Table 5 with some editing of rubrics in the default output to emphasize estimation of the sub-hazard. Age shows a significant effect, and disease is borderline (2 DF chi-square  $p=0.059$ ).

Table 5: Parameter Estimates from the Fine-Gray Model (COF=4)							
Parameter		Parameter Estimate	Standard Error	p-value	Sub-Hazard Ratio	95% Sub-Hazard Ratio Confidence Limits	
age		0.02780	0.01022	0.0065	1.028	1.008	1.049
gender	male	-0.08892	0.22287	0.6899	0.915	0.591	1.416
disease	PBC	-0.42420	0.27578	0.1240	0.654	0.381	1.123
disease	PSC	0.26454	0.22835	0.2467	1.303	0.833	2.038

We must resist comparisons between these estimates with those obtained from the multistate model (4). However, for the curious please see Collett (2015), Tables 12.5 and 12.6. One should also appreciate the convoluted verbal statements that would describe the estimate and 95% CI for these sub-hazards.

We now address estimation of cumulative incidence at a specified profile  $\mathbf{z}_0$ . The plug-in estimator is

$\hat{F}_4(t | \mathbf{z}_0) = 1 - \exp\left(-\exp(\mathbf{z}'_0 \hat{\beta}_4) \hat{H}_{4,0}(t)\right)$  where  $\hat{H}_{4,0}(t) = \sum_{i=1}^n \int_0^t \{S^{(0)}(u, \hat{\beta}_4)\}^{-1} w_i(u) dN_i(u)$  estimates the

cumulative sub-hazard function and  $S^{(0)}(t, \bar{\beta}_4) = \sum_{i=1}^n w_i(u) Y_i(u) \exp(\mathbf{z}'_i \bar{\beta}_4)$ . Note that all notations are from this section only despite symbolic similarities with usage in other parts of survival analyses. A confidence

interval for  $F_4(t|\mathbf{z}_0)$  is based on a monotone transformation  $g$  (CLYTPE=option) and estimate of approximate variance  $Var\left(g(\hat{F}_4(t|\mathbf{z}_0)) - g(F_4(t|\mathbf{z}_0))\right)$ . The theory employs a simulation with an optionally specified number of normally distributed variates to represent the asymptotic normal distribution of  $g(\hat{F}_4(t|\mathbf{z}_0))$ . See Martinussen and Scheike (2006) for the derivation of the distribution of  $\hat{F}_4(t|\mathbf{z}_0)$ .

Construct a COVAR data set for  $\mathbf{z}_0$  (Table 6). With age=55 years, six gender by disease profiles are created and a variable **legend** is added to identify the CIFs. The syntax below is analogous to deriving survival function estimates following estimation of a Cox model for the overall hazard (Gardiner, 2010).

