# Prevalence and the variance of state occupancy time

Tim Riffe*

Max Planck Institute for Demographic Research

March 13, 2019

### Abstract

**Background**   Markov reward methods have been proposed to calculate the variance of state occupancy time based on age-structured prevalence and survivorship.

**Objectives**   I aim to clarify the assumptions of this approach, give bounds to its reasonableness, and suggest improvements.

**Methods**   I use simplified toy data to demonstrate concepts, and then I use HRS data for different health conditions to derive some empirical guidelines. I calculate results for extreme cases to provide bounds, simulate variance under simple assumptions, and simulate a more natural inter-individual state distribution.

**Results**   I show that state occupancy variance for Sullivan-style inputs is not identified, and I show where previously proposed methods fall with respect to reasonable bounds, randomly generated variances, and my own opinion about what a reasonable inter-individual state distribution might look like.

**Conclusions**   The variance of state occupancy time is only identified if a) state life trajectories are directly observed or b) a process model that can generate an asymptotic life trajectory distribution, such as an incidence-based Markov model, is specified. Sullivan-calculations of life expectancy do not imply a single variance, and are therefore insufficient to make statements on within-group inter-individual disparities in state occupancies.

**Keywords**   Healthy life expectancy, Sullivan method, Healthy life variance, Lifetable, Health inequality

## Introduction

Health inequality is usually measured between populations by comparing life expectancies, health expectancies, or poor health expectancies. Within populations there is also inequality in health outcomes, either due to violations of the assumption of homogeneity of risk sets or due to random variation between people otherwise subject to the same risk. The life table and related statistics

---

*riffe@demogr.mpg.de

typically assume homogeneity of age trajectories of risk between individuals, and multistate models assume the same within states. Transition-based multistate models are flexible with respect to the patterns they can reproduce, but there is no consensus on best practices for state-spaces and which time dimensions ought to structure transitions. Most practitioners produce estimates with the Sullivan method (1) becuse it is often the only viable option, but some Markov skeptics accept the simpler Sullivan method despite its flaws due its process-free interpretation. For the present, let's take the Sullivan estimate of life expectancy for granted and instead focus on within-population estimates of the variance of state occupancy time implied by the same data inputs.

In this paper I argue that no single variance is implied, that in fact infinitely many variances are possible under minimal Sullivan constraints. Instead further assumptions are required to constrain as necessary and arrive at variance calculations. Two such approaches are available assuming that variance results from 1) a discrete time Bernoulli process and 2) from prevalence as a fixed property within age (2). I first translate these two present approaches to a more familiar demographic notation and interpretation, and illustrate their implications with a worked example. I argue that neither of these approaches reflects health prevalence patterns particularly well, and that this observation is not innocuous for variance estimation. I then give my own best approximation of Sullivan variance given what we know about how health processes work. This is at best a hackish patch for an otherwise intractable problem, but we make the best use of some stationary population properties and some observations of incidence-prevalence relationships to give some hopefully-useful rules of thumb.

At best these rules of thumb will stimulate further efforts to improve them into something that might be adopted in practice, although it would also be a perfectly good outcome if practitioners simply withheld from calculating the variance of state occupancy for Sullivan (health) inputs unless the health process and available approximations are well understood and judged apt. Presently available methods may indeed be apt in other domains, I simply argue here that for many health conditions the assumptions required are injurious to good measurement.

