

## Semi-Markov modeling for disease incidence risk and duration

Antoine Soetewey, Catherine Legrand, Michel Denuit & Geert Silversmit

**To cite this article:** Antoine Soetewey, Catherine Legrand, Michel Denuit & Geert Silversmit (2025) Semi-Markov modeling for disease incidence risk and duration, *Biostatistics & Epidemiology*, 9:1, e2517916, DOI: [10.1080/24709360.2025.2517916](https://doi.org/10.1080/24709360.2025.2517916)

**To link to this article:** <https://doi.org/10.1080/24709360.2025.2517916>



Published online: 15 Jun 2025.



Submit your article to this journal [↗](#)



View related articles [↗](#)



View Crossmark data [↗](#)



# Semi-Markov modeling for disease incidence risk and duration

Antoine Soetewey <sup>a</sup>, Catherine Legrand<sup>a</sup>, Michel Denuit<sup>a</sup> and Geert Silversmit<sup>b</sup>

<sup>a</sup>Institute of Statistics, Biostatistics and Actuarial Sciences, Louvain Institute of Data Analysis and Modeling in Economics and Statistics, UCLouvain, Louvain-la-Neuve, Belgium; <sup>b</sup>Cancer Research Department, Belgian Cancer Registry, Brussels, Belgium

## ABSTRACT

Over the last decade, the number of years of life lost (YLL) has become a popular metric in biostatistics for assessing mortality and life expectancy discrepancies between patient cohorts and the general population. Using data from the Belgian Cancer Registry (161,007 cases of melanoma, thyroid, and female breast cancer), a three-state (healthy–cancer–death) illness-death model is used to illustrate how it can be applied to cancer registry data to estimate the incidence risk and YLL due to cancer at various ages of diagnosis and survival times post-diagnosis. The probabilities of being diagnosed with cancer over the next 20 years remain low for melanoma and thyroid cancers for both sexes but considerably increase with age for female breast cancer. YLL before age 70 due to cancer is highest for early diagnoses of female breast cancer but peaks at later ages for melanoma and thyroid cancers. Additionally, male patients generally experience higher YLL before age 70 due to cancer than females for melanoma and thyroid cancers. For patients surviving 10 years post-diagnosis, YLL before age 70 due to cancer remains below one year for melanoma and thyroid cancers, suggesting a limited impact on life expectancy compared to the general population.

## ARTICLE HISTORY

Received 27 September 2024  
Accepted 2 June 2025

## KEYWORDS

Years of life lost; multi-state models; critical illness; cancer mortality

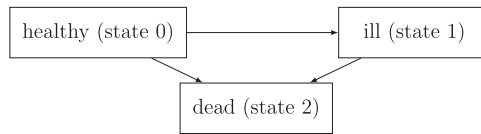
## 1. Introduction

Over the last decade, the number of years of life lost (YLL) has become a popular tool in biostatistics and epidemiology to measure discrepancies in life expectancy or mortality. The idea behind YLL is to quantify the number of years of life a specific cohort of patients has lost due, for example, to a given disease, compared to the general population. This measure, as defined by Andersen [1] and Andersen et al. [2], has the advantage (compared to others such as the hazard ratio or excess hazard) that it is measured on a time metric (usually in years) making its interpretation easy for policy-makers and meaningful for gauging public health outcomes [3].

It was first introduced to measure the reduction in life expectancy for a group of individuals compared to a hypothetical cohort where no one dies before a given age [1]. However, in most situations, it may seem more natural to measure the reduction in life expectancy for a group of individuals compared to a reference population (where some years of life are lost because of some standard or background mortality rates). In this sense, YLL can be used to estimate the number of years a specific cohort of patients (cancer patients, for instance) is expected to lose compared to the general population (i.e. the reference population to which the cancer cohort is compared). The difference between the life expectancy of the general population and that of the considered cohort of patients corresponds to YLL. This measure is sometimes referred to as excess YLL because it is the number of years of life patients lose in excess of that seen in the general population. The greater this measure, the more important the societal burden of the disease or condition.

Similarly to the excess hazard, information about the cause of death is not required to estimate YLL, making it a practical measure for population-based studies in which the cause of death is often unavailable or unreliable [4]. There are two types of YLL. First, the number of years of life lost by the entire cohort, which can be denoted  $YLL^c$ , and which is of interest if one wants to estimate at one point in time the global number of years of life lost due to a particular disease (see for instance Aragon et al. [5] who rank leading causes of premature death based on the total number of years of life lost due to each cause). This may be used to answer questions such as, ‘How many years of life are lost in the population due to cancer?’ [6]. It is of great interest to economists, governments and policy-makers to determine which condition or disease has the largest negative impact on citizens and society as a whole (for resource allocation, public health priorities, cancer control progress, etc.). Second, the number of years of life lost (on average) per individual, which we denote  $YLL^i$ , and which quantifies how many years of life a patient is expected to lose (see for example Belot et al. [7] or Latouche et al. [3]). It answers questions such as ‘How much does the life expectancy of an individual on average change if diagnosed with cancer?’. See examples with common cancers in Chu et al. [8] who measure health impacts on society using  $YLL^i$ . In this situation,  $YLL^i$  can be seen as an average per person, whereas  $YLL^c$  can be seen as the sum of the years of life lost for each individual in a patient cohort. See a comprehensive overview of the years difference measures in Manevski et al. [9]. Note that individuals do not necessarily lose years compared to the general population; they may also gain years. This is the case, for instance, in the study of the long-term survival of elite athletes for which survival may be better than that of the general population [10].