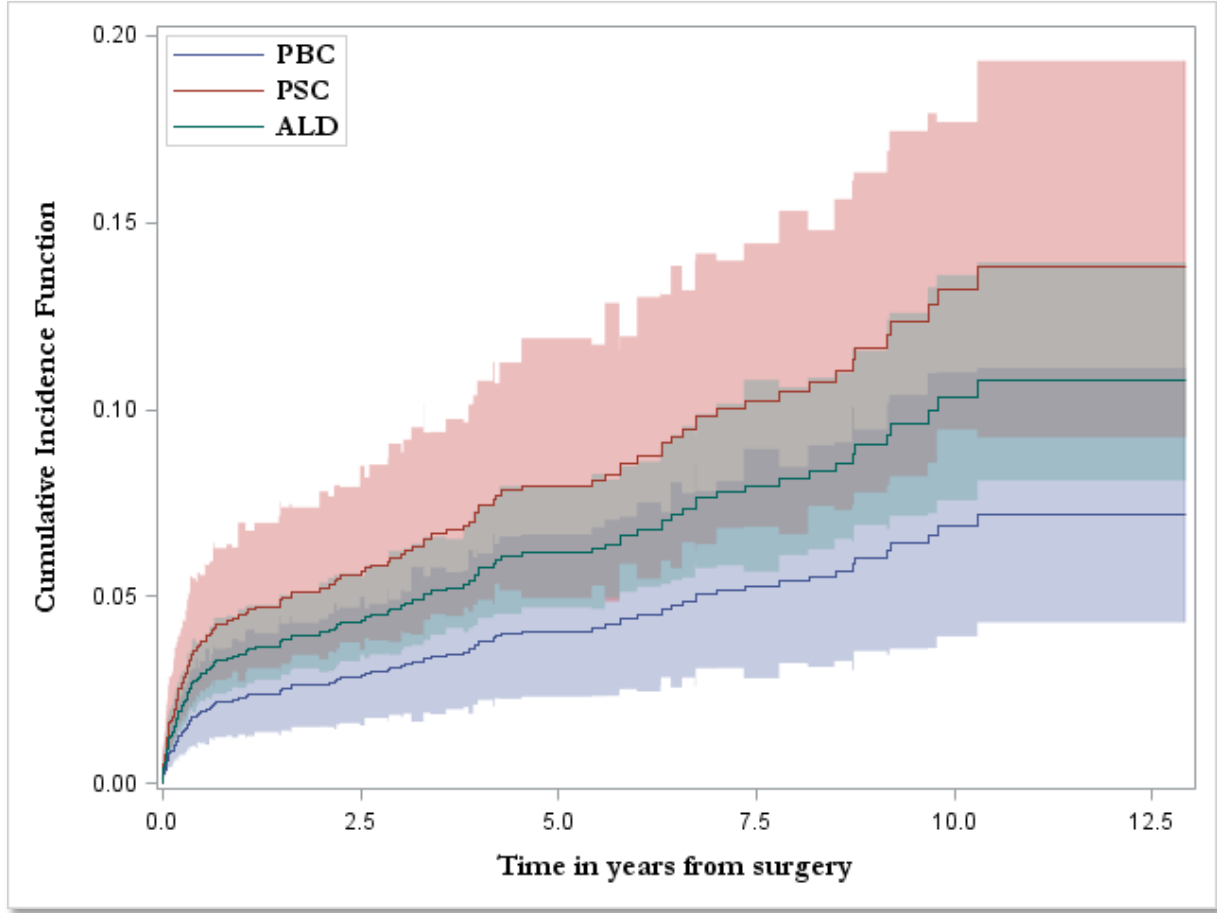
Table 6: Covariate profiles data set COVAR				
Obs	gender	disease	age	Legend
1	male	PBC	55	male -PBC
2	male	PSC	55	male -PSC
3	male	ALD	55	male -ALD
4	female	PBC	55	female-PBC
5	female	PSC	55	female-PSC
6	female	ALD	55	female-ALD

```
proc phreg data=c.liver plots(overlay=group cl)=cif;
  class gender(ref='female') disease(ref='ALD')/param=ref;
  model time_yr*cof(0)=age gender Disease / eventcode=4;
  baseline covariates=covar out=out_CIF cif=_all_ /
    rowid=legend group=gender seed=2416 cltype=loglog;
  format disease disease. gender gender. cof cof.;
  label time_yr='Time in years from surgery';
run;
```

The request `cif=_all_` in the `baseline` statement computes at each failure time due to cause 4, the CIF estimate, its standard error and 95% confidence limits saving the results in the data set `out_CIF`. Plots are created from the `plots` option. There are two sets of plots one for males and one for females: `group=gender` overlays the three CIFs by disease category and `rowid=legend` provides a label for identifying the plots. Execution time might be a few minutes on a standard desktop computer. The data set `out_cif_male2` used below subsets `out_CIF` to males only and extend the plots from the last failure time to the last observed time. The result is Figure 6.

```
%let opts=(Color=black Family=Garamond Size=12 Style=normal
Weight=Bold);
proc sgplot data=out_cif_male2;
  band x=time_yr upper=upperCIF lower=lowerCIF/type=step group=disease
  transparency=.8;
  step x=time_yr y=CIF/group=disease name="DISEASE";
  xaxis label='Time in years from surgery' LABELATTRS=&opts;
  yaxis label='Cumulative Incidence Function' LABELATTRS=&opts;
  keylegend "DISEASE" /location=inside position=topleft across=1
  VALUEATTRS=&opts;
run;
```

