

Further documentation on code and methods in the  
survival package

Terry Therneau

March 2024

# Chapter 1

## Introduction

For many years I have used the `noweb` package to intersperse technical documentation with the source code for the `survival` package. However, despite its advantages, the uptake of `noweb` by the R community in general has been nearly nil. It is increasingly clear that no future maintainer will continue the work — on the github page I have yet to receive a suggested update that actually fixed the `.Rmd` source file instead of the `.R` file derived from it. This means that I can't merge a suggested change automatically, but have to replicate it myself.

This document is a start at addressing this. As routines undergo maintenance, I will remove the relevant `.Rnw` file in the `noweb` directory and work directly on the C and R code, migrating the “extra” material into this document. In this vignette are discussions of design issues, algorithms, and detailed formulas. In the `.R` and `.c` code I will place comments of the form “See methods document, `abc:def`”, adding an `abc:def` entry to the index of this document. I am essentially splitting each `noweb` document into two parts. The three advantages are

- Ease the transition to community involvement and maintenance.
- The methods document will be a more ready source of documentation for those who want to know technical details, but are not currently modifying the code.
- (minor) I expect it will be useful to have the methods document and R or C code open simultaneously in 2 windows, when editing.

However, this conversion will be a long process.

**Notation** Throughout this document I will use formal counting process notation, so you may as well get used to it. The primary advantage is that it is completely precise, and part of the goal for this document is to give a fully accurate description of what is computed. In this notation we have

- $Y_i(t) = 1$  if observation  $i$  is at risk at time  $t$ , 0 otherwise.
- $N_i(t) =$  total number of events, up to time  $t$  for subject  $i$ .
- $dN_i(t) = 1$  if there was an event at exactly time  $t$

- $X$  = the matrix of covariates, with  $n$  rows, one per observation, and a column per covariate.
- $X(s)$  = the time-dependent covariate matrix, if there are time-dependent covariates

For multistate models the extends to  $Y_{ij}(t)$ , which is 1 if subject  $i$  is at risk and in state  $j$  at time  $t$ , and  $N_{ijk}(t)$  which is the number of  $j$  to  $k$  transitions that have occurred, for subject  $i$ , up to time  $t$ . The number at risk and number of events at time  $t$  can be written as

$$\begin{aligned} n(t) &= \sum_{i=1}^n w_i Y_i(t) \\ &= \bar{Y}(t) \\ d(t) &= \sum_{i=1}^n w_i dN_i(t) \\ &= d\bar{N}(t) \end{aligned}$$

where  $w_i$  are optional weights for each observation.  $\bar{N}(t)$  is the cumulative number of events up to time  $t$ . I will often also use  $r_i(t) = \exp(X_i(t)\beta)$  as the per observation risk score in a Cox model.

## Chapter 2

# Survival Curves

The `survfit` function was set up as a method so that we could apply the function to both formulas (to compute the Kaplan-Meier) and to `coxph` objects. The downside to this is that the manual pages get a little odd: `survfit` is generic and `survfit.formula` and `survfit.coxph` are the ones a user will want. But from a programming perspective it has mostly been a good idea. At one time, long long ago, we allowed the function to be called with “Surv(time, status)” as the formula, i.e., without a right hand side of  $\sim 1$ . That was a bad idea, now abandoned: `survfit.Surv` is now a simple stub that prints an error message.

### 2.1 Roundoff and tied times

One of the things that drove me nuts was the problem of “tied but not quite tied” times. It is a particular issue for both survival curves and the Cox model, as both treat a tied time differently than two close but untied values. As an example consider two values of  $24173 = 23805 + 368$ . These are values from an actual study with times in days: enrollment at age 23805 days and then 368 days of follow-up. However, the user chose to use age in years, and saved those values out in a CSV file, the left hand side of the above equation becomes 66.18206708000000 and the right hand side addition yeilds 66.18206708000001. The R phrase `unique(x)` sees these two values as distinct but `table(x)` and `tapply` see it as a single value since they first apply `factor` to the values, and that in turn uses `as.character`. A transition through CSV is not necessary to create the problem. Consider the small code chunk below. For someone born on 1960-03-10, it caclulates the time interval between a study enrollment on each date from 2010-01-01 to 2010-07-29 (200 unique dates) and a follow up that is exactly 29 days later, but doing so on age scale.

```
> tfun <- function(start, gap, birth= as.Date("1960-01-01")) {
  as.numeric(start-birth)/365.25 - as.numeric((start + gap)-birth)/365.25
}
> test <- logical(200)
> for (i in 1:200) {
  test[i] <- tfun(as.Date("2010/01/01"), 29) ==
    tfun(as.Date("2010/01/01") + i, 29)
```

```

    }
  > table(test)
test
FALSE TRUE
  148   52

```

The number of FALSE entries in the table depends on machine, compiler, and possibly several other issues. There is discussion of this general issue in the R FAQ: “why doesn’t R think these numbers are equal”. The Kaplan-Meier and Cox model both pay careful attention to ties, and so both now use the `aeqSurv` routine to first preprocess the time data. It uses the same rules as `all.equal` to adjudicate ties and near ties. See the vignette on tied times for more detail.

The `survfit` routine has been rewritten more times than any other in the package, as we trade off simplicity of the code with execution speed. The current version does all of the organizational work in S and calls a C routine for each separate curve. The first code did everything in C but was too hard to maintain and the most recent prior function did nearly everything in S. Introduction of robust variance prompted a movement of more of the code into C since that calculation is computationally intensive.

The `survfit.formula` routine does a number of data checks, then hands the actual work off to one of three computational routes: simple survival curves using the `survfitKM.R` and `survfitkm.c` functions, interval censored data to `survfitTurnbull.R`, and multi-state curves using the `survfitAJ.R` and `survfitaj.c` pair.

The addition of `+ cluster(id)` to the formula was the suggested form at one time, we now prefer `id` as a separate argument. Due to the long comet’s tail of usage we are unlikely to formally depreciate the older form any time soon. The `istate` argument applies to multi-state models (Aalen-Johansen), or any data set with multiple rows per subject; but the code refrains from complaint if it is present when not needed.

## 2.2 Single outcome

We begin with classic survival, where there is a single outcome. Multistate models will follow in a separate section.

At each event time we have

- $n(t)$  = weighted number at risk
- $d(t)$  = weighted number of events
- $e(t)$  = unweighted number of events

leading to the Nelson-Aalen estimate of cumulative hazard and the Kaplan-Meier estimate of survival, and the estimate of hazard at each time point.

$$\begin{aligned}
 KM(t) &= KM(t-)(1 - d(t)/n(t) \\
 NA(t) &= NA(t-) + d(t)/n(t)h(t) &= d(t)/n(t)
 \end{aligned}$$

An alternative estimate in the case of tied times is the Fleming-Harrington. When there are no case weights the FH idea is quite simple. The assumption is that the real data is not tied,

but we saw a coarsened version. If we see 3 events out of 10 subjects at risk the NA increment is  $3/10$ , but the FH is  $1/10 + 1/9 + 1/8$ , it is what we would have seen with the uncoarsened data. If there are case weights we give each of the 3 terms a  $1/3$  chance of being the first, second, or third event

$$\begin{aligned} KM(t) &= KM(t-)(1 - d(t)/n(t)) \\ NA(t) &= NA(t-) + d(t)/n(t) \\ FH(t) &= FH(t-) + \sum_{i=1}^3 \frac{(d(t)/3)}{n(t) - d(t)(i-1)/3} \end{aligned}$$

If we think of the size of the denominator as a random variable  $Z$ , an exact solution would use  $E(1/Z)$ , the FH uses  $1/E(Z)$  and the NA uses  $1/\max(Z)$  as the denominator for each of the 3 deaths. Although Fleming and Harrington were able to show that the FH has a lower MSE than the KM, it little used. The primary reasons for its inclusion in the package are first, that I was collaborating with them at the time so knew of these results, but second and more importantly, this is identical to the argument for the Efron approximation in a Cox model. (The NA hazard corresponds to the Breslow approximation.) The Efron approx is the default for the `coxph` function, so post `coxph` survival curves need to deal with it.

When one of these 3 subjects has an event but continues to be at risk, which can happen with start/stop data, then the argument gets trickier. Say that the first of the 3 continues, the others do not. We can argue that subject 1 remains at risk for all 3 denominators, or not. The first is more a mathematical viewpoint, the second more medical, e.g., in a study with repeated infections you will never have a second one recorded on the same day. For multi-state, we finally “tossed in the towel” and now use the Breslow approx as default for multi-state hazard (`coxph`) models. For single state, the FH estimate is rarely requested but we still do our best to handle all aspects of it (pride and history), but there would be few tears if it were dropped.

When there is (time1, time2) data the code uses a `position` vector, explained further below which is 1 if this is obs is the leftmost for a subset and 2 if it is the rightmost, 3 if it is both. A primary purpose was to not count a subject as an extra entry and censor in the middle of a string of times such as (0, 10), (20, 35), (35, 50), but we also use it to moderate the FH estimate: only those with an event at the end of their intervals participate in the special computation for ties.