From a general point of view, the major advantages of  $YLL^c$  and  $YLL^i$  are that (i) it is measured on a time metric (usually in years), facilitating its interpretation and communication [11,12], (ii) information on the cause of death is not needed to estimate it, and (iii) it can be computed for any time horizon and for a comprehensive list of causes of death. Andersen [13] suggested several measures of life years lost among patients with a given disease in the framework of a (Markov or non-Markov) illness-death model, illustrated using data on Danish male patients with bipolar disorder. The main goal of the present study is to demonstrate how  $YLL^i$  can be easily estimated from a multi-state model and what the advantages are of doing so, with a focus on two applications using data on Belgian cancer patients. Their use in the context of the right to be forgotten will also be discussed.



**Figure 1.** Visual representation of the ‘illness-death model’ without recovery for cancer patients.

Multi-state models (MSM) are a powerful statistical approach to study the evolution of individuals between several ‘states’ (see Andersen et al. [14] and Hougaard [15] for a general review). MSM can be seen as an extension of classical survival analysis, in which only the transition from being alive to being dead is considered [16–18]. Unlike classical survival models, MSM are used to model processes which go from an initial state (for instance ‘healthy’) to a terminal (also referred to as absorbing) state (for example ‘dead’), but where more than two states are considered, some being transient. For example, considering that the ‘healthy’ state is portioned into two or more intermediate states corresponding to specific stages of a disease [19]. Thus, MSM offer a complete and informative representation of the occurrence of intermediate events on the pathway to some final event, notably via transition probabilities which have a natural interpretation [20,21].

In this paper, a three-state model, assuming that an individual can either be ‘healthy’, ‘ill’ (diagnosed with cancer), or ‘dead’ is considered. We will see that in our context, we actually only need to consider transitions from healthy to ill, healthy to dead, and ill to dead. While excluding the possibility to transit from ill back to healthy can be interpreted as assuming that cancer is a permanent condition (which is debatable), we actually decided not to consider it following the parsimony principle since it would not bring any useful information in our context. Indeed, as it will be shown later, in our type of applications, distinguishing the health status of the patient between diagnosis and death is actually not required. This non-reversibility greatly simplifies the computations, as in this case, our three-state process is hierarchical and trajectories can be described in terms of just a few random variables [22]. See Figure 1 for a visual representation of the model, often referred to in the literature as the ‘(three-state) illness-death model’ without recovery. More advanced types of MSM (known as reversible MSM) can be used in case recoveries are possible and have to be taken into account for the application considered. Note that this three-state model is, in its mathematical concept, similar to the well-known SIR model (susceptible – infected – recovered) in epidemiology [23,24]. The difference with our three-state illness-death model is that a susceptible individual must go through the infectious state before being recovered; he/she cannot go directly from ‘susceptible’ to ‘recovered’.

The key contribution of this paper is thus to illustrate how disease incidence risk and  $YLL^i$  can be estimated based on a Semi-Markov three-state MSM using cancer registry data, and what type of useful information can be obtained out of it. Furthermore, the main advantage of computing these quantities in a Semi-Markov context is that it allows taking into account the number of years a patient survived after diagnosis. To the best of our knowledge, most studies refer to the number of years of life lost at the time of diagnosis, without taking the time survived since diagnosis into consideration. This is a major difference, given that time spent in the ill state is known to have an influence on survival for cancer patients.

**Table 1.** Number of persons diagnosed with melanoma, thyroid and female breast cancer in Belgium between 2004 and 2020 (BCR data) by sex, site and age group, together with the percentage lost to follow-up and the number of deaths.

Sex	Cancer site	Age at diagnosis	Lost to follow-up	Number of included cases	Number of deaths
Men	Melanoma	20–34	3.72%	969	94
		35–49	2.66%	3,266	404
		50–69	1.70%	7,460	1,583
Total				11,695	2,081
Men	Thyroid	20–34	4.10%	366	6
		35–49	3.12%	961	67
		50–69	2.14%	1,773	379
Total				3,100	452
Women	Melanoma	20–34	3.62%	2,488	78
		35–49	1.47%	6,137	382
		50–69	1.35%	8,893	1,112
Total				17,518	1,572
Women	Thyroid	20–34	3.80%	1,607	14
		35–49	2.67%	3,449	107
		50–69	2.06%	4,085	484
Total				9,141	605
Women	Breast	20–34	2.76%	3,112	502
		35–49	1.78%	32,743	4,058
		50–69	1.31%	83,698	15,946
Total				119,553	20,506

The remainder of this paper is laid out as follows. Section 2 presents the data used to perform the present study. Section 3 details the methods and tools, with a focus on Semi-Markov MSM. Section 4 illustrates two useful MSM-based health indices. The final section (Section 5) concludes the paper with a discussion.

## 2. Data

For these applications, the data available from the Belgian Cancer Registry (BCR) are considered. The BCR is a national population-based cancer registry collecting data on all new cancer diagnoses in Belgium since the incidence year 2004. For the execution of this main task, the BCR relies on its own specific legislation (more information can be found on the BCR website, at [kankerregister.org](http://kankerregister.org)).

To illustrate our work, the methods were applied to three cancer types: melanoma (ICD-10 C43), thyroid (ICD-10 C73) and female breast (ICD-10 C50) cancer. These three cancer sites have been selected to evaluate the proposed method in different scenarios. Melanoma and thyroid cancer patients are known to have a limited excess hazard compared to the general population and a high survival probability [25–29]. The situation for female breast cancer patients is different, with usually a high survival probability in the first years after the date of diagnosis before it eventually decreases due to late cancer recurrences [30]. Only female breast cancer is considered due to the limited number of male breast cancer cases.

Out of a total of 161,007 cases, melanoma, thyroid, and breast cancer represent, respectively, 29,213 (18.1%), 12,241 (7.6%), and 119,553 (74.3%) cases diagnosed between 2004 and 2020. Patients were followed up until April 11, 2022, resulting in follow-up periods ranging from 2 to 18 years. Only one record per patient (with the earliest incidence date) within each cancer site was kept for patients with multiple primary diagnoses. A minority of patients without a national security number were excluded from the analysis. Patients

lost to follow-up (mostly due to moving abroad) and patients still alive at the end of the follow-up period were treated as censored observations.