**Figure 6: Cumulative incidence functions for male gender and disease category from the Fine-Gray model for failure by other causes**



### *Parametric models*

We begin with the construction of the likelihood function in terms of the cause-specific hazards  $\{\alpha_j(t | \mathbf{z}) : j = 1, \dots, m\}$  and survival function  $S(t | \mathbf{z})$ . Data are comprised of a sample  $\{(T_i, J_i, \mathbf{z}_i) : 1 \leq i \leq n\}$  on event times and their cause, if observed or only right censored times (without cause). Construction of the likelihood is straight-forward. Conditionally on the covariates  $\{\mathbf{z}_i : 1 \leq i \leq n\}$ , for observed failure times  $C_j = \{T_i = t_i : J_i = j\}$  from cause  $j$ , the contribution to the likelihood is  $\prod_{C_j} \alpha_j(t_i | \mathbf{z}_i) S(t_i | \mathbf{z}_i)$ . For right-censored observations: the contribution to the likelihood is  $\prod_C S(t_i | \mathbf{z}_i)$  where  $C$  is the set of censored observations. Noting that  $S(t | \mathbf{z}) = \exp\left(-\int_0^t \sum_{j=1}^m \alpha_j(u | \mathbf{z}) du\right) = \prod_{j=1}^m \exp\left(-\int_0^t \alpha_j(u | \mathbf{z}) du\right)$ , the full likelihood is  $L = \left(\prod_{j=1}^m \prod_{C_j} \alpha_j(t_i | \mathbf{z}_i) S(t_i | \mathbf{z}_i)\right) \left(\prod_C S(t_i | \mathbf{z}_i)\right) = \prod_{j=1}^m L_j$  where

$$L_j = \prod_{C_j} \alpha_j(t_i | \mathbf{z}_i) \left\{ \prod_{i=1}^n \exp\left(-\int_0^{t_i} \alpha_j(u | \mathbf{z}_i) du\right) \right\}.$$

To complete the specification we must assume a parametric form for the cause-specific hazard, for example the Weibull form  $\alpha_j(t|\mathbf{z}) = \gamma_j \theta_j^{-\gamma_j} t^{\gamma_j-1}$ ,  $\log \theta_j(\mathbf{z}) = \mathbf{z}'\boldsymbol{\beta}_j$ . Therefore we have returned to the proportional cause-specific hazards model (4)  $\alpha_j(t|\mathbf{z}) = \gamma_j t^{\gamma_j-1} \exp(-\gamma_j \mathbf{z}'\boldsymbol{\beta}_j)$  with a parametrically specified baseline. If the model parameters are separable,  $(\gamma_j, \boldsymbol{\beta}_j)$  appear in the likelihood  $L_j$  only where failure is of type  $j$  and all other observations are regarded as right-censored. Each  $L_j$  may be optimized separately. If some parameters are shared, for example the shape  $\gamma_j = \gamma$  for all  $j$ , then the full likelihood  $L$  must be optimized. In general, a necessary and sufficient condition for independence of failure time  $T$  and cause of failure  $J$  is that the ratios  $\alpha_j(t|\mathbf{z})/\alpha_k(t|\mathbf{z})$  do not depend on  $t$  for all  $j, k$  (Crowder, 2012). The condition holds for the Weibull form with a common shape parameter.