The `survfitKM` call has arguments of

- `y`: a matrix `y` containing survival times, either 2 or 3 columns
- `weight`: vector of case weights
- `ctype`: computation for the cumulative hazard: either the Nelson-Aalen (1) or Efron approximation for ties (2) approach. Number 2 is very rarely used.
- `stype`: computation for the survival curve: either product-limit (1), also known as the Kaplan-Meier or `exp(-cumulative hazard)` (2). Use of `ctype=2, stype=2` matches a Cox model using the Efron approximation for ties, `ctype=1, stype=2` is the Fleming-Harrington estimate of survival or a Cox model with the Breslow approximation for ties.
- `type`: older form of the `ctype/stype` argument, retained for backwards compatibility. Type 1 = (`ctype` 1/ `stype` 1), 2 = (2, 1), 3 = (1,2), 4 = (2,2).

- `id`: subject id, used for data checks when there are multiple rows per subject
- `cluster`: clustering used for the robust variance. Clustering is based on `id` if `cluster` is missing (the usual case)
- `influence`: return the influence matrix for the survival
- `start.time`: optional starting time for the curve
- `entry`: return entry times for (time1, time2) data
- `time0`: include time 0 in the output
- `se.fit`: compute the standard errors
- `conf.type`, `conf.int`, `conf.lower`: options for confidence intervals

### 2.2.1 Confidence intervals

If  $p$  is the survival probability and  $s(p)$  its standard error, we can do confidence intervals on the simple scale of  $p \pm 1.96s(p)$ , but that does not have very good properties. Instead use a transformation  $y = f(p)$  for which the standard error is  $s(p)f'(p)$ , leading to the confidence interval

$$f^{-1}(f(p) + -1.96s(p)f'(p))$$

Here are the supported transformations.

	$f$	$f'$	$f^{-1}$
<code>log</code>	$\log(p)$	$1/p$	$\exp(y)$
<code>log-log</code>	$\log(-\log(p))$	$1/[p \log(p)]$	$\exp(-\exp(y))$
<code>logit</code>	$\log(p/1-p)$	$1/[p(1-p)]$	$1 - 1/[1 + \exp(y)]$
<code>arcsin</code>	$\arcsin(\sqrt{p})$	$1/(2\sqrt{p(1-p)})$	$\sin^2(y)$

Plain intervals can give limits outside of (0,1), we truncate them when this happens. The log intervals can give an upper limit greater than 1 (rare), again truncated to be  $\leq 1$ ; the lower limit is always valid. The log-log and logit are always valid. The arcsin requires some fiddling at 0 and  $\pi/2$  due to how the R function is defined, but that is only a nuisance and not a real flaw in the math. In practice, all the intervals except plain appear to work well. In all cases we return NA as the CI for survival=0: it makes the graphs look better.

Some of the underlying routines compute the standard error of  $p$  and some the standard error of  $\log(p)$ . The `selow` argument is used for the modified lower limits of Dory and Korn. When this is used for cumulative hazards the `ulimit` arg will be FALSE: we don't want to impose an upper limit of 1.

### 2.2.2 Robust variance

For an ordinary Kapan-Meier curve it can be shown that the infinitesimal jackknife (IJ) variance is identical to the Greenwood estimate, so the extra computational burden of a robust estimate is unnecessary. The proof does not carry over to curves with (time1, time2) data. The use of (time1, time2) data for a single outcomes arises in 3 cases, with illustrations shown below.

- Delayed entry for a subset of subjects. The robust variance in this case is not identical to the Greenwood formula, but is very close in all the situations I have seen.
- Multiple events of the same type, as arises in reliability data. In this case the cumulative hazard rather than  $P(\text{state})$  is the estimate of most interest; it is known as the mean cumulative function (MCF) in the reliability literature. If there are multiple events per subject (the usual) an id variable is required. A study with repeated glycemc events as the endpoint (not shown) had a mean of 48 events per subject over 5.5 months. In this case accounting for within-subject correlation was crucial; the ratio of robust/asymptotic standard error was 4.4. For the valveSeat data shown below the difference is minor.
- Survival curves for a time-dependent covariate, the so called “extended Kaplan-Meier” [2]. I have strong statistical reservations about this method.

Here are examples of each. In the myeloma data set the desired time scale was time since diagnosis, and some of the patients in the study had been diagnosed at another institution before referral to the study institution.

```

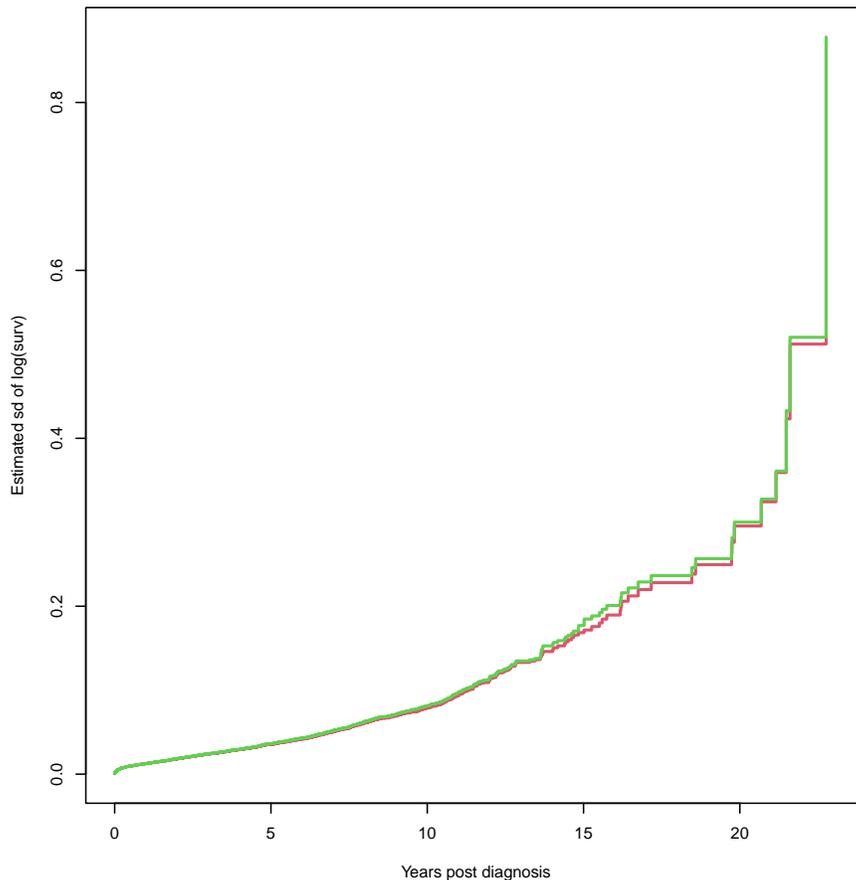
> fit1a <- survfit(Surv(entry, futime, death) ~ 1, myeloma)
> fit1b <- survfit(Surv(entry, futime, death) ~ 1, myeloma, id=id, robust=TRUE)
> matplot(fit1a$time/365.25, cbind(fit1a$std.err, fit1b$std.err/fit1b$surv),
          type='s', lwd=2, lty=1, col=2:3, #ylim=c(0, .6),
          xlab="Years post diagnosis", ylab="Estimated sd of log(surv)")
> #
> # when two valve seats failed at the same inspection, we need to jitter one
> # of the times, to avoid a (time1, time2) interval of length 0
> ties <- which(with(valveSeat, diff(id)==0 & diff(time)==0)) #first of a tie
> temp <- valveSeat$time
> temp[ties] <- temp[ties] - .1
> vdata <- valveSeat
> vdata$time1 <- ifelse(!duplicated(vdata$id), 0, c(0, temp[-length(temp)]))
> vdata$time2 <- temp
> fit2a <- survfit(Surv(time1, time2, status) ~1, vdata)
> fit2b <- survfit(Surv(time1, time2, status) ~1, vdata, id=id)
> plot(fit2a, cumhaz=TRUE, xscale=365.25, xlab="Years in service",
       ylab="Estimated number of repairs")
> lines(fit2b, cumhaz=TRUE, lty=c(1,3,3))
> legend(150, 1.5, c("Estimate", "asymptotic se", "robust se"), lty=1:3, bty='n')
> #
> # PBC data, categorized by most recent bilirubin
> # as an example of the EKM
> pdata <- tmerge(subset(pbcseq, !duplicated(id), c(id, trt, age, sex, stage)),
                 subset(pbcseq, !duplicated(id, fromLast=TRUE)), id,
                 death= event(futime, status==2))
> bcut <- cut(pbcseq$bili, c(0, 1.1, 5, 100), c('normal', 'moderate', 'high'))
> pdata <- tmerge(pdata, pbcseq, id, cbili = tdc(day, bcut))
> pdata$ibili <- pdata$cbili[match(pdata$id, pdata$id)] # initial bilirubin
> ekm <- survfit(Surv(tstart, tstop, death) ~ cbili, pdata, id=id)

```

```

> km <- survfit(Surv(tstart, tstop, death) ~ ibili, pdata, id=id)
> plot(ekm, fun='event', xscale=365.25, lwd=2, col=1:3, conf.int=TRUE,
      lty=2, conf.time=c(4,8,12)*365.25,
      xlab="Years post enrollment", ylab="Death")
> lines(km, fun='event', lwd=1, col=1:3, lty=1)
> #       conf.time= c(4.1, 8.1, 12.1)*365.25)
> text(c(4600, 4300, 2600), c(.23, .56, .78), c("Normal", "Moderate", "High"),
      col=1:3, adj=0)
> legend("topleft", c("KM", "EKM"), lty=1:2, col=1, lwd=2, bty='n')

```



The EKM plot is complex: dashed are the EKM along with confidence intervals, solid lines are the KM stratified by enrollment bilirubin. The EKM bias for normal and moderate is large, but confidence intervals differ little between the robust and asymptotic se (latter not shown).

