

Statistical Research and Applications Branch, NCI, Technical Report # 2003-03

Estimating Age Conditional Probability of Developing Cancer using a Piecewise Mid-Age Group Joinpoint Model for the Rates

Michael P. Fay

June 6, 2003

Abstract

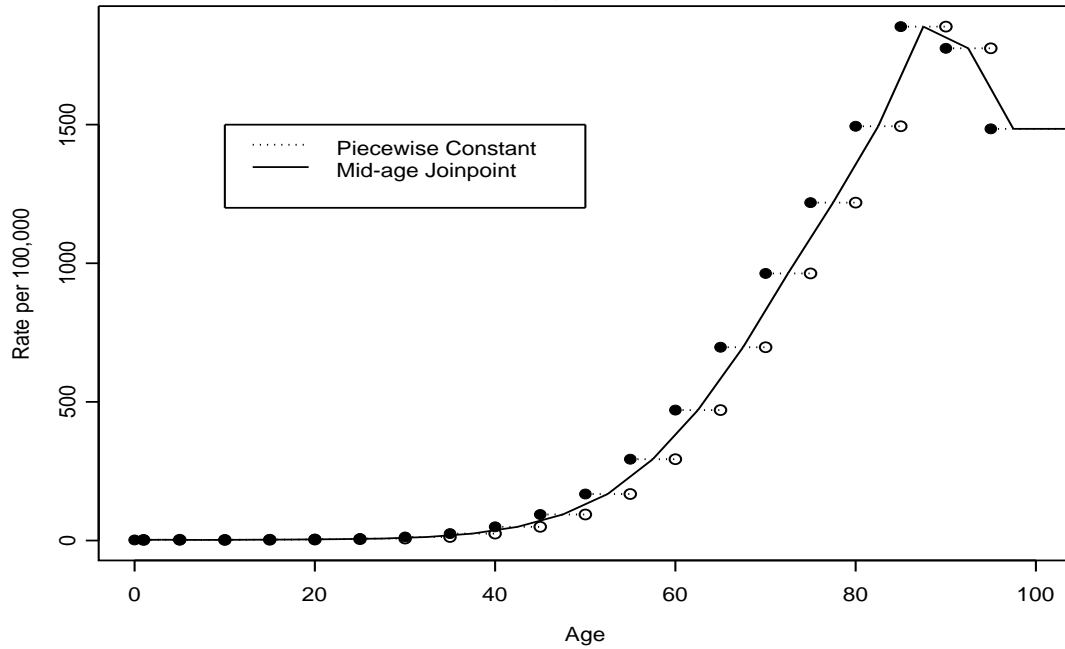
Fay, Pfeiffer, Cronin, Le, and Feuer (*Statistics in Medicine* 2003; **22**; 1837-1848) developed a formula to calculate the age-conditional probabilities of developing a disease (ACPDvD) from registry data. Fay et al. (2003) input into this formula a simple piecewise constant model for the rate functions, which have constant rates within each age group. In this paper, we detail a simple method for estimating rate function which does not have jumps at the beginning of age groupings. We call this method the mid-age group joinpoint (MAJ) model for the rates. The drawback of the MAJ model is that numerical integration must be used to estimate the resulting ACPDvD. To increase computational speed, we offer a piecewise approximation to the MAJ model, which we call the piecewise mid-age group joinpoint (PMAJ) model. This paper describes the PMAJ model for the rates input into the formula for ACPDvD described in Fay et al. (2003), which is the method used in version 5.0 of the freely available DevCan software (<http://srab.cancer.gov/DevCan/>).

1 Introduction

Fay, Pfeiffer, Cronin, Le, and Feuer (2003) showed how to calculate the age-conditional probabilities of developing a disease from registry data. Throughout this paper we use “cancer” as our disease of interest, but the method applies to specific types of cancer as well as other diseases where information is collected by population based surveillance methods. Fay et al. (2003) used a simple piecewise constant model for the rate functions, which have constant rates within each age group. Here we detail two more complicated models for the rates. The first model is a segmented regression model or joinpoint model for the rates, where the rate function is a series of linear functions that join at the mid-points of the age groups, and the rate function is constant before the first mid-point and after the last “mid-point” (because the last interval goes to infinity, the last “mid-point” is not really a mid-point at all, see below). We will call this model the MAJ (mid-age group joinpoint) model for the rates. In Figure 1 we show how both the piecewise constant model and the mid-age group joinpoint model

apply to all invasive cancer mortality in the U.S. in 1998-2000. Notice that the mid-age joinpoint model gives a more smoothly changing and probably a better modeled rate; we do not expect actual risks for people to jump on birthdays every 5 years like the piecewise constant model.

Figure 1: US All Invasive Cancer Mortality Rates, 1998-2000, All Races, Both Sexes



The problem with the mid-age group joinpoint model is that it requires numeric integration for its calculation. A faster method uses a series of piecewise constant values to approximate the mid-age group joinpoint model. We call this the PMAJ (piecewise mid-age group joinpoint) model. The PMAJ does not require numeric integration, so it is much faster than the MAJ model. Version 5.0 of the DevCan software (DevCan, 2003) uses the PMAJ method.

Here is an outline of this report. Section 2 gives the motivation for the MAJ estimator of age-conditional probability of developing cancer. The appendix shows how to calculate the integral needed for section 2. Section 3 describes the PMAJ model and how it is used to estimate the age-conditional probability of developing cancer. Section 4 compares the PMAJ method with the method of Wun, et al. (1998), since that method was the method used by previous versions of the DevCan software.