## Example 2

In a lung cancer clinical trial of two treatments (TREAT, A or B), Lagakos (1978) describes data on follow up times (TIME, years) in 194 patients with squamous cell carcinoma. There were two competing causes of death: 83 patients died from local spread of disease (STATUS=1); 44 from metastatic disease (STATUS=2), and 67 patients were alive at last follow up (STATUS=0). Activity performance index (PERF, ambulatory or non-ambulatory) an age in years (AGE, continuous) are recorded at baseline. The data are analyzed in Crowder (2012), Chapter 13. The following formats might be helpful:

```
proc format;
value treat 0='A' 1='B';
value perf 0='AMB' 1='Non-AMB';
run;
```

Assemble a stacked data set **lungcancer2** from the original data set **lungcancer** in exactly the same manner as **liver2** was created in Example 1.

```
data lungcancer2;
set lungcancer;
do i=1 to 2;
stratum=compress('0'|vvalue(i));
fstatus=(status=i);
output;
end;
run;
```

The competing risks model with cause-specific hazards in Weibull form is estimated with PROC LIFEREG with sorting of the input file **by stratum descending treat descending perf;**

```
ods output parameterestimates=parms;
proc lifereg data=lungcancer2 order=data;
by stratum;
class perf treat;
format perf perf. treat treat.;
model time*fstatus(0)=perf treat age/dist=weibull;
run;
```



PROC LIFEREG has estimated an accelerated failure time model (Klein and Moeschberger, 2003) for each failure type, regarding all other observations are right censored. Results are shown in the left hand side columns in Table 7. Since the scale parameters  $\gamma_j^{-1}, j = 1, 2$  appear to be equal a constrained model with common scale parameter is fitted. The 1 DF likelihood ratio test for common scale is not significant. Results are shown in the middle columns of Table 7.

Table 7: Competing risks analysis for lung cancer data										
		Separable parameter model			Constrained model†			Reduced model		
Stratum	Parm	Est	StdErr	p-value	Est	StdErr	p-value	Est	StdErr	p-value
01	Intercept	0.8339	0.5128	<.0001	0.8446	0.5029	<.0001	1.3176	0.1073	<.0001
01	Perf, Non AMB	-0.4679	0.1720	0.0065	-0.4670	0.1692	0.0058	-0.3791	0.1654	0.0219
01	Treat, B	-0.3175	0.2053	0.1220	-0.3182	0.2019	0.1151			
01	age	0.0129	0.0085	0.1316	0.0126	0.0084	0.1312			
01	Scale	0.7422	0.0585							
02	Intercept	2.5773	0.7307	<.0001	2.6194	0.7517	<.0001	1.7865	0.1529	<.0001
02	Perf, Non AMB	-0.4245	0.2230	0.0570	-0.4261	0.2304	0.0644	-0.3774	0.2270	0.0965
02	Treat, B	-0.3218	0.2709	0.2348	-0.3217	0.2800	0.2506			
02	age	-0.0093	0.0114	0.4151	-0.0096	0.0118	0.4170			
02	Scale	0.7031	0.0759		0.7288	0.0464		0.7399	0.0468	

† Model has common slope

The constrained model is fitted by

```
proc lifereg data=lungcancer2 order=data;
class stratum perf treat;
format perf perf. treat treat.;
model time*fstatus(0)=stratum stratum*perf stratum*treat
stratum*age/noint dist=weibull;
run;
```

We could go a step further and test jointly the significance of TREAT and AGE by adding the **estimate** statements below. The joint 4 DF chi-square test is not significant ( $p=0.17$ ). Therefore one could argue for a simple reduced model with PERF only (Table 7 right hand side columns).