## Sullivan in continuous time

Let's begin with some familiar notation, $\ell(x)$ denotes lifetable survivorship, with an arbitrary radix, or starting population. If $\ell(0) = 1$ then $\ell(x)$ can be interpreted as a probability of surviving from birth until age $x$. $\pi(x)$ is the prevalence of a given condition (healthiness or unhealthiness, disability) at age $x$, and it falls in the range $[0, 1]$, giving it a probability interpretation. The commonly used Sullivan method (1) of calculating healthy life expectancy, $e^H(x)$ is to integrate the product of these two functions.

$$e^H(x) = \frac{1}{\ell(x)} \int_x^\omega \ell(t)\pi(t)\mathrm{d}t \tag{1}$$

Replace $\pi$ with its complement to arrive at the complementary expectancy, $e^U(x)$, such that overall life expectancy for this age is $e(x) = e^U(x) + e^H(x)$, a rudimentary decomposition. Throughout this exposition, and for the sake of comparing methods, I stick to an odd choice of discretization of 1, which

assumes that an age entered is an age fully lived, i.e., that an age entered is a reward earned if one thinks in a markov rewards framework. For the case of (1), this translates to:

$$e^H(x) = \frac{1}{\ell_x} \sum_{t=x}^{\omega} \ell_t \pi_t \quad , \tag{2}$$

where subscripts denote the lower bound of the age interval, and function values are assumed constant through the interval. Age intervals are omitted from notation throughout for clarity. Under normal circumstances, one would prefer to replace $\ell_x$ inside the summation with $L_x$ to account for attrition over the course of the year, but we resist this temptation to keep comparisons simple.

## Two assumptions

If we would like to know something about the variance of state occupancy in a Sullivan setup, then we must make some assumptions about how the state is distributed over individuals. The Sullivan method does not state which between-individual state distribution underlies prevalence in a given age. If $\pi(x) = 0.5$, shall we assign each individual age $x$ a value of 0.5 a so-called *fixed reward*, or shall we assign half of them a value of 1, and the other half a value of 0, a *Bernoulli reward* (2)? Or something else entirely? There are infinitely many ways to distribute state $\pi$ among individuals in age $x$ such that the value $\pi(x)$ is maintained. By extension, the inter-individual distribution of total time spent in state $\pi$ is also infinitely variable, which implies that the variance and other moments of implied state occupancy are not uniquely identified.

Expressing the lifetable as a Markov chain with *rewards* defined as prevalence gives a unique definition for the variance of state occupancy for each of the above assumptions, but it is unclear to me how well-supported these are. Many other reward types are discussed by Caswell & Zarulli (2), and I do not discuss these, nor do I treat the case of multistate populations. Even so, the conclusions reached here will generalize to the case of multistate populations.

### Bernoulli rewards

I first give my own lifetable deconstruction of the matrix algebra approach given to Bernoulli state variance in (2). Average years lived in good health $e^H(x)$ is defined per (1), and it is the first moment of state occupancy, which we can denote $\eta^{(1)}$, where the superscript in parentheses denotes the moment number and is not a power. We continue to calculate the second moment of state occupancy, $\eta^{(2)}$ as:

$$\eta_x^{(2)} = \frac{1}{\ell_x} \sum_{t=x}^{\omega} \ell_t \left[ \pi_t + 2(1 - q_t)\pi_t \eta_t^{(1)} \right] \quad , \tag{3}$$

where $q_x$ is the probability of death in the interval and $\omega$ is the highest age of death. Here notation omits intervals, but we assume to be working in discrete bins of uniform width, where death and transition changes only happen in the moment of interval steps, implying a stepped survival curve and discrete life trajectories. We also assume that the same prevalence applies to each length-of-life bin within each age, ergo that all other unmodeled population strata have

the same prevalence. Then the variance is determined, and it can be calculated as:

$$Var_x = \eta_x^{(2)} - (\eta_x^{(1)})^2 \qquad (4)$$

I highlight two assumptions behind this expression: 1) mortality is independent of whether one is in the prevalent state, and 2) the same Bernoulli prevalence extends to any further stratification of the members in a given age class, $x$.