Table 1 summarizes the number of included cases, number and proportion of deaths and percentage lost to follow-up before April 11, 2022 per type of cancer, sex and age group. The fraction of patients lost to follow-up per subgroup varied from 1.31% for women with breast cancer aged 50–69 to 4.1% for male thyroid cancer patients aged 20–34. The total fraction of patients lost to follow-up cases, regardless of sex, site or age group, was 1.64%. Moreover, mean age at diagnosis was 50.5 years (standard deviation ( $SD$ ) = 12.1), 48.1 years ( $SD$  = 12.4) and 54.6 years ( $SD$  = 9.5) for melanoma, thyroid and breast cancer, respectively.

In order to estimate the number of years of life lost, mortality in the cancer cohort must be compared to the expected mortality in the general population. Mortality in the general population is therefore also needed. The complete Belgian population is also required to estimate the transition from healthy to ill (which cannot be estimated based on the cancer registry data). These general population data come from the Belgian population life tables, which are available from Statbel (the Belgian statistical office) and can be freely downloaded from the website [statbel.fgov.be](http://statbel.fgov.be).

Note that as population life tables take into account all deaths, those due to the cancer of interest are also included. Nonetheless, it is commonly assumed that the fact that population life tables include cancer mortality is not an issue since mortality for a given cancer represents only a small fraction of the overall mortality. Correcting for this mortality of the cancer being studied has, in practice, an insignificant effect on the survival of the general population [31,32].

### 3. MSM and YLL for cancer patients

A MSM, which is a model for time-to-event data, consists of states and transitions between pairs of states that reflect the disease and death mechanisms in medical applications. Main motivations for using a MSM are often to obtain (i) more biological insight into the disease or recovery process of a patient, and (ii) more accurate predictions than standard models neglecting intermediate states. Indeed, by incorporating intermediate events, predictions are adjusted in the course of time, giving more precise information about survival duration [16,17].

When considering MSM, the following concepts must be distinguished: (1) Markovian and Semi-Markovian, and (2) homogeneous and non-homogeneous models. These concepts can be defined as follows:

- Markovian: What happens next only depends on the current state, not on what happened before;
- Semi-Markovian: What happens next depends on the current state and how long ago it was reached (so the duration in that state);
- Homogeneous or time-homogeneous: Transition between states do not depend on time (but time seen as age and not duration in the state, hence the name *time*-homogeneous);

- Non-homogeneous or time-inhomogeneous: Transition between states may depend on time (seen as age, not duration).

For a non-homogeneous Markov model, the time until the next state is allowed to depend on the current state and the individual's age (i.e. time). For a homogeneous semi-Markov model, the time until the next state is allowed to depend on the current state and the time since he/she entered this state (i.e. duration). For a non-homogeneous Semi-Markov model, both aspects (time and duration) are combined: the time until the next state is allowed to depend on the current state, the time since he/she entered this state, and his/her age.

Thus, in our context, assuming a homogeneous Markov illness-death model would mean to consider that the expected length of stay in the ill state of a cancer patient depends only on the current state. In other words, it would assume that two cancer patients have the same expected length of stay in the ill state (and thus, the same mortality), even if one has been diagnosed for one year and the other for 10 years. However, it is known that mortality for cancer patients (and thus expected length of stay in the ill state) varies with time since diagnosis (and thus sojourn time) [33]. Therefore, the Markovian assumption does not hold for our situation, and a Semi-Markov assumption taking also into consideration the time spent in the ill state is preferable. Moreover, the non-homogeneous assumption is also preferable as transitions may depend on the patient's age. In this non-homogeneous Semi-Markov case (also known as general Semi-Markov), the expected length of stay in the ill state of a cancer patient will thus depend on both the age and the time since diagnosis. This assumption is important because it allows us to update the patient's life expectancy conditional on the fact that he/she survived up to that time and a given specific age. This is the reason why our calculations are performed in the context of a non-homogeneous Semi-Markov illness-death model.

The whole process from birth to death of any individual can be defined formally as a random process over time  $X = [X(t), t \geq 0]$ , where  $X(t)$  gives the state occupied at age  $t$ . Here,  $t$  corresponds to the time since birth. In the irreversible illness-death process depicted in Figure 1,  $X(t)$  has values in state space  $\mathcal{S} = \{0, 1, 2\}$  where state 0 corresponds to the 'healthy' state, state 1 to the 'ill' state, and state 2 to the 'dead' state. Individuals are initially with no cancer detected, thus considered as healthy. Then, they may be diagnosed with cancer and die, or they may die without having been diagnosed with cancer.

More formally, let's denote by  $T_{ij}$  the age at which the patient moves from state  $i$  to state  $j$ . For patients diagnosed with cancer at age  $T_{01}$  and who died at age  $T_{12}$ , we have

$$\begin{aligned} X(t) &= 0 & 0 \leq t < T_{01}, \\ X(t) &= 1 & T_{01} \leq t < T_{12} \text{ and} \\ X(t) &= 2 & t \geq T_{12}. \end{aligned}$$

For patients without cancer who died at age  $T_{02}$ , we have

$$\begin{aligned} X(t) &= 0 & 0 \leq t < T_{02} \text{ and} \\ X(t) &= 2 & t \geq T_{02}. \end{aligned}$$

Remember that it is assumed that a cancer patient stays in the ‘ill’ state until he/she dies (i.e. the transition from state 1 to state 0 is not allowed). So, in fact, the state ‘ill’ should rather be understood as ‘having been diagnosed with cancer’.