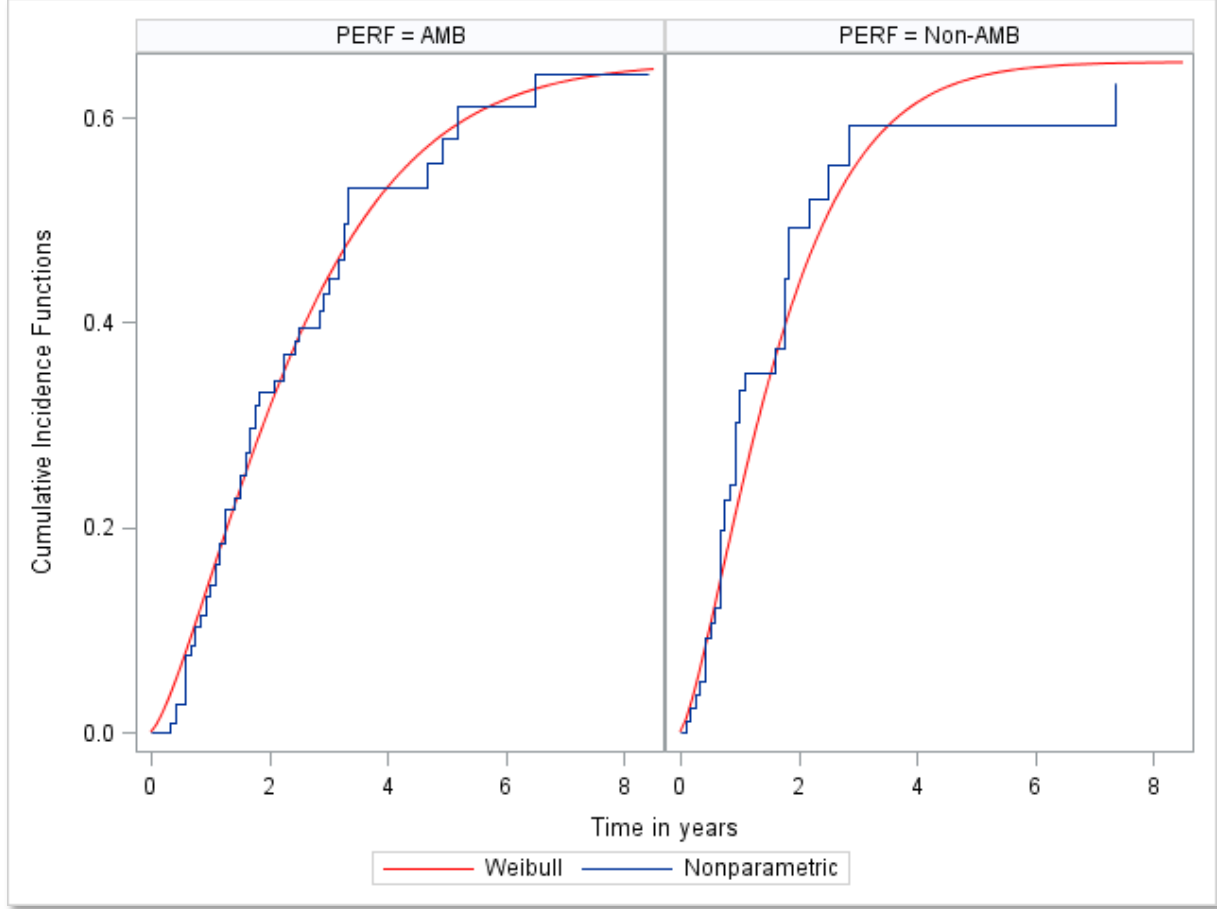
```
estimate "NO TREAT STRATUM=01" stratum*treat 1 -1,
      "NO TREAT STRATUM=02" stratum*treat 0 0 1 -1,
      "NO AGE STRATUM=01" stratum*age 1,
      "NO AGE STRATUM=02" stratum*age 0 1/e joint;
```

The common slope models imply independence of the death time and its cause. Calculation of the two CIFs is simple and are expressed closed form

$$F_j(t | \mathbf{z}_0) = \frac{K_j}{(K_1 + K_2)} \left( 1 - \exp(-(K_1 + K_2)t^\gamma) \right), j = 1, 2 \quad \text{where } K_j = \exp(-\gamma \mathbf{z}_0' \beta_j).$$

For deaths due to local spread of disease (STATUS=1) we use the reduced model (Table 7) to generate the CIFs by activity performance status. The CDF function for Weibull distribution helps this computation. We can superimpose the nonparametric estimates obtained from PROC LIFETEST. Figure 7 shows that we do quite well here.