Table 1: Notation
Random Variables and Parameters

Random Variables and Parameters	
$T =$ age at death	$T^* =$ age at first cancer or death before cancer
$J =$ type of death ($J = d$)=death from cancer ($J = o$)=death from other causes	$J^* =$ type of event ($J^* = c$)=first cancer ($J^* = c$)=death before first cancer
$\lambda_c(t) =$ rate at t for first cancer given alive	$\lambda_c^*(t) =$ rate at t for first cancer given alive and cancer-free
$\lambda_o(t) =$ rate at t for death before cancer given alive	$\lambda_o^*(t) =$ rate at t for death before cancer given alive and cancer-free
$\lambda_d(t) =$ rate at t for death from cancer given alive	
$\lambda_a(t) =$ rate at t for death given alive	$\lambda_a^*(t) =$ rate at t for first cancer or death before first cancer given alive and cancer-free
$S_j(t) = \exp \left\{ - \int_0^t \lambda_j(u) du \right\}$ for $j = a, c, o, d$	$S_j^*(t) = \exp \left\{ - \int_0^t \lambda_j^*(u) du \right\}$ for $j = a, c, o$
Observations	
Within the age interval, $[a_i, a_{i+1})$, and within the calendar interval of interest we observe...	
$c_i =$ number of first cancer incident cases	$n_i^{(j)} =$ estimate of person-years alive associated with $j = c, d, o$
$d_i =$ number of cancer deaths	(DevCan uses the sum of mid-year populations during the calendar interval of interest)
$o_i =$ number of other deaths	

2 Mid-Age group Joinpoint Estimator

Fay, Pfeiffer, Cronin, Le, and Feuer (2003) assumed that the hazard rate for other cause (i.e., non-cancer) mortality is the same for people with and without cancer. Fay et al. (2003) gave a formula for the age-conditional probability of developing cancer between the ages of x and y given alive and cancer-free just before age x as

$$A(x, y) = \frac{\int_x^y \lambda_c(u) S(u-) du}{S_o(x-) \{1 - \int_0^x \lambda_c(u) S_a(u-) du\}}$$

See Table 1 for the notation taken from Fay, et al. (2003). The only change in notation from Fay, et al. (2003) is that we use the subscript a to represent all causes of events instead of a blank subscript. For example, we let $S^*(u) = S_a^*(u)$. Other notation in this paper is defined as it is introduced.

In Fay et al (2003), the rates were estimated by a piecewise constant model. Here we use a mid-age group joinpoint (MAJ) model, where we draw lines connecting the midpoints of the intervals except the first and last interval. The first interval is constant until the midpoint, and the last interval is constant after a nominal “midpoint”. This nominal “midpoint” is half the length of the previous age interval from the beginning of the last interval, and would be the midpoint if the last age interval was the same length as the previous interval. See Figure 1 for a plot of age by hazard with both methods.

We introduce new notation for breaking up the ages. Previously, we used $0 = a_0 < a_1 < \dots < a_k < a_{k+1} = \infty$. Now we use a joinpoint model with joins at

$$\frac{a_1}{2} < \frac{a_1 + a_2}{2} < \dots < \frac{a_{k-1} + a_k}{2} < a_k + \frac{a_k - a_{k-1}}{2}.$$

So we let

$$0 = t_{-1} < t_0 = \frac{a_1}{2} < t_1 = \frac{a_1 + a_2}{2} < \dots < t_{k-1} = \frac{a_{k-1} + a_k}{2} < t_k = a_k + \frac{a_k - a_{k-1}}{2} < t_{k+1} = \infty$$

(The numbering starts at -1 so that the indices for λ match the indices for t .) The MAJ estimator for the rate of event j for $j = c, d$, or o , at t_i for $i = 0, 1, \dots, k$, is

$$\tilde{\lambda}_{ji} = \tilde{\lambda}_j(t_i) = \frac{j_i}{n_i^{(j)}} \quad (1)$$

(Note that $\tilde{\lambda}_j(t_i) = \hat{\lambda}_j(a_i) = \hat{\lambda}_j(t_i)$, where $\hat{\lambda}_j(\cdot)$ is the piecewise constant function used by Fay et al. [2003]). We define $\tilde{\lambda}_{j,-1} = \tilde{\lambda}_{j0}$ and $\tilde{\lambda}_{j,k+1} = \tilde{\lambda}_{jk}$. For $j = a$, MAJ estimator for the rate at t_i is

$$\tilde{\lambda}_{ai} = \tilde{\lambda}_a(t_i) = \frac{o_i}{n_i^{(o)}} + \frac{d_i}{n_i^{(d)}}. \quad (2)$$

Then for $t \in [t_i, t_{i+1})$ for $i = 1, \dots, k$, we define $\tilde{\lambda}_j(t)$ as the point on the line defined by connecting the points $(t_i, \tilde{\lambda}_{ji})$ and $(t_{i+1}, \tilde{\lambda}_{j,i+1})$. In other words,

$$\tilde{\lambda}_j(t) = \alpha_{ji} + \beta_{ji}t,$$

where

$$\alpha_{ji} = \frac{t_{i+1}\tilde{\lambda}_{ji} - t_i\tilde{\lambda}_{j,i+1}}{t_{i+1} - t_i} \quad (3)$$

and

$$\beta_{ji} = \left(\frac{\tilde{\lambda}_{j,i+1} - \tilde{\lambda}_{ji}}{t_{i+1} - t_i} \right). \quad (4)$$

Thus, $\alpha_{j,-1} = \tilde{\lambda}_{j0}$ and $\beta_{j,-1} = 0$, and similarly by taking limits as $t_{k+1} \rightarrow \infty$ then $\alpha_{j,k} = \tilde{\lambda}_{j,k}$ and $\beta_{j,k} = 0$.