## Fixed rewards

If instead of assuming that prevalence is partitioned in a binary fashion, and each individual experiences $\pi_x$ fraction of the year in the state, then equation (3) becomes

$$\eta_x^{(2)} = \frac{1}{\ell_x} \sum_{t=x}^{\omega} \ell_t \left[ \pi_t^2 + 2(1 - q_t)\pi_t \eta_t^{(1)} \right] \quad , \qquad (5)$$

and (4) is the same. State occupancy variance under the assumption of fixed rewards also has a more intuitive lifetable expression:

$$Var_x = \frac{1}{l_x} \sum_{x}^{\omega} (\mathcal{P}_t - \eta_x^{(1)})^2 d_t \quad , \qquad (6)$$

where $\mathcal{P}_t$ is the cumulative prevalence from age $x$ to age $t$, and $d_t$ is the probability of survival from age $x$ to age $t$, defined as $(q_t l_t)/l_x$.

# 1 Illustrations of Bernoulli and fixed prevalence

To add intuition to the previous equations I give a worked example of manageable size. Define a population with three discrete age groups of uniform width 1, and assume that state transitions occur only on birthdays (this is relaxable). The prevalent state, say sickness, is coded with 1, whereas its complement is coded with 0. Then all 14 possible life trajectories are given in Fig. 1a. Some prevalence $\pi_x$, its complement, and some lifetable functions are given in Tab 1.
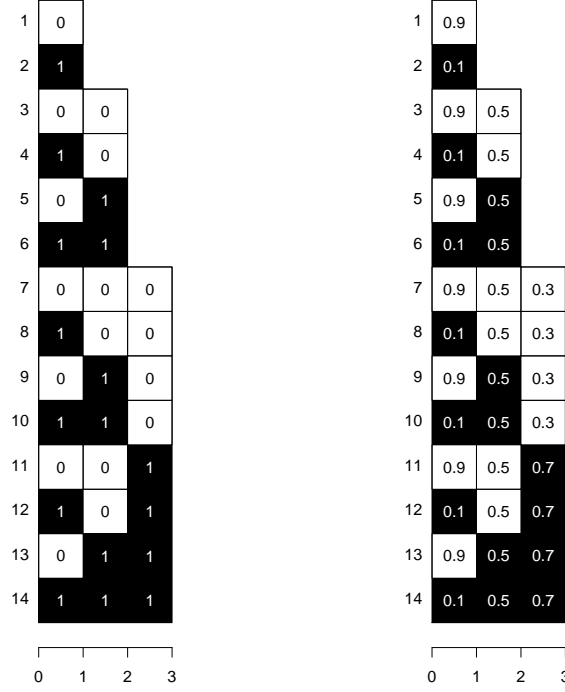
| age | $\pi_x$ | $1 - \pi_x$ | $l_x$ | $d_x$ | $q_x$ |
|---|---|---|---|---|---|
| 0 | 0.1 | 0.9 | 1.0 | 0.2 | 0.2 |
| 1 | 0.5 | 0.5 | 0.8 | 0.5 | 0.625 |
| 2 | 0.7 | 0.3 | 0.3 | 0.3 | 1.0 |