**Figure 7: Cumulative incidence functions for local spread of disease by activity performance**



If desired an analogous analysis can be performed with a log-normal form of the cause-specific hazards:

$$\alpha_j(t) = \gamma_j t^{-1} \frac{\phi(\log(t / \theta_j)^{\gamma_j})}{\Phi(-\log(t / \theta_j)^{\gamma_j})} \text{ where } \phi \text{ and } \Phi \text{ are the probability density and cumulative distribution}$$

function of the standard normal. Covariates enter through the scale  $\theta_j(\mathbf{z}) = \exp(\mathbf{z}'\beta_j)$ . In this case we do not have independence of the death time and its cause. Closed forms for the survival distribution and CIFs are not available. If there are only two competing risks the CIF expression for failure by cause=1 is

$$F_1(t | \mathbf{z}) = \int_{-\infty}^{t^*} \Phi(-\{\gamma_2 / \gamma_1\}u + \gamma_2 \log(\theta_1(\mathbf{z}) / \theta_2(\mathbf{z}))) \phi(u) du \text{ where } t^* = \log(t / \theta_1(\mathbf{z}))^{\gamma_1}. \text{ The proportion that fail by cause=1 is } P[J=1 | \mathbf{z}] = \Phi\left(\frac{\gamma_1 \gamma_2 \log(\theta_2(\mathbf{z}) / \theta_1(\mathbf{z}))}{\sqrt{\gamma_1^2 + \gamma_2^2}}\right).$$

## 5. Concluding remarks

Our paper has discussed approaches to analyzing competing risks data using SAS software with an emphasis on three approaches to regression modeling. First, the multi-state modelling framework based on cause-specific hazards. Second, the sub-hazards approach of Fine and Gray (1999) implemented in PROC PHREG, and third the classical parametric analyses based on a fully parameterized likelihood. In addition, for non-parametric analyses of cumulative incidence functions (CIFs) options are available in PROC LIFETEST. Although current SAS software does not offer a direct analysis of CIFs from multi-state models, they can be analyzed with PROC PHREG as demonstrated in this article. Constructing the estimator of the CIF and its standard error follows from statistics generated after estimation of regression parameters in the cause-specific hazards model. Two elegant macros by Rosthøj et al (2014) provide tools for the computations.

The sub-hazards model of Fine and Gray (1999) has received considerable attention and some criticism primarily on its interpretability and violation of the principle of strict exogeneity (Andersen and Keiding, 2012). Because the sub-hazards approach and cause-specific hazards are very different, comparisons should not be made between them in empirical applications. Users should be aware that the output from PHREG describes sub-hazard ratios and confidence intervals for the sub-hazard ratios and *not* for cause-specific hazard ratios. For emphasis, in our application we edited the rubrics from the default output (see Table 5).

There has been a long history of parametric analyses for competing risks. Essentially we need to specify parametrically the joint distribution of  $(T, J)$ , the event time and its cause. Although the probability statement for two events  $P(AB) = P(A)P(B|A)$  is symmetric in  $A$  and  $B$ , the conditional probability must be interpreted carefully. In mixture models, the conditioning is on cause of failure  $J$ : so we need to assume a parametric distribution for  $T$  conditional on cause-type and a multinomial model for  $J$  (Larson and Dinse, 1985; Lau et al, 2008). The model has limitations because it conditions on a future  $J$ . For instance, describing  $P[T > t | J = j]$  may be semantically challenging. A relatively recent approach, called vertical modelling reverses the conditioning event (Nicolaie et al, 2010). Here  $P[J = j | T > t]$  has a simple interpretation. In vertical modelling, it is the relative cause-specific hazards  $\pi_j(t) = \alpha_j(t) / \alpha_0(t)$  that are modelled, for example by a multinomial model with time-dependent basis functions. Incorporating covariates is a separate issue (Nicolaie et al, 2010, 2013).