Now $\tilde{S}_j(u)$ for $u \in [t_i, t_{i+1})$ is

$$\begin{aligned}\tilde{S}_j(u) &= \exp\left(-\int_0^u \tilde{\lambda}_j(t) dt\right) \\ &= \exp\left(-\sum_{\ell=0}^i \int_{t_{\ell-1}}^{t_\ell} \{\alpha_{j,\ell-1} + \beta_{j,\ell-1}t\} dt - \int_{t_i}^u \{\alpha_{j,i} + \beta_{j,i}t\} dt\right)\end{aligned}$$

Note that (for $\ell = 0, 1, \dots, k$)

$$\begin{aligned}\int_{t_{\ell-1}}^{t_\ell} \{\alpha_{j,\ell-1} + \beta_{j,\ell-1}t\} dt &= (t_\ell - t_{\ell-1})\alpha_{j,\ell-1} + (t_\ell^2 - t_{\ell-1}^2)\frac{\beta_{j,\ell-1}}{2} \\ &= t_\ell\tilde{\lambda}_{j,\ell-1} - t_{\ell-1}\tilde{\lambda}_{j,\ell} + (t_\ell - t_{\ell-1})(t_\ell + t_{\ell-1})\frac{\beta_{j,\ell-1}}{2} \\ &= t_\ell\tilde{\lambda}_{j,\ell-1} - t_{\ell-1}\tilde{\lambda}_{j,\ell} + (t_\ell + t_{\ell-1})\left(\frac{\tilde{\lambda}_{j,\ell} - \tilde{\lambda}_{j,\ell-1}}{2}\right) \\ &= (t_\ell - t_{\ell-1})\left(\frac{\tilde{\lambda}_{j,\ell-1} + \tilde{\lambda}_{j,\ell}}{2}\right)\end{aligned}$$

so that for $i = 0, 1, \dots, k$,

$$\tilde{S}_j(t_i) = \exp\left(-\sum_{\ell=0}^i (t_\ell - t_{\ell-1})\left(\frac{\tilde{\lambda}_{j,\ell-1} + \tilde{\lambda}_{j,\ell}}{2}\right)\right)$$

Also notice that (when $u < \infty$)

$$\int_{t_i}^u \{\alpha_{j,i} + \beta_{j,i}t\} dt = (u - t_i)\alpha_{j,i} + (u^2 - t_i^2)\frac{\beta_{j,i}}{2}$$

Therefore when $u \in [t_i, t_{i+1})$,

$$\begin{aligned}\tilde{S}_j(u) &= \exp\left(-\sum_{\ell=0}^i (t_\ell - t_{\ell-1})\left(\frac{\tilde{\lambda}_{j,\ell-1} + \tilde{\lambda}_{j,\ell}}{2}\right) - (u - t_i)\alpha_{j,i} - (u^2 - t_i^2)\frac{\beta_{j,i}}{2}\right) \\ &= \tilde{S}_j(t_i) \exp\left(-\left[(u - t_i)\alpha_{j,i} + (u^2 - t_i^2)\frac{\beta_{j,i}}{2}\right]\right)\end{aligned}$$

Let $\tilde{A}(x, y)$ be the estimator of $A(x, y)$ using the MAJ model. The two integrals we need to estimate for $\tilde{A}(x, y)$ are of the type,

$$\tilde{F}_{j,h}(t) = \int_0^t \tilde{\lambda}_j(u) \tilde{S}_h(u-) du, \quad (5)$$

where in the numerator of $\tilde{A}(x, y)$ we need $\tilde{F}_{c,a}$ (i.e., $j = c$ and $h = a$ in equation 5), and in the denominator of $\tilde{A}(x, y)$ we need $\tilde{F}_{c,d}$. Suppose without loss of generality that $t \in [t_i, t_{i+1})$, then

$$\begin{aligned} \tilde{F}_{j,h}(t) &= \sum_{\ell=-1}^{i-1} \int_{t_\ell}^{t_{\ell+1}} \tilde{\lambda}_j(u) \tilde{S}_h(u-) du + \int_{t_i}^t \tilde{\lambda}_j(u) \tilde{S}_h(u-) du \\ &= \sum_{\ell=-1}^{i-1} \tilde{S}_h(t_\ell) \int_{t_\ell}^{t_{\ell+1}} (\alpha_{j\ell} + \beta_{j\ell}u) \exp\left(-\left[(u - t_\ell)\alpha_{h\ell} + (u^2 - t_\ell^2)\frac{\beta_{h\ell}}{2}\right]\right) du \\ &\quad + \tilde{S}_h(t_i) \int_{t_i}^t (\alpha_{ji} + \beta_{ji}u) \exp\left(-\left[(u - t_i)\alpha_{hi} + (u^2 - t_i^2)\frac{\beta_{hi}}{2}\right]\right) du \\ &= \sum_{\ell=-1}^{i-1} \tilde{S}_h(t_\ell) R_{j,h}(t_\ell, t_{\ell+1}) + \tilde{S}_h(t_i) R_{j,h}(t_i, t) \end{aligned}$$

where $R_{j,h}(t_\ell, v)$ (for $\ell = -1, 0, 1, 2, \dots, i$ and $v \leq t_{\ell+1}$) is defined implicitly (see the Appendix). Then,

$$\tilde{A}(x, y) = \frac{\tilde{F}_{c,a}(y) - \tilde{F}_{c,a}(x)}{\tilde{S}_o(x) \{1 - \tilde{F}_{c,d}(x)\}}.$$