In our context, we have to assume that the time spent in state  $i$  influences transition to the next state. Therefore, the random variable  $Z(t)$  is introduced, and defined as the time spent in the state occupied at time  $t$ . Formally,

$$Z(t) = \max\{z \leq t | X(t) = X(t - h) \text{ for all } 0 \leq h \leq z\}.$$

For an individual in state  $i$  at time  $t$ ,  $Z(t)$  is the time since entry in the state (i.e. time from birth for  $i = 0$  and time from diagnosis for  $i = 1$ ). Henceforth, we work under the Semi-Markov assumption: the current state  $X(t)$  and the time  $Z(t)$  spent in the current state influence future transitions. This means that the stochastic process  $[(X(t), Z(t)), t \geq 0]$  is a Markov process.

A fundamental concept in multi-state models is the transition intensities, which govern movements between the different states depending on the state currently occupied and the sojourn time. The following transition intensities fully describe the process in an illness-death model:

$$\alpha_{01}(t) = \lim_{h \rightarrow 0} \frac{P[X(t+h) = 1 | X(t) = 0]}{h} \quad (1)$$

$$\alpha_{02}(t) = \lim_{h \rightarrow 0} \frac{P[X(t+h) = 2 | X(t) = 0]}{h} \quad (2)$$

$$\alpha_{12}(t; z) = \lim_{h \rightarrow 0} \frac{P[X(t+h) = 2 | X(t) = 1, Z(t) = z]}{h} \quad (3)$$

where  $\alpha_{ij}(\cdot)$  are the transition intensities between state  $i$  and state  $j$  ( $i = 0, 1; j = 1, 2$ ). Transition intensities from state 0 depend on the time spent in that initial state through attained age. Furthermore, there is an influence of the duration of stay in state 1 so that transition intensities from state 1 depend on both attained age and time  $z$  since diagnosis. In our context,  $\alpha_{01}(\cdot)$ ,  $\alpha_{02}(\cdot)$  and  $\alpha_{12}(\cdot; \cdot)$  are, respectively, the intensity of developing cancer, the death intensity without cancer, and the death intensity with cancer. Also, the exit intensity from state 0 is denoted  $\alpha_{0\bullet}(t)$ , that is,  $\alpha_{0\bullet}(t) = \alpha_{01}(t) + \alpha_{02}(t)$ .

Transition probabilities are meaningful to estimate in addition to transition intensities. Considering an individual who is healthy at age  $t$ , that is, who is in state 0 at time  $t$ , the probability of being in state 1 at time  $t+h$  is denoted as

$$p_{01}(t, t+h) = P[X(t+h) = 1 | X(t) = 0],$$

the probability of being in state 2 at time  $t+h$  is denoted as

$$p_{02}(t, t+h) = P[X(t+h) = 2 | X(t) = 0],$$

and the probability of still being in state 0 at time  $t+h$  is denoted as

$$p_{00}(t, t+h) = P[X(t+h) = 0 | X(t) = 0].$$

Since the time spent in state 1 influences future transitions, the random variable  $Z(t)$  also enters the transition probabilities from that state. Precisely, considering an ill individual



diagnosed at age  $T_{01}$  and aged  $t = T_{01} + z$ , that is, who is in state 1 since the last  $z = t - T_{01}$  years, the probability of being in state 2 at time  $t + h$  is denoted as

$$p_{12}(t, t + h; z) = P[X(t + h) = 2 | X(t) = 1, Z(t) = z]$$

and the probability of still being in state 1 at time  $t + h$  is denoted as

$$p_{11}(t, t + h; z) = P[X(t + h) = 1 | X(t) = 1, Z(t) = z].$$

As explained before, we do not need to consider the possibility to move back to the initial state, or to transition to an intermediate ‘recovery’ state for our applications. Hence, transition probabilities  $p_{00}(t, t + h)$  and  $p_{11}(t, t + h; z)$  are in reality sojourn probabilities, i.e.

$$p_{00}(t, t + h) = P[X(t + u) = 0 \text{ for all } 0 < u \leq h | X(t) = 0]$$

$$p_{11}(t, t + h; z) = P[X(t + u) = 1 \text{ for all } 0 < u \leq h | X(t) = 1, Z(t) = z].$$

More generally, transition probabilities can be rewritten as

$$p_{ij}(t, t + h; z) = P[X(t + h) = j | X(t) = i, Z(t) = z] \quad \forall i, j \in \mathcal{S}$$

and transition intensities can be rewritten as

$$\begin{aligned} \alpha_{ij}(t; z) &= \lim_{h \rightarrow 0} \frac{P[X(t + h) = j | X(t) = i, Z(t) = z]}{h} \quad \forall i, j \in \mathcal{S} \\ &= \lim_{h \rightarrow 0} \frac{p_{ij}(t, t + h; z)}{h} \quad \forall i, j \in \mathcal{S}. \end{aligned}$$

While these transition probabilities and transition intensities give useful information on the evolution of the individuals, obtaining information about survival duration is also of great interest for clinicians and patients. Life expectancy at birth is a metric widely used in demography to measure the length of survival present in a population, and corresponds to the average number of years an individual is expected to live from birth (given that mortality rates remain constant in the future) [34–36]. Moreover, remaining life expectancy is the average number of remaining years an individual is expected to live, starting from a certain age instead of birth. By computing remaining life expectancy starting at a certain age, it is meant to be conditional on survival to that certain age. If, in addition to estimate life expectancy from a given age instead of birth, it is also estimated up to a given time horizon, it is known as the restricted mean lifetime and it can be interpreted as the average number of years an individual is expected to live between two specific ages. In this paper, we will be particularly interested in taking into account both a starting age different than birth (so conditional on survival to some ages after birth) and a finite time horizon (so considering a given upper age  $\tau$ ). See Section 4 for more details about the choices of the starting age and  $\tau$ .