There are several excellent discussions of competing risks analyses (Klein et al, 2014; Beyersmann et al, 2012; Putter et al, 2007; Pintilie, 2006). A recent paper by Li (2016) extends the Fine-Gray model for interval censored data. Christian et al (2016) develop a frailty model for clustered competing risks data. Effraimidis and Dahl (2014) describe a nonparametric approach to CIF estimation with right censored data and missing at random cause of failure. A parametric regression model of Jeong and Fine (2006, 2007) uses a defective Gompertz distribution for the CIF. A quantile regression analysis for competing risks is discussed in Peng and Fine (2009).

Application of a methodology for analyses in competing risks data demands a careful interpretation of results. The three methods discussed in this paper are the standard approaches that can be analyzed with SAS software. Enhancements will be needed to extend their reach to the broader domain of biomedical applications.

## 6. References

- Aalen O, Borgan O, Gjessing H. *Survival and Event History Analysis*. Springer-Verlag, 2008.
- Allignol A, Schumacher M, Beyersmann J. A note on variance estimation of the Aalen-Johansen estimator of the cumulative incidence function in competing risks, with a view towards left-truncated data. *Biometrical J.* 2010; 52(1): 126-37.
- Andersen PK, Keiding N. Interpretability and importance of functionals in competing risks and multistate models. *Statistics in Med.* 2012; 31(11-12): 1074-88.
- Andersen PK, Abildstrom SZ, Rosthøj S. Competing risks as a multi-state model. *Statist Methods in Med Res.* 2002; 11(2): 203-15.
- Andersen PK, Borgan O, Gill RD, Keiding N. *Statistical Models Based on Counting Processes*. Springer-Verlag, 1993.
- Beyersmann J, Allignol A, Schumacher M. *Competing Risks and Multistate Models with R*. Springer-Verlag, 2012.
- Christian NJ, Ha ID, Jeong JH. Hierarchical likelihood inference on clustered competing risks data. *Statistics in Med.* 2016; 35(2): 251-67.
- Collett D. *Modelling Survival Data for Medical Research, 3rd edition*. CRC Press, 2015.
- Crowder M. *Multivariate Analysis and Competing Risks*. CRC Press, 2012.
- Effraïmidis G, Dahl CM. Nonparametric estimation of cumulative incidence functions for competing risks data with missing cause of failure. *Statist & Probab Letters.* 2014; 89: 1-7.
- Fine JP, Gray RJ. A proportional hazards model for the subdistribution of a competing risk. *J Amer Statist Assoc.* 1999; 94(446): 496-509.
- Gardiner JC. Survival Analysis: Overview of Parametric, Nonparametric and Semiparametric approaches and New Developments, Paper 252-2010. *SAS Global Forum*, 2010. SAS Institute Inc, Cary NC.
- Gray RJ. A class of k-sample tests for comparing the cumulative incidence of a competing risk. *Ann Statist.* 1988; 16(3): 1141-54.
- Jeong JH, Fine JP. Parametric regression on cumulative incidence function. *Biostatistics.* 2007; 8(2):184-96.
- Jeong JH, Fine JP. Direct parametric inference for the cumulative incidence function. *Applied Statistics, JRSS-C.* 2006; 55: 187-200.
- Klein JP, van Houwelingen HC, Ibrahim JG, Scheike TH (editors). *Handbook of Survival Analysis*. CRC Press, 2014.
- Klein JP, Moeschberger ML. *Survival Analysis: Techniques for Censored and Truncated Data, 2nd Edition*. Springer-Verlag, 2003.
- Lagakos SW. A covariate model for partially censored data subject to competing causes of failure. *Applied Statistics, JRSS-C.* 1978; 27(3): 235-41.
- Larson MG, Dinse GE. A mixture model for the regression-analysis of competing risks data. *Applied Statistics, JRSS-C.* 1985; 34(3): 201-11.