3 Piecewise Mid-Age group Joinpoint Estimator

In the MAJ model we divided up the age line into $k + 2$ intervals. Here we define those intervals in both the t_i notation and the a_i notation.

$$\begin{aligned} I_0 &= [t_{-1}, t_0) = \left[0, \frac{a_1}{2}\right) \\ I_1 &= [t_0, t_1) = \left[\frac{a_1}{2}, \frac{a_1 + a_2}{2}\right) \\ &\quad \vdots \quad \quad \quad \vdots \\ I_i &= [t_{i-1}, t_i) = \left[\frac{a_{i-1} + a_i}{2}, \frac{a_i + a_{i+1}}{2}\right) \\ &\quad \vdots \quad \quad \quad \vdots \\ I_k &= [t_{k-1}, t_k) = \left[\frac{a_{k-1} + a_k}{2}, a_k + \frac{a_k - a_{k-1}}{2}\right) \\ I_{k+1} &= [t_k, \infty) = \left[a_k + \frac{a_k - a_{k-1}}{2}, \infty\right) \end{aligned}$$

In the MAJ model the rates for the first and the last intervals are represented by lines with zero slope, and the rates for the i th interval ($i = 1, \dots, k$) for the j th rate type ($j = a, c, d, o$) is a line defined by connecting the points $(t_{i-1}, \tilde{\lambda}_{j,i-1})$ and $(t_i, \tilde{\lambda}_{ji})$

(see equations 1 and 2 for definition of $\tilde{\lambda}_{ji}$). In the PMAJ model we divide the i th interval into m_i equal sized intervals, and use a piecewise constant estimate on each of those m_i intervals. One way to define m_i is to chose m_i so that each equal sized interval is 1/2 year long. In other words, $m_i = 2(t_i - t_{i-1})$. This is the definition of m_i that we use for the DevCan software version 5.0 (see DevCan, 2003), but all the following holds for arbitrary m_i .

Here are the details. Consider the h th (for $h = 1, \dots, m_i$) of the m_i intervals within interval i (for $i = 1, \dots, k$) for rate type j (for $j = a, c, d, o$). This interval is

$$\left[t_{i-1} + \frac{(h-1)(t_i - t_{i-1})}{m_i}, t_{i-1} + \frac{h \cdot (t_i - t_{i-1})}{m_i} \right)$$

For convenience we introduce new notation for the ends of this interval, let

$$t_{i-1,h} = t_{i-1} + \frac{h \cdot (t_i - t_{i-1})}{m_i}$$

so that $t_{i-1,0} = t_{i-1}$ and $t_{i-1,m_i} = t_i$. At the beginning of this interval the value of the rate is

$$\begin{aligned} \tilde{\lambda}_j(t_{i-1,h-1}) &= \alpha_{j,i-1} + \beta_{j,i-1} \left(t_{i-1} + \frac{(h-1)(t_i - t_{i-1})}{m_i} \right) \\ &= \frac{t_i \tilde{\lambda}_{j,i-1} - t_{i-1} \tilde{\lambda}_{ji}}{t_i - t_{i-1}} + \frac{(\tilde{\lambda}_{ji} - \tilde{\lambda}_{j,i-1})t_{i-1}}{t_i - t_{i-1}} + \frac{(h-1)(\tilde{\lambda}_{ji} - \tilde{\lambda}_{j,i-1})}{m_i} \\ &= \tilde{\lambda}_{j,i-1} + \frac{(h-1)(\tilde{\lambda}_{ji} - \tilde{\lambda}_{j,i-1})}{m_i} \end{aligned}$$

(see equations 3 and 4 for definitions of $\alpha_{j,i-1}$ and $\beta_{j,i-1}$). Similarly at the end of this interval the rate is

$$\tilde{\lambda}_j(t_{i-1,h}) = \tilde{\lambda}_{j,i-1} + \frac{h(\tilde{\lambda}_{ji} - \tilde{\lambda}_{j,i-1})}{m_i}$$

For the PMAJ model we simply assume a constant rate equal to the average of the beginning and the end values of the rate over this interval. In other words, under the PMAJ model for any $t \in [t_{i-1,h-1}, t_{i-1,h})$ we estimate the rate with

$$\dot{\lambda}_j(t) = \tilde{\lambda}_{j,i-1} + \frac{(2h-1)(\tilde{\lambda}_{ji} - \tilde{\lambda}_{j,i-1})}{2m_i}$$

Since the PMAJ model is a piecewise model, we can use Appendix A of Fay *et al.* (2003) to express the estimator of age conditional probability of developing cancer. The

only hard part is correctly defining the starting and ending of each piecewise interval. The ends of these intervals are

$$0 \equiv t_{-1} < t_0 < t_{0,1} < t_{0,2} < \cdots < t_{0,m_1-1} < t_1 < t_{1,1} < \cdots < t_{k-1,m_k-1} < t_k < t_{k+1} \equiv \infty$$

For convenience write these interval ends with only a single index as

$$0 \equiv \tau_0 < \tau_1 < \tau_2 < \tau_3 < \cdots < \tau_{m_1+1} < \tau_{m_1+2} < \tau_{m_1+3} < \cdots < \tau_{M+1} < \tau_{M+2} < \tau_{M+3} \equiv \infty$$

where $M = \sum_{i=1}^k m_i$. In other words, for $i = 1, \dots, k$, then $t_i = \tau_g$ and $t_{i,h} = \tau_{g+h}$, where $g = \sum_{\ell=1}^i m_\ell + 2$.