As mentioned earlier, the number of years of life lost can be seen either at the cohort level ( $YLL^c$ ) or at the individual level ( $YLL^i$ ). When applied to cancer patients, on the one hand,  $YLL^c$  represents the total number of years of life lost by the cancer cohort. This is useful to compare, for instance, the societal burden of cancer with other diseases or between different countries. On the other hand,  $YLL^i$  can be interpreted as the average

number of years of life lost that a cancer patient experiences from the time of diagnosis in comparison to a healthy individual of the same age (and possibly sex, year and other covariates such as ethnicity or socio-economic factors). This latter definition resonates more in the patient-clinician communication. In this paper, it is the  $YLL^i$  which is chosen and illustrated.

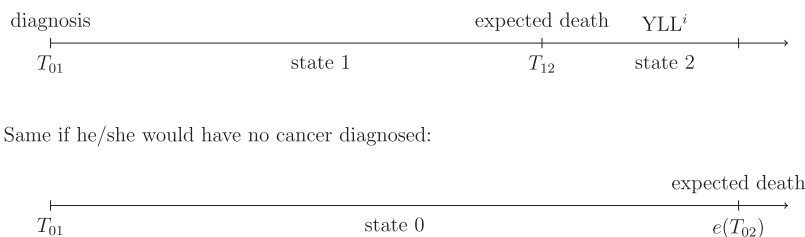
$YLL^i$  in a certain time interval is the sum of life years lost due to (i) population mortality (governed by mortality rates in that reference population) and due to (ii) the cancer of interest. This quantity can be computed based on the estimated survival observed in the general population minus the estimated survival in the cohort of cancer patients considered. Formally, the number of years of life lost due to cancer starting from the age at diagnosis  $T_{01}$  until age  $\tau$  is defined as

$$YLL^i(T_{01}) = \int_{T_{01}}^{\tau} \hat{S}_P(t) dt - \int_{T_{01}}^{\tau} \hat{S}_C(s) ds \quad (4)$$

where  $\hat{S}_P(\cdot)$  denotes the classical survival function estimated via the population mortality rates, and  $\hat{S}_C(\cdot)$  is the cancer survival curve (in general, estimated via the non-parametric Kaplan-Meier [37] method, but it could be estimated via another method as well) [7].

The lower bound  $T_{01}$  in the integrals represents age at diagnosis (so conditional on survival to age  $T_{01}$ ) and the upper bound  $\tau$  corresponds to the time horizon, chosen arbitrarily or such that it matches a certain cut-off. The number of years of life lost uses the age at diagnosis for each cancer patient as its starting point and estimates the expected remaining lifetime at that age using age-specific mortality rates. The number of years of life lost due to cancer is then estimated by matching the expected remaining lifetime for someone diagnosed with cancer with the life expectancy in the general population at that specific age. Age-specific mortality rates and life expectancy in the general population are generally available through life tables (as they are usually stratified by age). For life tables that are stratified by sex in addition to age, the number of years of life lost can be used to compare cancer patients to the general population of the same sex and age.

Our objective is to demonstrate how this quantity can be estimated from our MSM. The idea here is to start from our MSM to compute  $YLL^i$  using life expectancy, probabilities of developing the disease within a specific time period, and expected lengths of stay in each of the different states (also referred to in the literature as the mean sojourn time, see Jackson [38]). Following Equation (4), estimation of  $YLL^i$  via a MSM starting from the age at diagnosis is denoted  $YLL_{MSM}^i(T_{01})$  and corresponds to the number of years of life lost at the time of diagnosis for someone diagnosed at age  $T_{01}$ . Figure 2 illustrates the approach, where  $e(T_{02})$  is the remaining life expectancy until the expected death of a healthy individual. One could argue that it does not make sense to speak about age at diagnosis  $T_{01}$  if the person has no cancer. However, in fact, we compare what would have happened to a patient diagnosed at age  $T_{01}$  if he/she would not have had cancer at the time he/she was actually diagnosed. We are now considering the hypothetical trajectory that a patient diagnosed at age  $T_{01}$  would have had if he/she had not had cancer and therefore if he/she had remained in state 0.



**Figure 2.** Representation of MSM to estimate  $YLL^i$  from diagnosis.

In the context of a Semi-Markov multi-state model, the remaining life expectancy for a cancer patient diagnosed at age  $T_{01}$ , given the time  $z$  elapsed since diagnosis is

$$e_{11}^{\tau}(T_{01} + z; z) = \int_{T_{01}+z}^{\tau} p_{11}(T_{01} + z, s; z) ds. \quad (5)$$

Since  $t = T_{01} + z$ , Equation (5) becomes

$$e_{11}^{\tau}(t; z) = \int_t^{\tau} p_{11}(t, s; z) ds. \quad (6)$$

Following Figure 2, to define  $YLL_{MSM}^i(T_{01})$  in a Semi-Markov context we add the conditioning on  $z$  to have the number of YLL for someone diagnosed at age  $T_{01}$  but who would have already survived with his/her cancer for  $z$  years. In that case, we obviously have to update the life expectancy for the cancer patient (the fact that he/she lived already for  $z$  years gives information on his/her life expectancy) and do the same for his 'healthy' counterpart. This is denoted  $YLL_{MSM}^i(T_{01}; z)$  and is defined as follows

$$\begin{aligned} YLL_{MSM}^i(T_{01}; z) &= \text{remaining life expectancy at age } T_{01} \text{ for a healthy individual} \\ &\quad - \text{remaining life expectancy for a cancer patient diagnosed at age } T_{01}, \\ &\quad \text{given the time } z \text{ elapsed since diagnosis} \\ &= e(T_{01}) - e_{11}^{\tau}(T_{01} + z; z). \end{aligned} \quad (7)$$

Remaining life expectancy at age  $T_{01}$  for a healthy individual is usually found with life tables and population mortality rates. Here, the expected remaining lifetime until age  $\tau$  for someone diagnosed with cancer is matched with the  $\tau$ -restricted life expectancy in the general population at that specific age.