As shown above, more often than not the robust IJ variance is not needed for the KM curves themselves. However, they are the central computation for psuedovalues, which are steadily increasing in popularity. Let  $U_k(t)$  be the IJ for observation  $k$  at time  $t$ . This is a vector, but in section ?? for multi-state data it will be a matrix with 1 column per state. Likewise let  $C(t)$  and  $A(t)$  be influence vectors for the cumulative hazard and the area under the curve at time  $t$ .

Let  $h(t)$  be the hazard increment at time  $t$ . Then

$$h(t) = \frac{\sum_i w_i dN_i(t)}{\sum_i w_i Y_i(t)} \quad (2.1)$$

$$\begin{aligned} \frac{\partial h(t)}{\partial w_k} &= \frac{dN_i(t) - Y_k(t)h(t)}{\sum_i w_i Y_i(t)} \\ &= C_k(t) - C_k(t-) \end{aligned} \quad (2.2)$$

$$\begin{aligned} S(t) &= \prod_{s \leq t} 1 - h(s) \\ &= S(t-)[1 - h(t)] \end{aligned} \quad (2.3)$$

$$U_k(t) = \frac{\partial S(t)}{\partial w_k} \quad (2.4)$$

$$= U_k(t-)[1 - h(t)] - S(t-) \frac{dN_i(t) - Y_k(t)h(t)}{\sum_i w_i Y_i(t)} \quad (2.5)$$

$$\begin{aligned} U_k(t) &= \frac{\partial \prod_{s \leq t} 1 - h(s)}{\partial w_k} \\ &= \sum_{s \leq t} \frac{S(t)}{1 - h(s)} \frac{dN_i(s) - Y_k(s)h(s)}{\sum_i w_i Y_i(s)} \\ &= S(t) \sum_{s \leq t} \frac{dN_i(s) - Y_k(s)h(s)}{[1 - h(s)] \sum_i w_i Y_i(s)} \end{aligned} \quad (2.6)$$

$$= S(t) \int_0^t \frac{\partial \log[1 - h_k(s)]}{w_k} d\bar{N}(s) \quad (2.7)$$

The simplest case is  $C_k(\tau)$  at a single user requested reporting time  $\tau$ , i.e., a simple use of the pseudo routine. The obvious code will update all  $n$  elements of  $C$  at each event time, an  $O(nd)$  computation where  $d$  is the number of unique event times  $\leq \tau$ . For most data  $d$  and  $n$  grow together, and  $O(n^2)$  computations are something we want to avoid at all costs, a large data set will essentially freeze the computer.

The code solution is to create a running total  $z$  of the right hand term of (2.2), which applies to all subjects at risk. Then the leverage for an observation at risk over the interval  $(a, b)$  is

$$\begin{aligned} z(t) &= \sum_{s \leq t} \frac{h(s)}{\sum_i w_i Y_i(s)} \\ C_k(\tau) &= \text{frac} N_k(\tau) \sum_i w_i Y_i(b) + z(\min(a, \tau)) - z(\min(b, \tau)) \end{aligned}$$

In R code this becomes an indexing problem. We can create  $\mathbf{3}$   $n$  by  $m =$  number of reporting times matrices that point to the correct elements of the vectors  $\mathbf{c}(0, 1/\mathbf{n.risk})$  and  $\mathbf{c}(0, \mathbf{n.event}/\mathbf{n.risk}^2)$  obtained from the `survfit` object.

The computation for survival starts out exactly the same if using the exponential estimate  $S(t) = \exp(-\Lambda(t))$ , otherwise do the same computation but using the rightmost term of equation (2.6). At the end multiply the column for reporting time  $\tau$  by  $-S(\tau)$ , for each  $\tau$ .

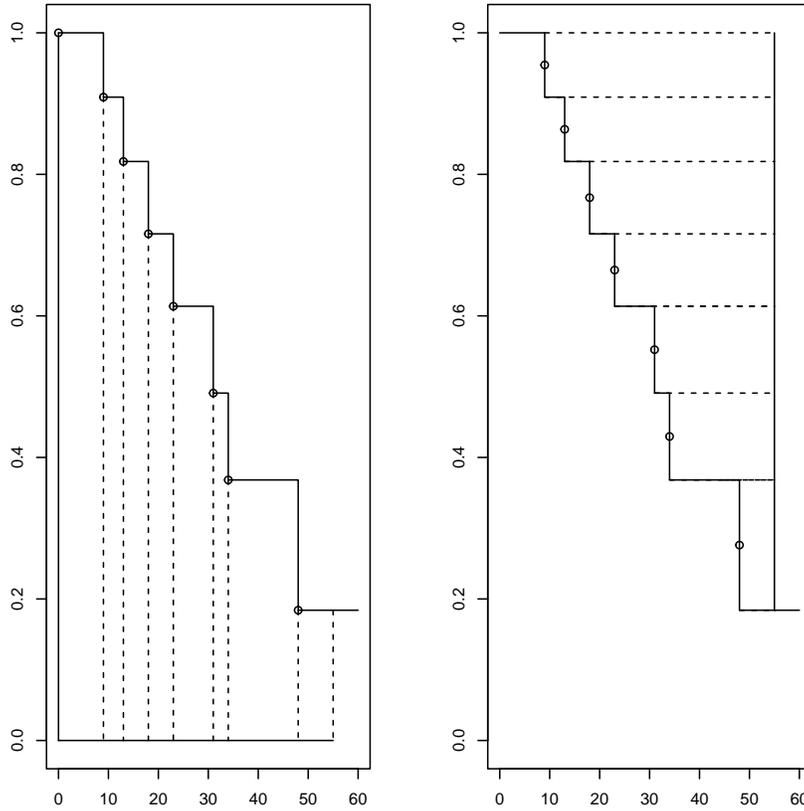


Figure 2.1: Two ways of looking at the AUC

There are two ways to look at the influence on the sojourn time, which is calculated as the area under the curve. They are shown in the figure 2.1 for an AUC at time 55.

The first uses the obvious rectangle rule. The leverage of observation  $i$  on  $\text{AUC}(55)$  is the sum of the leverage of observation  $i$  on the height of each of the circles times the width of the associated rectangle. The leverage on the height of each circle are the elements of  $U$ . The second figure depends on the fact that the influence on the AUC must be the same as the influence on 55-AUC. It will be the influence of each observation on the length of each circled segment, times the distance to 55. This is closely related to the asymptotic method used for the ordinary KM, which assumes that the increments are uncorrelated.

Using the first approach, let  $w_0, w_1, \dots, w_m$  be the widths of the  $m + 1$  rectangles under the survival curve  $S(t)$  from 0 to  $\tau$ , and  $d_1, \dots, d_m$  the set of event times that are  $\leq \tau$ . Further,

let  $u_{ij}$  be the additive elements that make up  $U$  as found in the prior formulas, i.e.,

$$\begin{aligned}
U_{ik} &= S(d_k) \sum_{j=1}^k u_{ij} \\
u_{ij} &= \frac{\partial \log(1 - h(d_j))}{\partial w_i} \text{ KM} \\
&= - \frac{dN_i(d_j) - Y_i(d_j)h(d_j)}{[1 - h(d_j)]\bar{Y}(d_j)} \\
u_{ij}^* &= \frac{\partial - h(d_j)}{\partial w_i} \exp(\text{chaz}) \\
&= - \frac{dN_i(d_j) - Y_i(d_j)h(d_j)}{\bar{Y}(d_j)}
\end{aligned}$$

The second variant holds when using the exponential form of  $S$ .

Finally, let  $A(a, b)$  be the integral under  $S$  from  $a$  to  $b$ .

$$A(0, \tau) = w_0 1 + \sum_{j=1}^m w_j S(d_j) \tag{2.8}$$

$$\text{nonumber} \tag{2.9}$$

$$\begin{aligned}
\frac{\partial A(0, \tau)}{\partial w_k} &= \frac{\partial A(d_1, \tau)}{\partial w_k} \\
\frac{\partial A(d_1, \tau)}{\partial w_k} &= \sum_{j=1}^m w_j U_{kj} \\
&= \sum_{j=1}^m w_j S(d_j) \sum_{i=1}^j u_{ki} \\
&= \sum_{i=1}^m u_{ki} \left[ \sum_{j=i}^m w_j S(d_j) \right] \\
&= \sum_{i=1}^m A(d_i, \tau) u_{ki}
\end{aligned} \tag{2.10}$$

This is a similar sum to before, with a new set of weights. Any given (time1, time2) observation involves  $u$  terms over that (time1,time2) range, we can form a single cumulative sum and use the value at time2 minus the value at time1. A single running total will not work for multiple  $\tau$  values, however, since the weights depend on  $\tau$ . Write  $A(d_i, \tau) = A(d_1, \tau) - A(d_1, d_i)$  and separate the two sums. For the first  $A(d_1, \tau)$  can be moved outside the sum, and so will play the same role as  $S(d)$  in  $U$ , the second becomes a weighted cumulative sum with weights that are independent of tau.

Before declaring victory with this last modification take a moment assess the computational impact of replacing  $A(d_i, \tau)$  with two terms. Assume  $m$  time points,  $d$  deaths, and  $n$  observations.