Now we can follow very similar notation to Appendix A of Fay *et al.* (2003). We now repeat that Appendix with the modifications to notation required for the PMAJ model. Let the estimator of $A(x, y)$ under the PMAJ model be denoted $\dot{A}(x, y)$. Let $\tau_i \leq x < \tau_{i+1}$ and $\tau_j < y \leq \tau_{j+1}$ for $x < y, i \leq j$, and $j \leq M + 2$. For convenience we regroup the ages after inserting group delimiters at x and y . Let the new delimiters be $0 = b_0 \leq b_1 \leq b_2 \leq \cdots \leq b_{M+5} = \infty$ where $b_0 = \tau_0, \dots, b_i = \tau_i, b_{i+1} = x, b_{i+2} = \tau_{i+1}, \dots, b_{j+1} = \tau_j, b_{j+2} = y, b_{j+3} = \tau_{j+1}, \dots, b_{M+5} = \tau_{M+3} = \infty$. We let

$$\dot{S}_a(b_\ell) = \exp \left\{ - \int_0^{b_\ell} \dot{\lambda}_a(u) du \right\} = \exp \left\{ - \sum_{u=0}^{\ell-1} \dot{\lambda}_a(b_u) (b_{u+1} - b_u) \right\},$$

and similarly $\dot{S}_d(b_\ell) = \exp \left\{ - \int_0^{b_\ell} \dot{\lambda}_d(u) du \right\}$ and $\dot{S}_o(b_\ell) = \exp \left\{ - \int_0^{b_\ell} \dot{\lambda}_o(u) du \right\}$. In this notation, the probability of developing cancer by age y given survival until age x is $A(x, y) = A(b_{i+1}, b_{j+2})$, and under the PMAJ model we estimate it with

$$\begin{aligned} \dot{A}(b_{i+1}, b_{j+2}) &= \frac{\sum_{\ell=i+1}^{j+1} \int_{b_\ell}^{b_{\ell+1}} \dot{\lambda}_c(b_\ell) \dot{S}_a(b_\ell) \exp \left(- \int_{b_\ell}^u \dot{\lambda}_a(b_\ell) dt \right) du}{\dot{S}_o(b_{i+1}) \left\{ 1 - \sum_{\ell=0}^i \int_{b_\ell}^{b_{\ell+1}} \dot{\lambda}_c(b_\ell) \dot{S}_d(b_\ell) \exp \left(- \int_{b_\ell}^u \dot{\lambda}_d(b_\ell) dt \right) du \right\}} \\ &= \frac{\sum_{\ell=i+1}^{j+1} \dot{\lambda}_c(b_\ell) \dot{S}_a(b_\ell) \int_{b_\ell}^{b_{\ell+1}} \exp \left(-(u - b_\ell) \dot{\lambda}_a(b_\ell) \right) du}{\dot{S}_o(b_{i+1}) \left\{ 1 - \sum_{\ell=0}^i \dot{\lambda}_c(b_\ell) \dot{S}_d(b_\ell) \int_{b_\ell}^{b_{\ell+1}} \exp \left(-(u - b_\ell) \dot{\lambda}_d(b_\ell) \right) du \right\}}. \end{aligned}$$

Because $\dot{\lambda}_a(b_\ell)$ or $\dot{\lambda}_d(b_\ell)$ may equal zero and $b_{\ell+1}$ may equal infinity, we let $\phi(\lambda, \ell) = \int_{b_\ell}^{b_{\ell+1}} \exp \left(-(u - b_\ell) \lambda \right) du$. These integrals are

$$\phi(\lambda, \ell) = \begin{cases} \frac{1 - \exp[-(b_{\ell+1} - b_\ell)\lambda]}{\lambda} & \text{if } \lambda > 0 \text{ and } b_{\ell+1} \neq \infty \\ b_{\ell+1} - b_\ell & \text{if } \lambda = 0 \text{ and } b_{\ell+1} \neq \infty \\ \frac{1}{\lambda} & \text{if } \lambda > 0 \text{ and } b_{\ell+1} = \infty \\ \infty & \text{if } \lambda = 0 \text{ and } b_{\ell+1} = \infty \end{cases}$$

where the case $\lambda = 0$ and $b_{\ell+1} = \infty$ is one of the “impossible” hypothetical cohorts (see Section 3.1 of Fay *et al.* 2003). Thus, we obtain,

$$\dot{A}(b_{i+1}, b_{j+2}) = \frac{\sum_{\ell=i+1}^{j+1} \dot{\lambda}_c(b_\ell) \dot{S}_a(b_\ell) \phi(\dot{\lambda}_a(b_\ell), \ell)}{\dot{S}_o(b_{i+1}) \left\{ 1 - \sum_{\ell=0}^i \dot{\lambda}_c(b_\ell) \dot{S}_d(b_\ell) \phi(\dot{\lambda}_d(b_\ell), \ell) \right\}}.$$

4 Comparing the Method of Wun, Merrill, and Feuer (1998) to the PMAJ Method

Since versions of the DevCan software prior to 5.0 used the method described in Wun, Merrill, and Feuer (1998), here we compare that method to the PMAJ method. The bulk of the comparison has previously been done (see Fay *et al.* 2002). That comparison assumed the simple piecewise hazards models using the method described in Fay *et al.* (2003). The only difference between the method described in Fay *et al.* (2003) and that described in this paper is that in this paper we estimate the hazard functions with the PMAJ method.