As often the case in practice, transition intensities are assumed to be piecewise constant in order to ease calculations but also given the information available in cancer registries. In that case, transition intensities are easily estimated by the ratio of the observed number of transitions (diagnosis or death) to the corresponding exposure (in the state to be left) [33]. When (annual) piecewise constant transition intensities are considered, we get

$$e_{11}^{\tau}(t; z) = \sum_{k=0}^{\tau-t-1} \exp\left(-\sum_{l=0}^{k-1} \alpha_{12}(t+l; z+l)\right) \frac{1 - \exp(-\alpha_{12}(t+k; z+k))}{\alpha_{12}(t+k; z+k)} \quad (8)$$

with  $\sum_{l=0}^{k-1} \alpha_{12}(t+l; z+l) = 0$  if  $l = 0$ . The development of  $e_{11}^{\tau}(t; z)$  is explained in Appendix.

Remember that  $YLL_{MSM}^i(T_{01}; z)$  is defined at an individual level. In this sense,  $YLL_{MSM}^i(T_{01}; z)$  quantifies the number of years of life a patient diagnosed with cancer  $z$  years ago is expected to lose compared to someone who will never develop the disease. It can be seen as an insightful health indicator, complementary to other health indicators already used by clinicians and policy-makers. Indeed, it can be used to communicate about a patient's survival, but it can also serve as a measure of the burden of cancer for the whole society (with comparisons between diseases, countries or throughout the years, for example).

#### 4. Derived health indices: case studies for three cancer types

One of the main advantages of estimating  $YLL^i$  from a MSM is that several health indicators could be derived from it. The focus here is put on two different applications to illustrate its potential uses; (i) the cancer incidence risk and (ii) the number of years of life lost due to cancer given a certain time spent after diagnosis. Note that the first health indicator requires the 3 states. However, regarding the second one, we consider an individual of age  $T_{01}$  at diagnosis. This means that state 0 is no longer needed, since we are already in state 1. Also note that incidence refers to the number of new cases of a disease over a specified period, and can be expressed as a risk or an incidence rate [39]. We are interested in the former, that is, the incidence risk that a subject within a population will develop a given cancer, over a specified follow-up period. This incidence risk, expressed as a probability, can be interpreted as an estimation of the risk of cancer in an individual subject over a certain time frame.

For these applications, our analyses are limited to patients aged 20 to 69 years old at the time of diagnosis for two main reasons. First, childhood cancers can be seen as a category of cancer on their own and are often studied separately because they differ greatly from adult cancers. Second,  $\tau$  has been set to 70 years, an age at which persons were censored if they had not died before to focus on active life from a public policy perspective. The estimate of  $YLL^i$  has therefore to be interpreted as the number of years of life lost before that specific age. This is analogous to the  $\tau$ -restricted mean lifetime, which can be interpreted as the average number of years lived before time  $\tau$ . Note that the choice of  $\tau$  is arbitrary. In some settings, researchers may be interested in  $YLL^i$  before retirement age applicable in a country. In our case, we are interested in potential implications for insurers in the context of the right to be forgotten, hence the upper limit of 70 years (people aged above are unlikely to contract a loan). Note the distinction between the maximum age at diagnosis (69 years) and the upper age limit  $\tau$  (70 years). This difference is explained by the fact that we include patients who have been diagnosed before their 70<sup>th</sup> birthday (and thus who are still 69 years old at the time of diagnosis), while we are interested in the number of years of life lost before the age of 70 due to cancer. This is to avoid the possibility that a patient is diagnosed between his or her 70<sup>th</sup> and 71<sup>th</sup> birthday, while computing the number of years of life lost before he or she has reached the age of 70 years.

To display our results, the time since diagnosis  $z$  is set to 0, 5, and 10 years.  $YLL^i(T_{01}; 0)$  corresponds to the number of years of life lost due to cancer at the time of diagnosis for a patient diagnosed at age  $T_{01}$ .  $YLL^i(T_{01}; 5)$  and  $YLL^i(T_{01}; 10)$  correspond to the same quantity computed after having survived to the cancer for respectively 5 and 10 years. A  $z$  of 5 and 10 years after diagnosis has been chosen to cover a relatively large period of time

after diagnosis, while we refrain from setting it higher due to the limited follow-up period in our data.

#### 4.1. Incidence risk

We start the applications with the estimation of the probability for the population of age  $t$  to be diagnosed with each of the three types of cancer we consider between age  $t$  and  $t + n$ . In other words, the probability of being diagnosed with cancer for a healthy individual aged  $t$  over the next  $n$  years is computed. This measure, similar to the incidence risk and again assuming yearly-constant intensities, is defined based on a MSM as follows

$$p_{01}(t, t + n) = \sum_{k=0}^{n-1} \alpha_{01}(t + k) \exp\left(-\sum_{l=0}^{k-1} \alpha_{0\bullet}(t + l)\right) \frac{1 - \exp(-\alpha_{0\bullet}(t + k))}{\alpha_{0\bullet}(t + k)} \quad (9)$$