- Li CX. The Fine-Gray model under interval censored competing risks data. *J Mult Anal.* 2016; 143: 327-44.
- Martinussen T, Scheike TH. *Dynamic Regression Models for Survival Data.* Springer-Verlag, 2006.
- Nicolaie MA, van Houwelingen JC, de Witte TM, Putter H. Dynamic prediction by landmarking in competing risks. *Statistics in Med.* 2013; 32(12): 2031-47.
- Nicolaie MA, van Houwelingen HC, Putter H. Vertical modeling: A pattern mixture approach for competing risks modeling. *Statistics in Med.* 2010; 29(11): 1190-205.
- Peng LM, Fine JP. Competing risks quantile regression. *J Amer Statist Assoc.* 2009; 104(488): 1440-53.
- Pintilie M. *Competing Risks, A Practical Perspective.* John Wiley & Sons, 2006.
- Putter H, Fiocco M, Geskus RB. Tutorial in biostatistics: Competing risks and multi-state models. *Statistics in Med.* 2007; 26(11): 2389-430.
- Rosthøj S, Andersen PK, Abildstrom SZ. SAS macros for estimation of the cumulative incidence functions based on a Cox regression model for competing risks survival data. *Computer Methods and Programs in Biomedicine.* 2004; 74(1): 69-75.

## CONTACT INFORMATION

We welcome your comments and questions. Please contact

Joseph C. Gardiner  
Department of Epidemiology and Biostatistics  
Michigan State University  
East Lansing, MI 48824  
[jgardiner@epi.msu.edu](mailto:jgardiner@epi.msu.edu)

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.

## Appendix

The cumulative incidence functions are estimated from the multistate model using the data set from the BASELINE statement in PHREG. Begin with the JOIN data set that merges by failure time the cumulative hazard estimates  $\hat{A}_{0j}(t | \mathbf{z}_0)$  by each cause of failure  $j=1, \dots, m$ . This data set will show several ‘missing’ values because practically cause-specific failures times will have very few ties. Therefore for cosmetic reasons, JOIN2 will flesh out the estimates at all failure times. One additional record is added to extend all estimates to the last observed time (if censored) rather than ending at the last failure time. This occurs in the example for graft failure due to recurrent disease. The change was made to generate figure 5.

```

data join2;
set join end=last;
array s{4} surv01-surv04;
array c{4} cumhaz01-cumhaz04;
array s_{4} s_1-s_4 (4*1);
array c_{4} c_1-c_4 (4*0);
do i=1 to 4;
if s{i}=. then s{i}=s_{i};
else s_{i}=s{i};
if c{i}=. then c{i}=c_{i};
else c_{i}=c{i};
end;
P00=surv01*surv02*surv03*surv04; /*---overall survival---*/
output;
if last then do;
time=4720; /*-- extended estimates to last observed time --*/
output;
end;
drop i s_ : c_ ;
run;

```

Data set JOIN3 performs the calculation of the CIFs  $F_j(t|z_0), j=1, \dots, 4$  named P01, P02, P03, P04.

```

data join3;
set join2;
array c{4} cumhaz01-cumhaz04;
P00_=max(0, lag(P00));
array a_{4} a_1-a_4 (4*0);
array P{4} P01-P04 (4*0);
do i=1 to 4;
a_{i}= max(0, lag(c{i}));
P{i}=P00_*(c{i}-a_{i});
time_yr=time/365.25;
end;
drop i a_ ;
run;

```

Figure 5 is produced by

```

%let opts=(Color=black Family=Garamond Size=11 Style=normal
Weight=Bold);
proc sgplot data=join3;
step x=time_yr y=P01/name="P01" legendlabel="Reject";
step x=time_yr y=P02/name="P02" legendlabel="Thrombosis";
step x=time_yr y=P03/name="P03" legendlabel="Recurrent Disease";
step x=time_yr y=P04/name="P04" legendlabel="Other causes";
yaxis label='Cumulative Incidence Functions' LABELATTRS=&opts;
keylegend "P01" "P02" "P03" "P04"/location=inside across=1
VALUEATTRS=&opts;
xaxis label='Time in years' LABELATTRS=&opts;
footnote justify=left 'Age=55, Disease=ALD';run;

```