Method 1 will need to  $O(d)$  to compute the weights,  $O(d)$  for the weighted sum, and  $O(2n)$  to create each column of the output for (time1, time2) data, for a total of  $O(2m(n + d))$ . The second needs to create two running sums, each  $O(d)$  and  $O(4n)$  for each column of output for  $O(2d + 4mn)$ . The more 'clever' approach is always slightly slower! (The naive method of adding up rectangles is much slower than either, however, because it needs to compute  $U$  at every death time.)

**Variance** Efficient computation of a the variance of the cumulative hazard within survfit is more complex, since this is computed at every event time. We consider three cases below.

Case 1: The simplest case is unclustered data without delayed entry; not (time1,time2) data.

- Let  $W = \sum_i w_i^2$  be the sum of case weights for all those at risk. At the start everyone is in the risk set. Let  $v(t) = 0$  and  $z(t) = 0$  be running sums.
- At each event time  $t$ 
  1. Update  $z(t) = z(t-) - h(t)/\bar{Y}(t)$ , the running sum of the right hand term in equation (2.2).
  2. For all  $i =$  (event or censored at  $t$ ), fill in their element in  $C_i$ , add  $(w_i C_i)^2$  to  $v(t)$ , and subtract  $w_i^2$  from  $W$ .
  3. The variance at this time point is  $v(t) + Wz^2(t)$ ; all those still at risk have the same value of  $C_k$  at this time point.

Case 2: Unclustered with delayed entry. Let the time interval for each observation  $i$  be  $(s_i, t_i)$ . In this case the contribution at each event time, in step 3 above, for those currently at risk is

$$\sum_j w_j^2 [z(t) - z(s_j)]^2 = z^2(t) \sum_j w_j^2 + \sum_j w_j^2 z^2(s_j) - 2z(t) \sum_j (w_j z(s_j))$$

This requires two more running sums. When an observation leaves the risk set all three of them are updated. In both case 1 and 2 our algorithm is  $O(d + n)$ .

Case 3: Clustered variance. The variance at each time point is now

$$\text{varNA}(t) = \sum_g \left( \sum_{i \in g} w_i \frac{\partial \text{NA}(t)}{\partial w_i} \right)^2 \quad (2.11)$$

The observations within a cluster do not necessarily have the same case weight, so this does not collapse to one of the prior cases. It is also more difficult to identify when a group is no longer at risk and could be moved over to the  $v(t)$  sum. In the common case of (time1, time2) data the cluster is usually a single subject and the weight will stay constant, but we can not count on that. In a marginal structural model (Robbins) for instance the inverse probability of censoring weights (IPCW) will change over time.

The above, plus the fact that the number of groups may be far less than the number of observations, suggests a different approach. Keep the elements of the weighted grouped  $C$  vector, one row per group rather than one row per observation. This corresponds to the inner term of equation (2.11). Now update each element at each event time. This leads to an  $O(gd)$

algorithm. A slight improvement comes from realizing that the increment for group  $g$  at a given time involves the sum of  $Y_i(t)w_i$ , and only updating those groups with a non-zero weight.

For the Fleming-Harrington update, the key realization is that if there are  $d$  death at a given timepoint, the increment in the cumulative hazard  $\text{NA}(t)$  is a sum of  $d$  'normal' updates, but with perturbed weights for the tied deaths. Everyone else at risk gets a common update to the hazard. The basic computation and equations for  $C$  involve an small loop over  $d$  to create the increments, sort of like an aside in a stage play, but then adding them into  $C$  is the same as before. The basic steps for case 1-3 do not change.

The C code for `survfit` has adapted case 3 above, for all three of the hazard, survival, and AUC. A rationale is that if the data set is very large the user can choose `std.err = FALSE` as an option, then later use the residuals function to get values at selected times. When the number of deaths gets over 1000, which is where we begin to really notice the  $O(nd)$  slowdown, a plot only needs 100 points to look smooth. Most use cases will only need variance at a few points. Another reason is that for ordinary survival, robust variance is almost never requested unless there is clustering.

## 2.3 Aalen-Johansen

In a multistate model let the hazard for the  $jk$  transition be

$$\begin{aligned}\hat{\lambda}_{jk}(t) &= \frac{\sum_i dN_{ijk}(t)}{\sum_i Y_{ij}(t)} \\ &= \frac{\bar{N}_{jk}(t)}{\bar{Y}(t)}\end{aligned}$$

and gather terms together into the nstate by nstate matrix  $A(t)$  with  $A_{jk}(t) = \hat{\lambda}_{jk}(t)$ , with  $A_{jj}(t) = -\sum_{k \neq j} A_{jk}(t)$ . That is, the row sums of  $A$  are constrained to be 0.

The cumulative hazard estimate for the  $j : k$  transition is the cumulative sum of  $\hat{\lambda}_{jk}(t)$ . Each is treated separately, there is no change in methods or formula from the single state model.

The estimated probability in state  $p(t)$  is a vector with one element per state and the natural constraint that the elements sum to 1; it is a probability distribution over the states. Two estimates are

$$p(t) = p(0) \prod_{s \leq t} e^{A(s)} \tag{2.12}$$

$$\begin{aligned}p(t) &= p(0) \prod_{s \leq t} (I + A(s)) \\ &= p(0) \prod_{s \leq t} H(s)\end{aligned} \tag{2.13}$$

Here  $p(0)$  is the estimate at the starting time point. Both  $\exp(A(s))$  and  $I + A(s)$  are transition matrices, i.e., all elements are positive and rows sum to 1. They encode the state transition probabilities at time  $s$ , the diagonal is the probability of no transition for observations in the given state. At any time point with no observed transitions  $A(s) = 0$  and  $H(s) = I$ ; wlog the

product is only over the discrete times  $s$  at which an event (transition) actually occurred. For matrices  $\exp(C)\exp(D) \neq \exp(C+D)$  unless  $DC = CD$ , i.e., equation (2.12) does not simplify.

Equation (2.13) is known as the Aalen-Johansen estimate. In the case of a single state it reduces to the Kaplan-Meier, for competing risks it reduces to the cumulative incidence (CI) estimator. Interestingly, the exponential form (2.12) is almost never used for `survfit.formula`, but is always used for the curves after a Cox model. For `survfit.formula`, but of the forms obey the basic probability rules that  $0 \leq p_k(t) \leq 1$  and  $\sum_k p_k(t) = 1$ ; all probabilities are between 0 and 1 and the sum is 1. The analog of (2.13) has been proposed by some authors for Cox model curves, but it can lead to negative elements of  $p(t)$  in certain cases, so the survival package does not support it.

The computations are straightforward. The code makes two passes through the data, the first to create all the counts: number at risk, number of events of each type, number censored, etc. Since for most data sets the total number at risk decreases over time this pass is done in reverse time order since it leads to more additions than subtractions in the running number at risk and so less prone to round off error. (Unless there are extreme weights this is gilding the lily - there will not be round off error for either case.) The rule for tied times is that events happen first, censors second, and entries third; imagine the censors at  $t + \epsilon$  and entries at  $t + 2\epsilon$ .

When a particular subject has multiple observations, say times of (0,10), (10,15) and (15,24), we don't want the output to count this as a "censoring" and/or entry at time 10 or 15. The `position` vector is 1= first obs for a subject, 2= last obs, 3= both (someone with only one row of data), 0=neither. This subject would be 1,0,0,2. The position vector is created by the `survflag` routine. If the same subject id were used in two curves the counts are separate for each curve; for example curves by enrolling institution in a multi-center trial, where two centers happened to have an overlapping id value.

A variant of this is if a given subject changed curves at time 15, which occurs when users are estimating the "extended Kaplan-Meier" curves proposed by Snapinn [2], in which case the subject will be counted as censored at time 15 wrt the first curve, and an entry at time 15 for the second. (I myself consider this approach to be statistically unsound.)

If a subject changed case weights between the (0,10) and (10,15) interval we do not count them as separate entry and exits for the weighted `n.enter` and `n.censor` counts. This means that the running total of weighted `n.enter`, `n.event`, and `n.censor` will not recreate the weighted `n.risk` value; the latter *does* change when weights change in this way. The rationale for this is that the entry and censoring counts are mostly used for display, they do not participate in further computations.

### 2.3.1 Data

The `survfit` routine uses the `survcheck` routine, internally, to verify that the data set follows an overall rule that each subject in the data set follows a path which is physically possible:

1. Cannot be in two places at once (no overlaps)
2. Over the interval from first entry to last follow-up, they have to be somewhere (no gaps).
3. Any given state, if entered, must be occupied for a finite amount of time (no zero length intervals).

4. States must be consistent. For an interval  $(t1, t2, B)$  = entered state B at time t2, the current state for the next interval must be B (no teleporting).

I spent some time thinking about whether I should allow the `survfit` routine to bend the rules, and allow some of these. First off, number 1 and 3 above are not negotiable; otherwise it becomes impossible to clearly define the number at risk. But perhaps there are use cases for 2 and/or 4?

One data example I have used in the past for the Cox model is a subject on a research protocol, over 20 years ago now, who was completely lost to follow-up for nearly 2 years and then reappeared. Should they be counted as “at risk” in that interim? I argue no, on the grounds that they were not at risk for an “event that would have been captured” in our data set — we would never have learned of a disease progression. Someone should not be in the denominator of the Cox model (or KM) at a given time point if they cannot be in the numerator. We decided to put them into the data set with a gap in their follow-up time.