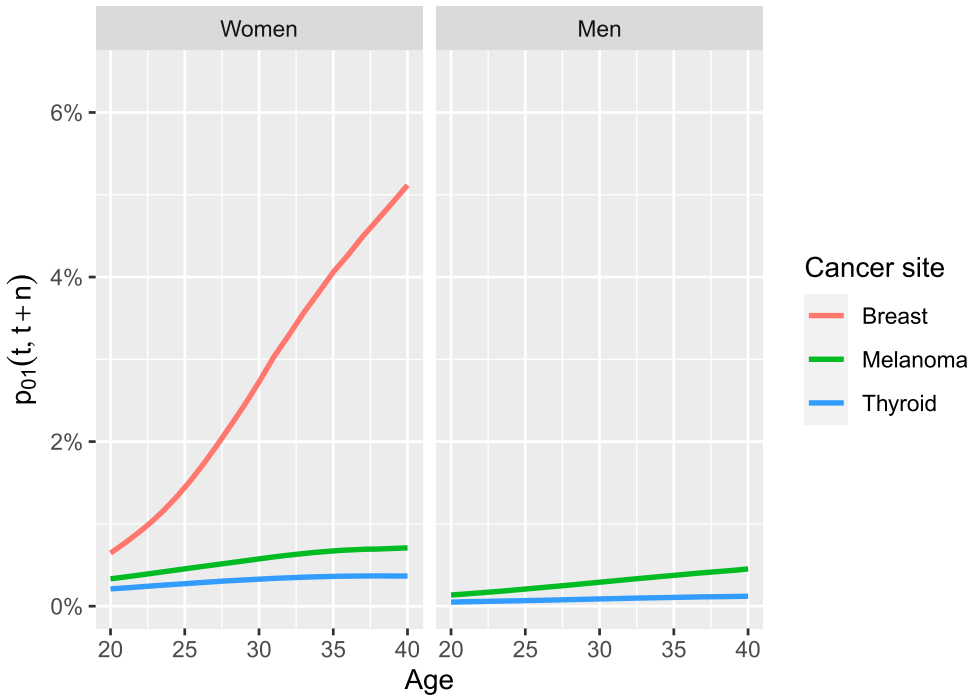
with  $\sum_{l=0}^{k-1} \alpha_{0\bullet}(t + l) = 0$  if  $l = 0$ .

The probabilities of being diagnosed with breast, melanoma, and thyroid cancer over the next 20 years for a healthy individual obtained via the Semi-Markov three-state model are displayed in Figure 3, for ages  $t \in \{20, 21, \dots, 40\}$  and for each sex separately. Figure 3 shows that incidence risk over a 20-year period remains rather low ( $< 0.71\%$ ) for melanoma and thyroid cancers for both sexes, but considerably increases with age for female breast cancer (culminating at 5.12% at age 40).

#### 4.2. Years of life lost from diagnosis

Results of  $YLL_{MSM}^i(T_{01}; z)$  as functions of age at diagnosis ( $T_{01} \in \{20, 21, \dots, 69\}$ ) and for  $z = 0, 5$  and 10 years after diagnosis are presented by sex and cancer site in Figure 4.

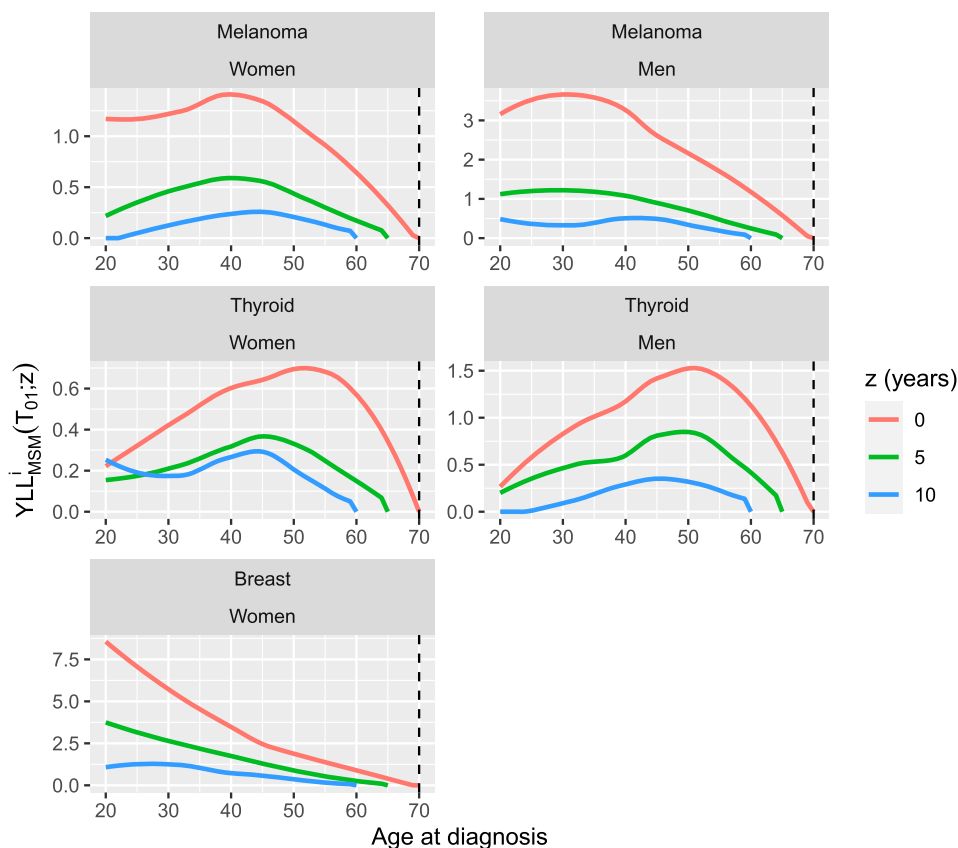
We can see that, for both sexes and all three cancers of interest, the longer the time survived after diagnosis (i.e. the greater the  $z$ ), the lower  $YLL_{MSM}^i(T_{01}; z)$  (with an exception for women diagnosed with thyroid cancer at the age of 25 and below). For female breast cancer,  $YLL_{MSM}^i(T_{01}; z)$  is highest when diagnosed at the age of 20 and then decreases with age at diagnosis, whereas for melanoma and thyroid cancers, it peaks when diagnosed at later ages (between 35 and 55 years depending on the cancer and sex). For both melanoma and thyroid cancers,  $YLL_{MSM}^i(T_{01}; z)$  is larger for men than for women. Botta et al. [40], who describe the impact of cancer during patients' entire lives, found a similar pattern between women and men. Comparisons between sexes cannot be made for breast cancer as only female breast cancer is included. Among men,  $YLL_{MSM}^i(T_{01}; z)$  is globally lower for thyroid cancer than for melanoma cancer. Among women,  $YLL_{MSM}^i(T_{01}; z)$  is lowest for thyroid cancer and highest for breast cancer. Note also that, for patients diagnosed with melanoma or thyroid cancer at all considered ages,  $YLL_{MSM}^i(T_{01}; 10)$  remains below one year. This indicates that, once they have survived their cancer for 10 years, they lose (compared to the general population and up to the age of 70 years) a limited number of years of life due to cancer.