In Table 2 (see pages 12-15) we recalculate Table I-15 from Ries *et al.* (2003) which gives lifetime risks of developing certain cancers for different race and sex combinations. We give the old method of Wun, Merrill, and Feuer (1998), the new method presented in this paper, and the percent differences. For the the Wun, Merrill, and Feuer (1998) method the age groups of the data must be in 5 year intervals except the last open ended interval. For the new method the data can be input with any age intervals, and for the example in Table 2 the first age interval is 1 year, the second is 4 years, and all subsequent intervals except the last are 5 years. Thus, the input data are slightly different for the two methods.

In general the two methods agree to within about 2 percent (see Table 2). The only cancer type with larger than about 2 percent in absolute difference is acute lymphocytic leukemia (ALL). For ALL the absolute percent differences are as large as 4.5 percent (for black males). One reason for that large absolute percent difference is the small absolute size of the ALL lifetime risk, so small absolute changes in risk translate to large absolute percentage changes. Another reason may be that ALL is a pediatric cancer, so the differences in the input data may be part of the cause of the differences.

For age conditional probabilities of developing cancers, the methods give similar answers. Although for very small probabilities the absolute percent difference can be very large, in those cases the absolute difference is small. For large probabilities where the absolute difference between the methods may be larger, the absolute percent difference is small.

References

- DEVCAN: Probability of DEveloping CANcer software* Version 5.0, National Cancer Institute and Information Management Services, Inc., 2003. (Accessed at <http://srab.cancer.gov/DevCan/> on May 30, 2003).
- Fay, M.P., Pfeiffer, R., Cronin, K.A., Le, C. and Feuer, E.J. Comparison of Two Methods for Calculating Age-Conditional Probabilities of Developing Cancer. Technical Report #2002-01, Statistical Research and Applications Branch, National Cancer Institute 2002. (Accessed at <http://srab.cancer.gov/reports> on September 3, 2002).
- Fay, M.P., Pfeiffer, R., Cronin, K.A., Le, C., Feuer, E.J. (2003). Age-Conditional Probabilities of Developing Cancer. *Statistics in Medicine* **22**(11) 1837-1848.
- Lange, K. (1999). *Numerical Analysis for Statisticians* Springer:New York.
- Ries LAG, Eisner MP, Kosary CL, Hankey BF, Miller BA, Clegg L, Mariotto A, Fay MP, Feuer EJ, Edwards BK (eds). SEER Cancer Statistics Review, 1975-2000, National Cancer Institute. Bethesda, MD, http://seer.cancer.gov/csr/1975_2000, 2003.
- Wun, L-M, Merrill, R.M., and Feuer, E.J. Estimating lifetime and age-conditional probabilities of developing cancer. *Lifetime Data Analysis* 1998; **4**, 169-186.

Appendix: Calculation of R function

Recall that $R_{j,h}(t_\ell, v)$ represents an integral with 4 parameters. We can write it as

$$R(t_\ell, v, \alpha_{j\ell}, \beta_{j\ell}, \alpha_{h\ell}, \beta_{h\ell}) = \int_{t_\ell}^v (\alpha_{j\ell} + \beta_{j\ell}x) \exp\left(-\left[(x - t_\ell)\alpha_{h\ell} + (x^2 - t_\ell^2)\frac{\beta_{h\ell}}{2}\right]\right) dx$$

To simplify notation substitute let $t_\ell = u$ and $\alpha_{j\ell} = a_j, \beta_{j\ell} = b_j, \alpha_{h\ell} = a_h$, and $\beta_{h\ell} = b_h$. Thus,

$$R(u, v, a_j, b_j, a_h, b_h) = \int_u^v (a_j + b_jx) \exp\left(-\left[(x - u)a_h + (x^2 - u^2)\frac{b_h}{2}\right]\right) dx$$

Case 1: $b_j = 0$ and $b_h = 0$

For our application, whenever $v \rightarrow \infty$ then $b_j = 0$ and $b_h = 0$, so this is an important special case.

When $b_j = 0$ and $b_h = 0$ and $a_h = 0$ and we obtain

$$R(u, v, a_j, 0, a_h, 0) = \int_u^v a_j dx = (v - u)a_j$$

which goes to ∞ when $v \rightarrow \infty$.

When $b_j = 0$ and $b_h = 0$ and $a_h \neq 0$ and we obtain

$$\begin{aligned} R(u, v, a_j, 0, a_h, 0) &= \int_u^v a_j \exp(-[(x - u)a_h]) dx \\ &= \frac{a_j}{a_h} [1 - \exp(-[(v - u)a_h])] \end{aligned}$$

which goes to a_j/a_h when $v \rightarrow \infty$.

Case 2: General Case with $v < \infty$

To calculate the integral, $R(u, v, a_j, b_j, a_h, b_h)$ for finite v , we can use an adaptive use of Romberg's algorithm for numeric integration (we follow closely Lange, 1999, pp. 210-211).