But in the end I’ve decided no to bending the rules. The largest reason is that in my own experience a data set that breaks any of 1–4 above is almost universally the outcome of a programming mistake. The above example is one of only 2 non-error cases I can think of, in nearly 40 years of clinical research: I *want* the routine to complain. (Remember, I’d like the package to be helpful to others, but I wrote the code for me.) Second is that a gap can easily be accommodated by creating extra state which plays the role of a waiting room. The  $P(\text{state})$  estimate for that new state is not interesting, but the others will be identical to what we would have obtained by allowing a gap. Instances of case 4 can often be solved by relabeling the states. I cannot think of an actual use case.

### 2.3.2 Counting the number at risk

For a model with  $k$  states, counting the number at risk is fairly straightforward, and results in a matrix with  $k$  columns and one row per time point. Each element contains the number who are currently in that state and not censored. Consider the simple data set shown in figure 2.2.

The obvious algorithm is simple: at any given time point draw a vertical line and count the number in each state who intersect it. At a given time  $t$  the rule is that events happen first, then censor, then entry. At time 2 subject 4 has an event and subject 3 enters: the number at risk for the four states at time point 2 is  $(4, 0, 0, 0)$ . At time point  $2 + \epsilon$  it is  $(4, 1, 0, 1)$ . At time 9 subject 3 is still at risk.

The default is to report a line of output at each unique censoring and event time, adding rows for the unique entry time is an optional argument. Subject 5 has 4 intervals of  $(1,3b)$ ,  $(3,6+)$ ,  $(6,8b)$  and  $(8,11+)$ : time 6 is not reported as a censoring time, nor as an entry time. Sometimes a data set will have been preprocessed by the user with a tool like `survSplit`, resulting in hundreds of rows for some subjects, most of which are ‘no event here’, and there is no reason to include all these intermediate times in the output. We make use of an ancillary vector `position` which is 1 if this interval is the leftmost of a subject sequence, 2 if rightmost, 3 if both. If someone had a gap in followup we would treat their re-entry to the risk sets as an entry, however; the `survcheck` routine will have already generated an error in that case.

The code does its summation starting at the largest time and moving left, adding and subtracting from the number at risk (by state) as it goes. In most data sets the number at risk decreases over time, going from right to left results in more additions than subtractions and

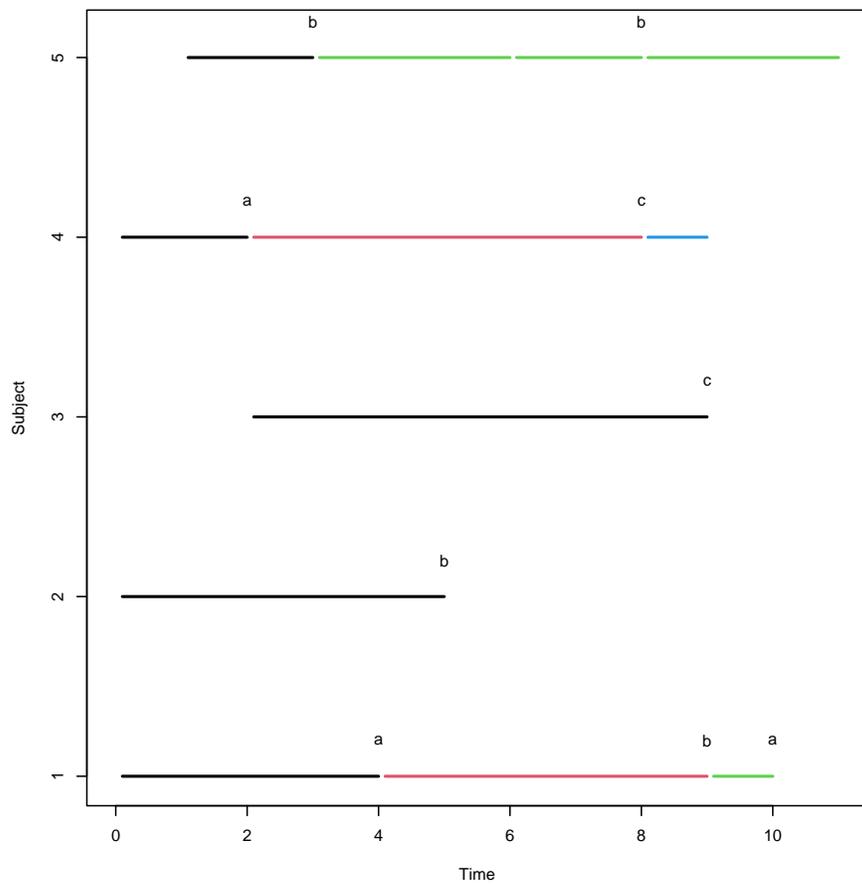


Figure 2.2: A simple multi-state data set. Colors show the current state of entry (black), a (red), b(green) or c (blue), letters show the occurrence of an event.

thus a smidge less potential roundoff error. Since it is possible for a subject's intervals to have different case weights, the number at risk *does* need to be updated at time 6 for subject 5, though that time point is not in the output.

In the example outcome *b* can be a repeated event, e.g., something like repeated infections. At time 8 the cumulative hazard for event *b* will change, but the probability in state will not. It would be quite unusual to have a data set with some repeated states and others not, but the code allows it.

Consider the simple illness-death model in figure ???. There are 3 states and 4 transitions in the model. The number at risk (`n.risk`), number censored (`n.censor`), probability in state (`pstate`) and number of entries (`n.enter`) components of the `survfit` call will have 3 columns, one per state. The number of events (`n.event`) and cumulative hazard (`cumhaz`) components will have 4 columns, one per transition, labeled as 1.2, 1.3, 2.1, and 2.3. In a matrix of counts with row = starting state and column = ending state, this is the order in which the non-zero elements appear. The order of the transitions is yet another consequence of the fact that R stores matrices in column major order.

There are two tasks for which the standard output is insufficient: drawing the curves and computing the area under the curve. For both these we need to know  $t_0$ , the starting point for the curves. There are 4 cases:

1. This was specified by the user using `start.time`
2. All subjects start at the same time
3. All subjects start in the same state
4. Staggered entry in diverse states

For case 1 use what the user specified, of course. For (time, status) data, i.e. competing risks use the same rule as for simple survival, which is  $\min(\text{time}, 0)$ . Curves are assumed to start at 0 unless there are negative time values. Otherwise for case 2 and 3 use  $\min(\text{time}_1)$ , the first time point found; in many cases this will be 0 as well.

The hardest case most often arises for data that is on age scale where there may be a smattering of subjects at early ages. Imagine that we had subjects at ages 19, 23, 28, another dozen from 30-40, and the first transition at age 41, 3 states A, B, C, and the user did not specify a starting time. Suppose we start at age=19 and that subject is in state B. Then  $p(19) = c(0,1,0)$ , and AJ  $p(t)$  estimate will remain (0,1,0) until there is a B:A or B:C transition, no matter what the overall prevalence of initial states as enrollment progresses. Such an outcome is not useful. The default for case 4 is to use the smallest event time at time 0, with initial prevalence the distribution of state for those observations which are at risk at that time. The prevalence at the starting time is saved in the fit object as `p0`. The best choice for case 4 is a user specified time, however.

Since we do not have an estimate of  $p$  = probability in state before time 0, the output will not contain any rows before that time point. Note that if there is an event exactly at the starting time: a death at time 0 for competing risks, a transition at `start.time`, or case 4 above, that the first row of the `pstate` matrix will not be the same as `p0`. The output never has repeats of the same time value (for any given curve). One side effect of this is that a plotted curve never begins with a vertical drop.

## 2.4 Influence

Let  $C(t)$  be the influence matrix for the cumulative hazard  $\Lambda(t)$  and  $U(t)$  that for the probability in state  $p(t)$ . In a multistate model with  $s$  states there are  $s^2$  potential transitions, but only a few are observed in any given data set. The code returns the estimated cumulative hazard  $\hat{\Lambda}$  as a matrix with one column for each  $j : k$  transition that is actually observed, the same will be true for  $C$ . We will slightly abuse the usual subscript notation; read  $C_{ijk}$  as the  $i$ th row and the  $j : k$  transition (column). Remember that a transition from state  $j$  to the same state can occur with multiple events of the same type. Let

$$\bar{Y}_j(t) = \sum_i Y_{ij}(t)$$

be the number at risk in state  $j$  at time  $t$ . Then

$$\begin{aligned} C_i(t) &= \frac{\partial \Lambda(t)}{\partial w_i} \\ &= C_i(t-) + \frac{\partial \lambda(t)}{\partial w_i} \\ \lambda_{jk}(t) &= \frac{\sum w_i dN_{ijk}}{\bar{Y}_j(t)} \end{aligned} \tag{2.14}$$

$$\frac{\partial \lambda_{jk}(t)}{\partial w_i} = \frac{dN_{ijk}(t) - Y_{ij}(t)\lambda_{jk}(t)}{\bar{Y}_j(t)} \tag{2.15}$$

$$V_c(t) = \sum_i [w_i C_i(t)]^2 \tag{2.16}$$

At any time  $t$ , formula (2.15) captures the fact that an observation has influence only on potential transitions from its current state at time  $t$ , and only over the (time1, time2) interval it spans. At any given event time, each non-zero  $h_{jk}$  at that time will add an increment to only those rows of  $C$  representing observations currently at risk and in state  $j$ . The IJ variance is the weighted sum of squares. (For replication weights, where  $w=2$  means two actual observations, the weight is outside the brackets. We will only deal with sampling weights.)