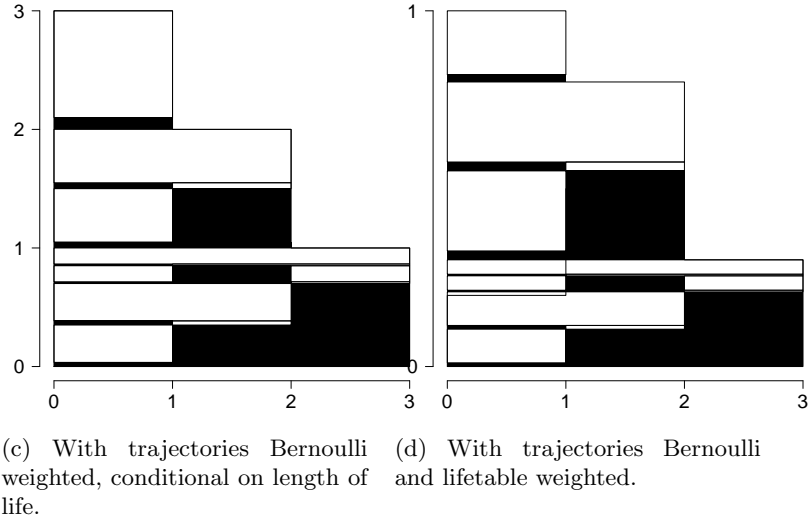Table 1: Parameters for toy examples

## 1.1 Bernoulli illustration

Fig. 1a imputes values of $\pi_x$ and $1 - \pi_x$ from Tab 1 into their respective positions for all 14 possible discrete life trajectories. The product of the proportions in each life sequence gives the probability of experiencing the given sequence of states conditional on total length of life. Fig . 1c shows the same sequences, where height is scaled to the length-of-life conditional probability of observing

the trajectory. The set of sequence probabilities within each length-of-life bin sum to 1, and this is why Fig . 1c forms a set of even steps, each of height 1. These probabilities are not the probabilities-at-birth of observing the given sequence; for this the conditional probability of each sequence must be multiplied by the probability of surviving the respective length of time, in our case given by $d_x$ in Tab. 1. For example, the probability of the 14th trajectory is given by $0.1 \times 0.5 \times 0.7 \times 0.3 = 0.0105$

(a) The binary 2-state trajectory space.

(b) With Bernoulli prevalence cell probabilities.

(c) With trajectories Bernoulli weighted, conditional on length of life.

(d) With trajectories Bernoulli and lifetable weighted.

Figure 1: A depiction of the Bernoulli life trajectories implied by the prevalence and lifetable of Tab. 1. 1a shows all binary discrete life sequences that are possible under these conditions. 1b shows the conditional cell probabilities as drawn from $\pi_x$ and $1 - \pi_x$ in Tab. 1. 1c shows the same sequences where height is defined as the probability of observing the sequence conditional on eventual length of life. 1d shows the same sequences duly weighted by the length of life distribution $d_x$. This is the stationary sequence distribution implied by the discrete Bernoulli process of Tab. 1.

The stationary sequence distribution depicted in Fig. 1d underlies the occupancy time distribution implied by our Bernoulli prevalence example. The heights of each sequence in this graph give the weights to assign to each trajectory, the item we wish to weight is the total number of black boxes per trajectory in Fig 1a, which gives the total time spent in the state of interest. For example, three of the trajectories (1,3,7) do not include time spent in the state, and only one of the sequences (14) consists in three years spent in the state. Fig. 2 shows the stationary distribution of state occupancy, with contributions from each trajectory stacked and indicated with labelled tick marks. This distribution gives an alternate path to the expectation of state occupancy of (1) and the variance of (4).
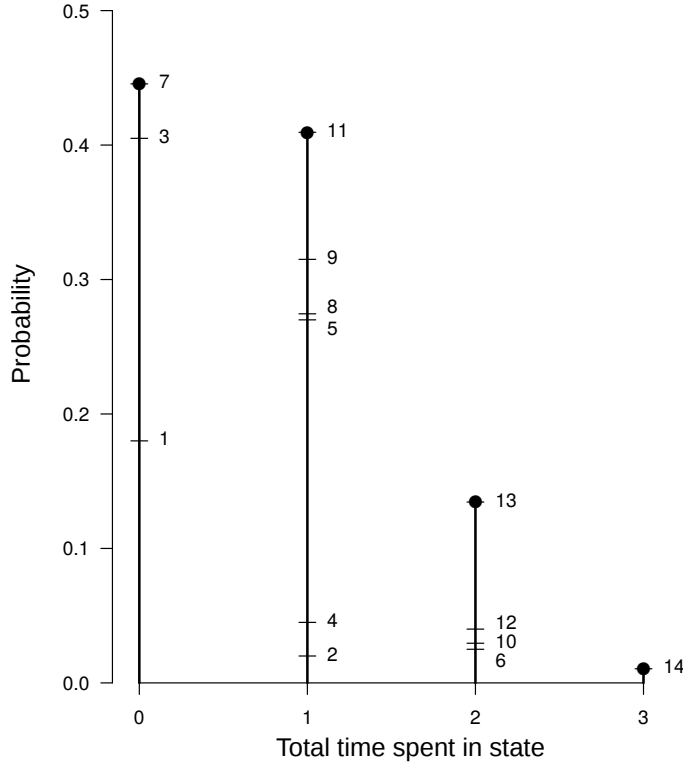