Let

$$f(x) = f(x, u, a_j, b_j, a_h, b_h) = (a_j + b_j x) \exp\left(-\left[(x - u)a_h + (x^2 - u^2)\frac{b_h}{2}\right]\right)$$

Divide the interval $[u, v]$ into n equal subintervals of length $(v - u)/n$, and let

$$T_n = \frac{(v - u)}{n} \left[\frac{1}{2}f(u) + \frac{1}{2}f(v) + \sum_{i=1}^{n-1} f\left(u + \frac{i(v - u)}{n}\right) \right]$$

Then $\lim_{n \rightarrow \infty} T_n = R(u, v, a_j, b_j, a_h, b_h)$.

A more accurate approximation uses Romberg's algorithm,

$$R(u, v, a_j, b_j, a_h, b_h) \approx \frac{4T_{2n} - T_n}{3}$$

Let \hat{R} be our estimate of R . The algorithm we use to calculate \hat{R} is as follows:

1. Choose n .
2. Calculate T_n .
3. Calculate T_{2n} .
4. For $i=1$ to I_{max} do:
 - If $|T_{2^i n} - T_{2^{i-1} n}| < \delta$ then let $\hat{R} = \frac{4T_{2^i n} - T_{2^{i-1} n}}{3}$ and stop.
 - Otherwise calculate $T_{2^{i+1} n}$, and continue.

For example, one could use $n = 100$ and $\delta = 10^{-5}$ and $I_{max} = 100$.

Table 2: Lifetime Risk (percent) of Being Diagnosed with Cancer by Site, Race and Sex. 12 SEER Areas, 1998-2000. (Compare to Ries, et al. 2003, Table I-15). Each cell has 3 values: PMAJ method, Wun, Merrill, and Feuer (1998) method (WMF), and percent difference= $100(PMAJ - WMF)/WMF$.

Site	All Races		Whites		Blacks	
	Males	Females	Males	Females	Males	Females
All Sites	45.19	38.67	45.40	39.99	42.45	32.09
	44.88	38.65	45.05	39.90	42.47	32.38
	0.70	0.08	0.78	0.22	-0.03	-0.91
Invasive and In Situ	46.40	41.98	46.70	43.44	42.83	34.32
	46.03	41.87	46.30	43.26	42.83	34.59
	0.80	0.26	0.88	0.41	0.00	-0.77
Oral Cavity and Pharynx	1.40	0.67	1.40	0.69	1.38	0.50
	1.41	0.68	1.42	0.69	1.41	0.50
	-1.05	-0.93	-0.94	-0.87	-1.93	-1.30
Esophagus	0.75	0.26	0.77	0.25	0.77	0.38
	0.76	0.26	0.77	0.25	0.79	0.39
	-1.21	-1.05	-1.09	-0.93	-2.05	-1.81
Stomach	1.25	0.79	1.09	0.66	1.32	1.03
	1.27	0.80	1.10	0.66	1.35	1.05
	-1.09	-1.00	-0.99	-0.90	-1.91	-1.63
Colon and Rectum	5.97	5.66	6.03	5.66	4.94	5.38
	6.01	5.71	6.07	5.70	5.02	5.46
	-0.79	-0.74	-0.67	-0.63	-1.65	-1.48
Invasive and In Situ	6.31	5.94	6.36	5.92	5.29	5.71
	6.36	5.98	6.40	5.96	5.38	5.79
	-0.76	-0.72	-0.64	-0.61	-1.60	-1.46
Liver and Intrahepatic Bile Duct	0.86	0.42	0.72	0.35	0.80	0.36
	0.87	0.42	0.73	0.36	0.82	0.37
	-1.10	-0.72	-0.97	-0.48	-2.04	-1.96
Pancreas	1.23	1.24	1.23	1.22	1.27	1.34
	1.25	1.25	1.24	1.23	1.30	1.36
	-1.19	-1.08	-1.10	-0.97	-1.97	-1.79
Larynx	0.65	0.16	0.65	0.17	0.86	0.23
	0.65	0.17	0.66	0.17	0.88	0.24
	-1.18	-1.04	-1.06	-0.94	-2.00	-1.81
Invasive and In Situ	0.70	0.18	0.71	0.18	0.89	0.24
	0.71	0.18	0.71	0.18	0.91	0.25
	-1.18	-1.04	-1.07	-0.93	-2.00	-1.81

Note: Invasive cancer only unless specified otherwise

Table 2: (continued) Lifetime Risk (percent) of Being Diagnosed with Cancer by Site, Race and Sex. 12 SEER Areas, 1998-2000. (Compare to Ries, et al. 2003, Table I-15). Each cell has 3 values: PMAJ method, Wun, Merrill, and Feuer (1998) method (WMF), and percent difference= $100(PMAJ - WMF)/WMF$.