The (weighted) grouped estimates are

$$\begin{aligned} C_{gjk}(t) &= \sum_{i \in g} w_i C_{ijk}(t) \\ &= C_{gjk}(t-) + \sum_{i \in g} w_i \frac{dN_{ijk}(t) - Y_{ij}(t)\lambda_{jk}(t)}{\bar{Y}_j(t)} \end{aligned} \tag{2.17}$$

$$= C_{gjk}(t-) + \sum_{i \in g} w_i \frac{dN_{ijk}(t)}{n_j(t)} - \left( \sum_{i \in g} Y_{ij}(t)w_i \right) \frac{\lambda_{jk}(t)}{\bar{Y}_j(t)} \tag{2.18}$$

$$\tag{2.19}$$

A  $j : k$  event at time  $t$  adds only to the  $jk$  column of  $C$ . The last line above (2.18) spits out the  $dN$  portion, which will normally be 1 or at most a few events at this time point, from the

second, which is identical for subjects at risk in group  $g$ . This allows us to split out the sum of weights. Let  $w_{gj}(t)$  be the sum of weights by group and state, which is kept in a matrix of the same form as  $U$  and can be efficiently updated. Each row of the grouped  $C$  is computed with the same effort as a row of the ungrouped matrix.

Let  $C$  be the per observation influence matrix,  $D$  a diagonal matrix of per-observation weights, and  $B$  the  $n$  by  $g$  01/ design matrix that collapses observations by groups. For instance, `B = model.matrix( factor(grp) -1)`. The the weighted and collapsed influence is  $V = B'DC$ , and the infinitesimal jackknife variances are the column sums of the squared elements of  $V$ , `colSums(V*V)`. A feature of the IJ is that the column sums of  $DC$  and of  $B'DC$  are zero. In the special case where each subject is a group and the weight for a subject is constant over time, then  $B'D = WB'$  where  $W$  is a diagonal matrix of per-subject weights, i.e, one can collapse and then weight rather than weight and then collapse. The `survfitci.c` routine (now replaced by `survfitaj.c`) made this assumption, but unfortunately the code had no test cases with disparate weights within a group and so the error was not caught for several years. There is now a test case.

For  $C$  above, where each row is a simple sum over event times, formula (2.18) essentially does the scale + sum step at each event time and so only needs to keep one row per group. The parent routine will warn if any subject id is in two different groups. Doing otherwise makes no statistical sense, IMHO, so any instance if this is almost certainly an error in the data setup. However, someone may one day find a counterexample, hence a warning rather than an error.

For the probability in state  $p(t)$ , if the starting estimate  $p(0)$  is provided by the user, then  $U(0) = 0$ . If not,  $p(0)$  is defined as the distribution of states among those at risk at the first event time  $\tau$ .

$$p_j = \frac{\sum w_i Y_{ij}(\tau)}{\sum w_i Y_i(\tau)}$$

$$\frac{\partial p_j}{\partial w_k} = \frac{Y_{kj}(\tau) - Y_k(\tau)p_j}{\sum w_i Y_i(\tau)}$$

Assume 4 observations at risk at the starting point, with weights of (1, 4, 6, 9) in states (a, b, a, c), respectively. Then  $p(0) = (7/20, 4/20, 9/20)$  and the unweighted  $U$  matrix for those observations is

$$U(0) = \begin{pmatrix} (1 - 7/20)/20 & (0 - 4/20)/20 & (0 - 9/20)/20 \\ (0 - 7/20)/20 & (1 - 4/20)/20 & (0 - 9/20)/20 \\ (1 - 7/20)/20 & (0 - 4/20)/20 & (0 - 9/20)/20 \\ (0 - 7/20)/20 & (0 - 4/20)/20 & (1 - 9/20)/20 \end{pmatrix}$$

with rows of 0 for all observations not at risk at the starting time. Weighted column sums are  $wU = 0$ .

The AJ estimate of the probability in state vector  $p(t)$  is defined by the recursive formula  $p(t) = p(t-)H(t)$ . Remember that the derivative of a matrix product  $AB$  is  $d(A)B + Ad(B)$  where  $d(A)$  is the elementwise derivative of  $A$  and similarly for  $B$ . (Write out each element of

the matrix product.) Then  $i$ th row of  $U$  satisfies

$$\begin{aligned}
U_i(t) &= \frac{\partial p(t)}{\partial w_i} \\
&= \frac{\partial p(t-)}{\partial w_i} H(t) + p(t-) \frac{\partial H(t)}{\partial w_i} \\
&= U_i(t-)H(t) + p(t-) \frac{\partial H(t)}{\partial w_i}
\end{aligned} \tag{2.20}$$

The first term of 2.20 collapses to ordinary matrix multiplication, the second to a sparse multiplication. Consider the second term for any chosen observation  $i$ , which is in state  $j$ . This observation appears only in row  $j$  of  $H(t)$ , and thus  $dH$  is zero for all other rows of  $H(t)$  by definition. If observation  $i$  is not at risk at time  $t$  then  $dH$  is zero. The derivative vector thus collapses an elementwise multiplication of  $p(t-)$  and the appropriate row of  $dH$ , or 0 for those not at risk.

Let  $Q(t)$  be the matrix containing  $p(t-)Y_i(t)dH_{j(i)}(t)/dw_i$  as its  $i$ th row,  $j(i)$  the state for row  $j$ . Then

$$\begin{aligned}
U(t_2) &= U(t_1)H(t_2) + Q(t_2) \\
U(t_3) &= (U(t_1)H(t_2) + Q(t_2))H(t_3) + Q(t_3) \\
&\vdots = \vdots
\end{aligned}$$

Can we collapse this in the same way as  $C$ , retaining only one row of  $U$  per group containing the weighted sum? The equations above are more complex due to matrix multiplication. The answer is yes, because the weight+collapse operation can be viewed as multiplication by a design matrix  $B$  on the left, and matrix multiplication is associative and additive. That is  $B(U+Q) = BU+BQ$  and  $B(UH) = (BU)H$ . However, we cannot do the rearrangement that was used in the single endpoint case to turn this into a sum, equation (2.6), as matrix multiplication is not commutative.

For the grouped IJ, we have

$$\begin{aligned}
U_g(t) &= \sum_{i \in g} w_i U_i(t) \\
&= U_g(t-)H(t) + \sum_{i \in g} w_i Y_{ij}(t) \frac{\partial H(t)}{\partial w_i}
\end{aligned}$$

The above has the potential to be a very slow computation. If  $U$  has  $g$  rows and  $p$  columns, at each event time there is a matrix multiplication  $UH$ , creation of the  $H$  and  $H'$  matrices, addition of the correct row of  $H'$  to each row of  $U$ , and finally adding up the squared elements of  $U$  to get a variance. This gives  $O(d(gp^2 + 2p^2 + gp + gp))$ . Due to the matrix multiplication, every element of  $U$  is modified at every event time, so there is no escaping the final  $O(gp)$  sum of squares for each column. The primary gain is to make  $g$  small.

At each event time  $t$  the leverage for the AUC must be updated by  $\delta U(t-)$ , where  $\delta$  is the amount of time between this event and the last. At each reporting which does not coincide with an event time, there is a further update with  $\delta$  the time between the last event and the reporting time, before the sum of squares is computed. The AUC at the left endpoint of the curve is 0 and likewise the leverage. If a reporting time and event time coincide, do the AUC first.

At each event time

1. Update  $w_{gj}(\tau) = \sum_{i \in g} Y_{ij}(\tau)w_i$ , the number currently at risk in each group.
2. Loop across the events ( $dN_{ijk}(\tau) = 1$ ). Update  $C$  using equation (2.21) and create the transformation matrix  $H$ .
3. Compute  $U(\tau) = U(\tau-)H$  using sparse methods.
4. Loop a second time over the  $dN_{ijk}(\tau)$  terms, and for all those with  $j \neq k$  apply (2.22) and (2.23).
5. For each type of transition  $jk$  at this time, apply equation (2.24), and if  $j \neq k$ , (2.25) and (2.26).
6. Update the variance estimates, which are column sums of squared elements of  $C$ ,  $U$ , and AUC leverage.

$$\lambda_{jk}(\tau) = \frac{\sum_i dN_{ijk}(\tau)}{\sum_i w_i Y_{ij}(\tau)}$$

$$C_{g(i)jk}(\tau) = C_{gjk}(\tau-) + w_{g(i)} \frac{dN_{ijk}(\tau)}{\sum_i w_i Y_{ij}(\tau)} \quad (2.21)$$

$$U_{g(i)j}(\tau) = U_{gj}(\tau-) - w_{g(i)} \frac{dN_{ijk}(\tau)p_j(\tau-)}{\sum_i w_i Y_{ij}(\tau)} \quad (2.22)$$

$$U_{g(i)k}(\tau) = U_{gk}(\tau-) + w_{g(i)} \frac{dN_{ijk}(\tau)p_j(\tau-)}{\sum_i w_i Y_{ij}(\tau)} \quad (2.23)$$

$$C_{ghk}(\tau) = C_{gjk}(\tau-) - w_{gj}(\tau) \frac{\lambda_{jk}(\tau)}{\sum_i w_i Y_{ij}(\tau)} \quad (2.24)$$

$$U_{gj}(\tau) = U_{gj}(\tau-) + w_{gj}(\tau) \frac{\lambda_{jk}(\tau)p_j(\tau-)}{\sum_i w_i Y_{ij}(\tau)} \quad (2.25)$$

$$U_{gk}(\tau) = U_{gk}(\tau-) - w_{gj}(\tau) \frac{\lambda_{jk}(\tau)p_j(\tau-)}{\sum_i w_i Y_{ij}(\tau)} \quad (2.26)$$

In equations (2.21)–(2.23) above  $g(i)$  is the group to which observed event  $dN_i(\tau)$  belongs. At any reporting time there is often 1 or at most a handful of events. In equations (2.24)–(2.26) there is an update for each row of  $C$  and  $U$ , each row a group. Most often, the routine will be called with the default where each unique subject is their own group: if only 10% of the subjects are in group  $j$  at time  $\tau$  then 90% of the  $w_{gj}$  weights will be 0. It may be slightly faster to test for  $w$  positive before multiplying by 0?