Figure 2: The distribution of occupancy time defined by a Bernoulli interpretation of the prevalence and lifetable given in Tab. 1. The height of each segment gives the proportion of lives the experience the total time given on the $x$ axis. Sub-segments indicated with labelled tick marks indicate the contributions from each sequence in Fig. 1.

In our case the expectation at birth of time spent in the state is 0.71 years under Bernoulli prevalence assumptions, and the variance is 0.5379. These are identical results to those obtained from the matrix methods of (2) if one assumes that a full reward is assigned on entrance into an age/state.

(a) Fixed rewards uniformly distributed within age classes.

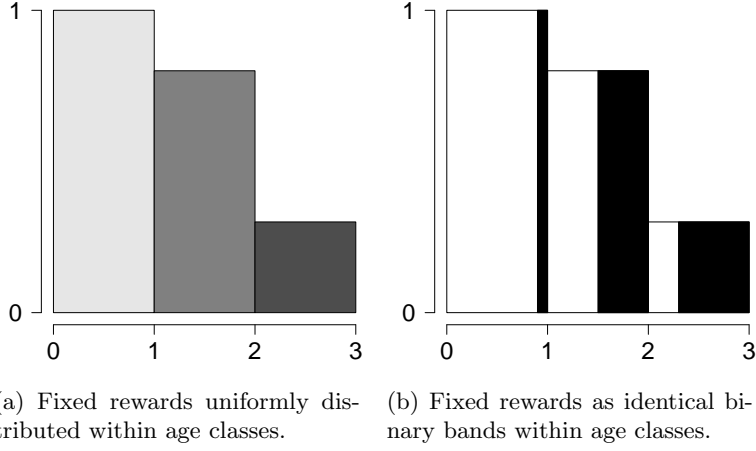(b) Fixed rewards as identical binary bands within age classes.

Figure 3: Two interpretations of fixed rewards following the data of Tab. 1. In Fig. 3a grayscale is proportional to the prevalence value, and in Fig. 3b the state is binary but assigned to an equal fraction of the year within age classes. There are also infinitely may ways to distribute states within age classes such that the same fraction is held

## 1.2 Fixed reward illustration

We can also recycle values from Tab 1 to illustrate the meaning of fixed rewards. In this case $\mathcal{P}$ obtains the values $[0.1, 0.6, 1.3]$, and $d_x$ are our weights used to average these, giving us the same expectation of 0.71 years, but (6) implies a variance of 0.1849, quite different from the Bernoulli assumption. Under fixed rewards, one might redraw Fig. 1d as of a 2d uniform distribution under $l_x$ and within age bins, as in Fig. 3a, or else as binary year fractions, as in Fig. 3b. Indeed there are infinitely many ways to satisfy the constraints of fixed rewards by moving prevalence mass around within horizontal strips of age bins for this discrete case. Key for this case is that the cumulative value along any horizontal bisection is identical on birthday transitions.

# 2 Alternative prevalence scenarios and their variances

The purpose of

# Competing interests

The authors declare that they have no competing interests.

# Author's contributions

TR did everything so far.

# Acknowledgements

# References

[1] Sullivan, D.F.: A single index of mortality and morbidity. HSMHA health reports **86**(4), 347 (1971)

[2] Caswell, H., Zarulli, V.: Matrix methods in health demography: a new approach to the stochastic analysis of healthy longevity and dalys. Population health metrics **16**(1), 8 (2018)