Site	All Races		Whites		Blacks	
	Males	Females	Males	Females	Males	Females
Lung and Bronchus	7.75	5.79	7.75	6.08	8.29	5.37
	7.84	5.85	7.83	6.13	8.45	5.47
	-1.11	-1.01	-1.01	-0.90	-1.88	-1.75
Melanomas of Skin	1.83	1.23	2.16	1.48	0.10	0.08
	1.84	1.24	2.18	1.49	0.11	0.08
	-1.02	-0.88	-0.89	-0.77	-1.80	-0.21
Invasive and In Situ	2.88	1.96	3.38	2.32	0.13	0.12
	2.91	1.98	3.41	2.34	0.13	0.12
	-0.96	-0.87	-0.81	-0.75	-1.84	-0.65
Breast	0.11	13.51	0.12	14.28	0.14	10.14
	0.11	13.56	0.12	14.31	0.14	10.27
	-1.12	-0.36	-1.03	-0.20	-1.86	-1.31
Invasive and In Situ	0.13	16.06	0.13	16.92	0.15	12.14
	0.13	16.09	0.13	16.92	0.15	12.28
	-1.12	-0.20	-1.04	-0.04	-1.87	-1.17
Cervix	-	0.79	-	0.75	-	0.97
	-	0.80	-	0.76	-	0.99
	-	-0.88	-	-0.74	-	-1.63
Corpus and Uterus, NOS	-	2.62	-	2.82	-	1.67
	-	2.65	-	2.85	-	1.70
	-	-0.89	-	-0.76	-	-1.73
Invasive and In Situ	-	2.66	-	2.87	-	1.69
	-	2.69	-	2.90	-	1.72
	-	-0.89	-	-0.76	-	-1.73
Ovary	-	1.72	-	1.85	-	1.10
	-	1.73	-	1.86	-	1.12
	-	-0.98	-	-0.86	-	-1.75
Prostate	17.28	-	16.90	-	20.40	-
	17.22	-	16.83	-	20.39	-
	0.35	-	0.42	-	0.06	-
Testis	0.35	-	0.42	-	0.10	-
	0.36	-	0.42	-	0.10	-
	-0.53	-	-0.52	-	-0.72	-

Note: Invasive cancer only unless specified otherwise

Table 2: (continued) Lifetime Risk (percent) of Being Diagnosed with Cancer by Site, Race and Sex. 12 SEER Areas, 1998-2000. (Compare to Ries, et al. 2003, Table I-15). Each cell has 3 values: PMAJ method, Wun, Merrill, and Feuer (1998) method (WMF), and percent difference= $100(PMAJ - WMF)/WMF$.

Site	All Races		Whites		Blacks	
	Males	Females	Males	Females	Males	Females
Urinary Bladder (In Situ and Inv)	3.52	1.13	3.91	1.22	1.42	0.78
	3.55	1.14	3.94	1.23	1.45	0.79
	-0.86	-1.01	-0.72	-0.89	-1.84	-1.69
Kidney and Renal Pelvis	1.46	0.87	1.53	0.91	1.23	0.86
	1.48	0.88	1.55	0.92	1.26	0.87
	-1.06	-0.89	-0.93	-0.87	-2.01	-1.00
Brain and Other Nervous System	0.66	0.53	0.75	0.59	0.32	0.30
	0.67	0.53	0.75	0.59	0.33	0.31
	-1.05	-0.63	-0.89	-0.62	-1.97	-1.62
Thyroid	0.30	0.84	0.32	0.87	0.14	0.46
	0.30	0.85	0.32	0.88	0.15	0.47
	-1.08	-0.86	-0.95	-0.73	-2.00	-1.67
Hodgkin's Disease	0.23	0.20	0.26	0.22	0.19	0.15
	0.24	0.20	0.26	0.22	0.19	0.15
	-1.12	-0.88	-1.01	-0.76	-1.92	-1.63
Non-Hodgkin's Lymphomas	2.12	1.79	2.26	1.91	1.17	1.04
	2.15	1.81	2.28	1.93	1.19	1.06
	-1.10	-1.03	-0.99	-0.92	-1.93	-1.73
Myeloma	0.66	0.54	0.65	0.49	0.89	0.93
	0.67	0.54	0.66	0.50	0.91	0.94
	-1.19	-1.13	-1.09	-1.02	-2.00	-1.87
Leukemias	1.45	1.03	1.55	1.09	0.89	0.74
	1.47	1.04	1.56	1.10	0.91	0.75
	-1.05	-0.91	-0.98	-0.86	-1.96	-1.46
Acute Lymphocytic Leukemia	0.12	0.11	0.13	0.12	0.06	0.06
	0.12	0.11	0.13	0.12	0.07	0.06
	-2.92	-1.42	-2.98	-1.99	-4.54	2.65
Chronic Lymphocytic Leukemia	0.47	0.29	0.51	0.31	0.30	0.19
	0.47	0.29	0.52	0.32	0.30	0.20
	-1.11	-0.96	-0.99	-0.83	-1.99	-1.84
Acute Myeloid Leukemia	0.45	0.35	0.47	0.36	0.29	0.28
	0.45	0.36	0.47	0.37	0.29	0.28
	-0.67	-0.87	-0.62	-0.78	-1.33	-1.80

Note: Invasive cancer only unless specified otherwise

Table 2: (continued) Lifetime Risk (percent) of Being Diagnosed with Cancer by Site, Race and Sex. 12 SEER Areas, 1998-2000. (Compare to Ries, et al. 2003, Table I-15). Each cell has 3 values: PMAJ method, Wun, Merrill, and Feuer (1998) method (WMF), and percent difference= $100(PMAJ - WMF)/WMF$.

Site	All Races		Whites		Blacks	
	Males	Females	Males	Females	Males	Females
Chronic Myeloid Leukemia	0.19	0.14	0.20	0.14	0.13	0.10
	0.20	0.14	0.20	0.14	0.14	0.10
	-1.00	-1.04	-0.84	-0.95	-2.09	-1.70

Note: Invasive cancer only unless specified otherwise