For the ‘sparse multiplication’ above we employ simple approach: if the  $H - I$  matrix has only 1 non-zero diagonal, then  $U$  can be updated in place. For the NAFLD example in the survival vignette, for instance, this is true for all but 161 out of the 6108 unique event times (5%). For the remainder of the event times it suffices to create a temporary copy of  $U(t-)$ .

### 2.4.1 Residuals

The residuals.survfit function will return the IJ values at selected times, and is the more useful interface for a user. One reason is that since the user specifies the times, the result can be returned as an array with (observation, state, time) as the dimensions, even for a fit that has multiple curves. Full influence from the survfit routine, on the other hand, are returned as a list with one per curve. The reason is that each curve might have a different set of event times.

Because only a few times are reported, an important question is whether a more efficient algorithm can be created. For the cumulative hazard, each transition is a separate computation, so we can adapt the prior fast computation. We now have separate hazards for each observed transition  $jk$ , with numerator in the n.transitions matrix and denominator in the n.risk matrix. An observation in state  $j$  is at risk for all  $j : k$  transitions for the entire (time1, time2) interval; the code can use a single set of yindex, tindex, sindex and dindex intervals. The largest change is to iterate over the transitions, and the 0/1 for each  $j:k$  transition replaces simple status.

The code for this section is a bit opaque, here is some further explanation. Any given observation is in a single state over the (time1, time2) interval for that transition and accumulates the  $-\lambda_{jk}(t)/n_j(t)$  portion of the IJ for all  $j : k$  transitions that it overlaps, but only those that occur before the chosen reporting time  $\tau$ . Create a matrix hsum, one column for each observed  $jk$  transition, each column contains the cumulative sum of  $-\lambda_{jk}(t)/n_j(t)$  for that transition, and one row for each event time in the survfit object.

To index the hsum matrix create a matrix ymin which has a row per observation and a column per reporting time (so is the same size as the desired output), whose  $i,j$  element contains the index of the largest event time which is  $\leq \min(\text{time2}[i], \tau_j)$ , that is, the last row of hsum that would apply to this (observation time, reporting time) pair. Likewise let smin be a matrix which points to the largest event time which does *not* apply. The ymin and smin matrices apply to all the transitions; they are created using the outer and pmin functions. The desired  $d\lambda$  portion of the residual for the  $jk$  transition will be  $\text{hsum}[\text{smin}, \text{jk}] - \text{hsum}[\text{ymin}, \text{jk}]$ , where  $jk$  is a shorthand for the column of hmat that encodes the  $j : k$  transition. A similar sort of trickery is used for the  $dN_{jk}$  part of the residual. Matrix subscripts could be used to make it slightly faster, at the price of further impenetrability.

Efficient computation of residuals for the probability in state are more difficult. The key issue is that at every event time there is a matrix multiplication  $U(t-)H(t)$  which visits all  $n$  by  $p$  elements of  $U$ . For illustration assume 3 event times at 1,2,3, and the user only wants to report the leverage at time 3. Let  $B$  be the additions at each time point. Then

$$\begin{aligned}
 U(1) &= U(0)H(1) + B(1) \\
 U(2) &= U(1)H(2) + B(2) \\
 U(3) &= U(2)H(3) + B(3) \\
 U(4) &= U(3)H(4) + B(4) \\
 &= ([U(0)H(1) + B(1)]H(2) + B(2))H(3) + B(3) \\
 &= U(0)[H(1)H(2)H(3)] + B(1)[H(2)H(3)] + B(2)H(3) + B(3)
 \end{aligned}$$

The first four rows are the standard iteration. The  $U$  matrix will quickly become dense, since any  $j : k$  transition adds to any at-risk row in state  $j$ . Reporting back only at only a few time point does not change the computational burden of the matrix multiplications. We cannot reuse

the logarithm trick from the KM residuals, as  $\log(AB) \neq \log(A) + \log(B)$  when  $A$  and  $B$  are matrices.

The expansion in the last row above shows another way to arrange the computation. For any time  $t$  with only a  $j : k$  event,  $B(t)$  is non-zero only in columns  $j$  and  $k$ , and only for rows with  $Y_j(t) = 1$ , those at risk and in state  $j$  at time  $t$ . For example, assume that  $d_{50}$  were the largest death time less than or equal to the first reporting time. Accumulate the product in reverse order as  $B(50) + B(49)H(50) + B(48)[H(49)H(50)], \dots$ , updating the product of  $H$  terms at each step. Most updates to  $H$  involve a sparse matrix (only one event) so are  $O(s^2)$  where  $s$  is the number of states. Each  $B$  multiplication is also sparse. We can then step ahead to the next reporting time using the same algorithm, along with a final matrix update to carry forward  $U(50)$ . At the end there will have been  $d$ =number of event times matrix multiplications of a sparse update  $H(t)$  times the dense product of later  $H$  matrices, and each row  $i$  of  $U$  will have been 'touched' once for every death time  $k$  such that  $Y_i(d_k) = 1$  (i.e., each  $B$  term that it is part of), and once more for each prior reporting time.

I have not explored whether this idea will work in practice, nor do I see how to extend it to the AUC computation.

As a first example, consider the mgus2 data set, set, a case with 2 competing risks of death and plasma cell malignancy (PCM), and use a reporting times of 10, 20, and 30 years. There are  $n = \text{r nrow(m2)}$  observations,  $\text{r etot}$  events and  $\text{r sum(etime)}$  unique event times before 30 years, with an average of  $r = \text{r round(nisk,1)}$  subjects at risk at each of these event times. As part of data blinding follow up times in the MGUS data set were rounded to months, and as a consequence there are very few singleton event times.

Both the  $H$  matrix and products of  $H$  remain sparse for any row corresponding to an absorbing state (a single 1 on the diagonal), so for a competing risk problem the  $UH$  multiplication is  $O(3n)$  multiplications, those for the nested algorithm will be  $O(3r)$  where  $r$  is the average number at risk. In this case the improvement for the nested algorithm is modest, and at an additional price of  $9d$  multiplications to accumulate  $H$ .

```
> ndata <- tmerge(nafld1[,1:8], nafld1, id=id, death= event(futime, status))
> ndata <- tmerge(ndata, subset(nafld3, event=="nafld"), id,
  nafld= tdc(days))
> ndata <- tmerge(ndata, subset(nafld3, event=="diabetes"), id = id,
  diabetes = tdc(days), e1= cumevent(days))
> ndata <- tmerge(ndata, subset(nafld3, event=="htn"), id = id,
  htn = tdc(days), e2 = cumevent(days))
> ndata <- tmerge(ndata, subset(nafld3, event=="dyslipidemia"), id=id,
  lipid = tdc(days), e3= cumevent(days))
> ndata <- tmerge(ndata, subset(nafld3, event %in% c("diabetes", "htn",
  "dyslipidemia")),
  id=id, comorbid= cumevent(days))
> ndata$cstate <- with(ndata, factor(diabetes + htn + lipid, 0:3,
  c("0mc", "1mc", "2mc", "3mc")))
> temp <- with(ndata, ifelse(death, 4, comorbid))
> ndata$event <- factor(temp, 0:4,
  c("censored", "1mc", "2mc", "3mc", "death"))
> ndata$age1 <- ndata$age + ndata$tstart/365.25 # analysis on age scale
```

```

> ndata$age2 <- ndata$age + ndata$tstop/365.25
> ndata2 <- subset(ndata, age2 > 50 & age1 < 90)
> nfit <- survfit(Surv(age1, age2, event) ~1, ndata, id=id, start.time=50,
                 p0=c(1,0,0,0,0), istate=cstate, se.fit = FALSE)
> netime <- (nfit$time <=90 & rowSums(nfit$n.event) > 0)
> # the number at risk at any time is all those in the intial state of a transition
> # at that time
> from <- as.numeric(sub("\\.[0-9]*", "", colnames(nfit$n.transition)))
> fmat <- model.matrix(~ factor(from, 1:5) -1)
> temp <- (nfit$n.transition %*% fmat) >0 # TRUE for any transition 'from' state
> frisk <- (nfit$n.risk * ifelse(temp,1, 0))
> nrisk <- rowSums(frisk[netime,])
> maxrisk <- apply(frisk[netime,],2,max)

```

As a second case consider the NAFLD data used as an example in the survival vignette, a 5 state model with death and 0, 1, 2, or 3 metabolic comorbidities as the states. For a survival curve on age scale, using age 50 as a starting point and computing the influence at age 90, the data set has ‘`r nrow(ndata2)`’ rows that overlap the age 50-90 interval, while the average risk set size is  $r = \text{‘r round(mean(nrisk))’}$ , < 12% of the data rows. There are ‘`r sum(netime)`’ unique event times between 50 and 90 of which ‘`r sum(rowSums(nfit$n.transition[netime,]) ==1)`’ have only a single event type and the remainder 2.

In this case the first algorithm can take advantage of sparse  $H$  matrices and the second of sparse  $B$  matrices. The big gain is rows of  $B$  that are 0. To take maximal advantage of this we can keep an updated vector of the indices of the rows that are at risk, see the section on skiplists. Note: this has not yet been implemented.

For the AUC we can use a rearrangement. Below we write out the leverage on the AUC at time  $d_4$ , the fourth event, with  $w_i = d_{i+1} - d_i$  the width of the adjacent rectangles that make up the area under the curve.

$$\begin{aligned}
A(d_4) = & U(0) [w_0 + H(d_1)w_1 + H(d_1)H(d_2)w_2 + H(d_1)H(d_2)H(d_3)w_3] + \\
& B(d_1) [w_1 + H(d_2)w_2 + H(d_2)H(d_3)w_3] + \\
& B(d_2) [w_2 + H(d_3)w_3] + B(d_3)w_3
\end{aligned}$$

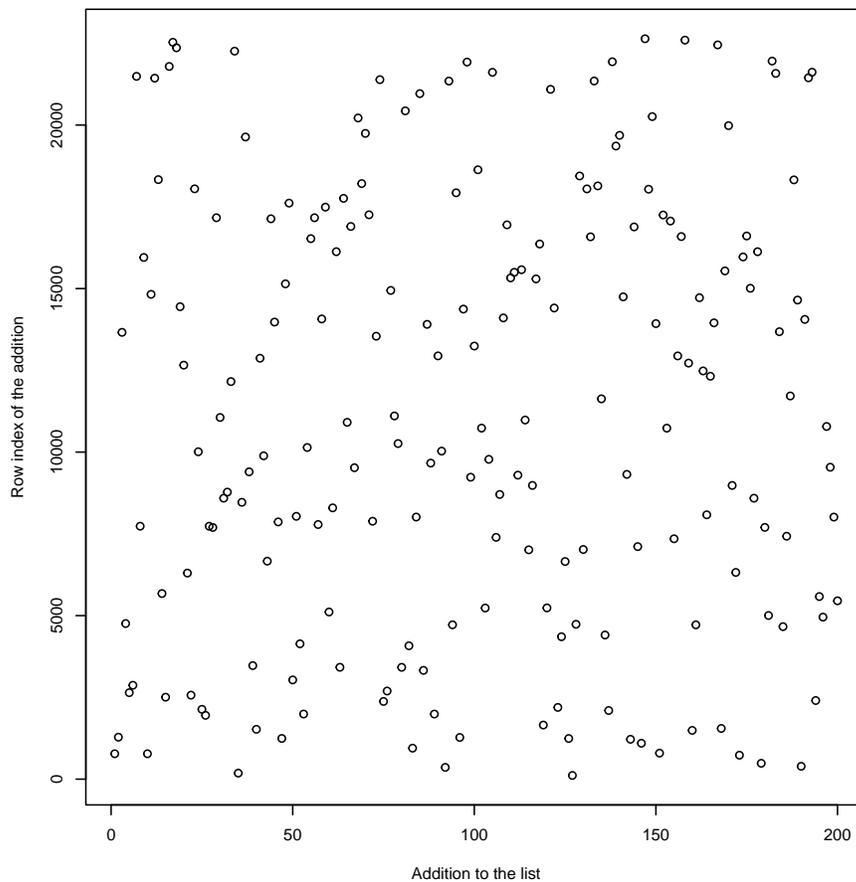
For  $U$  the matrix sequence was  $I, H_3I, H_2H_3I, \dots$ , it is now starts with  $w_3I$ , at the next step we multiply the prior matrix by  $H_3$  and add  $w_2I$ , at the third multiply the prior matrix by  $H_2$  and add  $w_1$ , etc.

## 2.5 Skiplists

An efficient structure for a continually updated index to the rows currently at risk appears to be a skiplist, see table II of Pugh [1]. Each observation is added (at time2) and deleted (at time1) just once from the risk set so we have  $n$  additions and  $n$  deletions, which are faster for skip lists than binary trees. Reading out the entire list, which is done at each event time, has the same cost for each structure. Searching is faster for trees, but we don’t need that operation. If there are  $k$  states we keep a separate skiplist containing those at risk for each state.

When used in the residuals function, we can modify the usual skiplist algorithm to take advantage of two things: we know the maximum size of the list from the `n.risk` component, and that the additions are nearly in random order. The last follows from the fact that most data sets are in subject id order, so that ordering by observation's ending time skips wildly across data rows. For the `naflid` data example above, the next figure shows the row number of the first 200 additions to the risk set, when going from largest time to smallest. For the example given above, the maximum number at risk in each of the 5 states is (`r paste(apply(nfit$n.risk, 2, max), collapse=', ')`).

```
> sort2 <- order(ndata$age2)
> plot(1:200, rev(sort2)[1:200], xlab="Addition to the list",
      ylab="Row index of the addition")
```



Using  $p = 4$  as recommended by Pugh [1] leads to  $\log_4(1420) = 5$  linked lists; 1420 is the maximum number of the  $n = \text{'r nrow(ndata)}$ ' rows at risk in any one state at any one time. The first is a forward linked list containing all those currently at risk, the second has  $1/4$  the elements, the third  $1/16$ ,  $1/64$ ,  $1/256$ . More precisely,  $3/4$  of the nodes have a single forward pointer,  $3/16$  have 2 forward pointers corresponding to levels 1 and 2,  $3/64$  have 3 forward pointer, etc. All told there are  $4/3$  as many pointer as a singly linked list, a modest memory cost over a linked

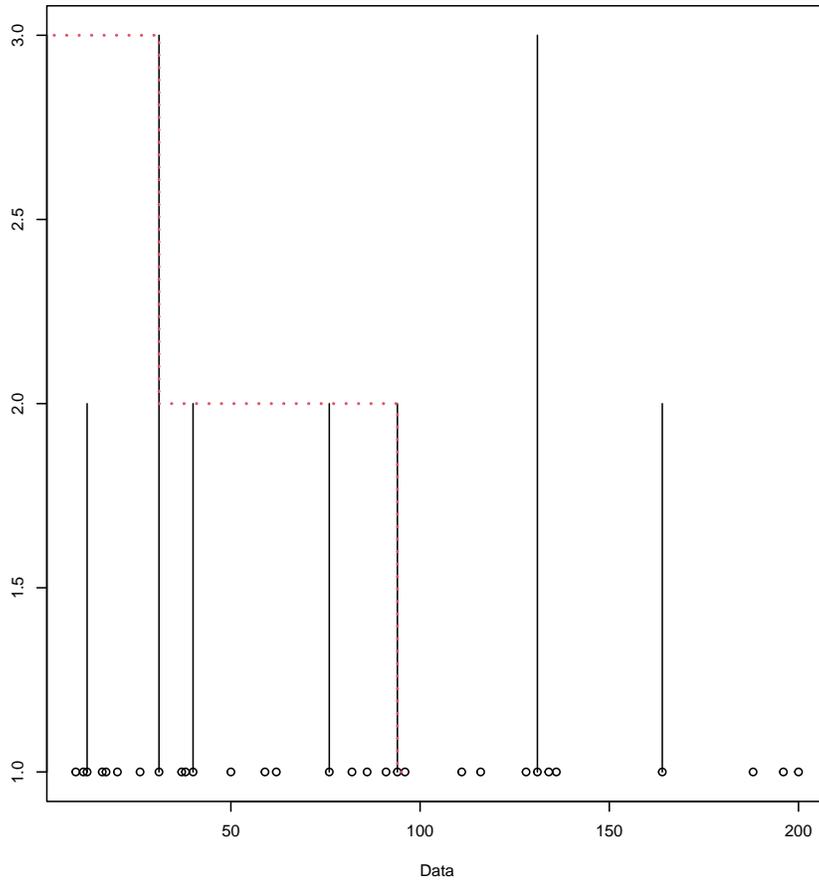


Figure 2.3: A simplified skiplist with 40 elements and 3 levels, showing the path used to insert a new element with value 100.

list. To add a new element to the list first find the point of insertion in list 4, list 3, list 2, and list 1 sequentially; each of which requires on average 2 forward steps. Randomly choose whether the new entry should participate in 1, 2, 3 or 4 of the levels, and insert it.

Figure 2.3 shows a simplified list with 40 elements and 3 levels, along with the search path for adding a new entry with value 100. The search starts at the head (shown here at  $x=0$ ), takes one step at level 3 to find the largest element  $< 100$  at level 3, then 3 steps at level 2, then one more at level 1 to find the insertion point. This is compared to 20 steps using only level 1. Since the height of a inserted node is chosen randomly there is always the possibility of a long final search on level 1, but even a string of 10 is unlikely:  $(3/4)^{10} < .06$ .

# Index

Aalen-Johansen  
  residuals, 22  
Andersen-Gill  
  influence, 18  
  
Kaplan-Meir, 4  
  
multi-state  
  influence, 18  
  
skiplist, 24  
surfit  
  number at risk, 15  
survfit, 3  
  AUC, 10  
survfitkm  
  robust variance, 6  
  
tied times, 3

# Bibliography

- [1] W. Pugh. Skip lists: a probabilistic alternative to balanced trees. *Communications ACM*, 33:668–676, 1990.
- [2] S. M. Snapinn, Q. Jiang, and B. Iglewicz. Illustrating the impact of a time-varying covariate with an extended Kaplan-Meier estimator. *Amstat News*, 59:301–7, 2005.