**Figure 3.** Probabilities of being diagnosed with breast, melanoma and thyroid cancer over the next  $n = 20$  years for a healthy individual as a function of age  $t \in \{20, 21, \dots, 40\}$ .

Remember that  $YLL_{MSM}^i(T_{01}; z)$  is computed at the individual level with  $\tau = 70$  years, so these figures give the number of years of life a patient diagnosed with cancer is expected to lose due to the disease before the age of 70 years (at the time of diagnosis, 5 and 10 years after diagnosis). This health indicator can, however, also be analyzed in relative terms, that is, in comparison with other cancers, diseases, or conditions rather than in absolute terms. Indeed, knowing that a group of patients has more to lose (up to a certain age) in terms of years of life due to a specific disease compared to another one is more meaningful for policy-makers and clinicians. This comparison would allow, for example, to rank diseases in terms of burden to society, that is, highlight those which are, until a chosen age, the most (or least) lethal.

It is also worth noting that curves displayed in Figure 4 would be different if another age was chosen for  $\tau$ . Indeed, the higher the upper age limit  $\tau$ , the more years of life an individual can lose. The decreasing trend of  $YLL_{MSM}^i(T_{01}; z)$  at older ages can be explained partly by the fact that the survival of a cancer patient is approaching that of the general population, and partly by the fact that a cancer patient has, comparatively to the general population, simply fewer years of life to lose before the age of 70 years as he or she approaches that age. Unfortunately, it is not possible to distinguish between these two reasons without explicitly calculating the number of years of life lost with a much higher age limit. However, this higher limit was not applied in our calculations, as we are determining the number of years of life lost within the context of mortgage insurance (which we recall, is an insurance product mainly taken out by young adults).



**Figure 4.** Number of life years lost at the individual level before the age of 70 years due to cancer, estimated from  $z = 0, 5$  and 10 years after diagnosis, as a function of age at diagnosis.

## 5. Discussion

As it has been highlighted on several occasions in the literature, there are several approaches and methods to estimate the number of years of life lost due to cancer [13]. Sometimes, it even has different definitions and meanings depending on the context and the audience [7]. It is therefore hard to compare  $YLL^i$  due to cancer across different studies, in particular when the upper age limit  $\tau$  is different. In the present study, it is set to 70 years to focus on young adults and active life, while most studies set it at a higher age to consider the number of years of life lost during the entire lifetime [41,42]. As mentioned above, the number of years of life lost before a given time horizon (70 years in our illustration) obviously depends on how far this time horizon is. Therefore, it is important to note that results found for the number of years of life lost from diagnosis until age 70 should not be taken as an evaluation of the risk from a medical or biological point of view. Such information could however still be very useful in a situation where this time horizon would be meaningful, as could for example be the case from an actuarial or economical point of view. Indeed, in the context of the right to be forgotten for instance, the insurer is mainly interested in the survival until the end of the loan contracted. More generally, from a public

policy perspective, one may be interested in the number of years of life lost before the age of retirement.

Although it is hard to compare results with existing literature, our results could be considered as in line with Silversmit et al. [43], who, also using Belgian data, found a  $YLL^i$  of 3.2 years for female breast cancer, 2.5 and 3.6 years for female and male melanoma cancer, and 1.5 and 2.5 years for female and male thyroid cancer, respectively. These results are obtained using the life expectancy from the general population at the age of diagnosis as the reference age, which is mostly larger than 78 years.

There has been a proliferation of research on the topic of the number of years of life lost, in several countries and for several conditions or diseases. Findings from other studies cannot be compared to each other easily, nor to our results due to the fact that time horizons are different. Nonetheless, for the sake of completeness and for the interested reader, we highlight the main findings related to cancer research. Andersen et al. [2] found a  $YLL^i$  due to cancer (all types of cancer) ranging from 0.11 to 3.68 years for Danish males and from 0.21 to 1.62 years for Russian males, depending on the age at diagnosis and the method used. Baade et al. [11] found a  $YLL^i$  due to, respectively, melanoma and female breast cancer ranging from 3 years (at 40 years old) to 1 year (at 80 years old) and from 12.1 years (at 40 years old) to 1.6 years (at 80 years old). Capocaccia et al. [44] obtained a  $YLL^i$  due to female breast cancer ranging from 8.7 years at age 40–44 to 2.4 years at ages 70–74. For patients diagnosed at age 45 years, Botta et al. [40] found a  $YLL^i$  below 6 years for thyroid cancer in women and melanoma in men. Andersson et al. [6] arrived at a  $YLL^i$  for female breast cancer ranging from 13 years (50–59 age group) to 2.2 years (80+ age group) and from 9.13 years (50–59 age group) to 1.84 years (80+ age group) for melanoma cancer. Finally, Belot et al. [7] found a  $YLL^i$  due to colon cancer over a 10-year time window ranging from 4.14 to 4.77 years depending on the socioeconomic group.

There is a vast literature on YLL and MSM in biostatistical and medical studies. The present paper illustrates their relevance for computing a measure of the number of years of life lost before a given age, chosen depending on the situation or the research question. Arik et al. [45] have shown the implementation of years of life lost in the context of a multi-state model. However, it differs from the present study on several points: (i) it uses a Markov model (so transition intensities do not depend on the duration of stay in the current state), (ii) it is targeted to another age group as it uses data on women diagnosed with breast cancer aged 65–89 years, and (iii) it focuses on the number of years of life lost by the entire cohort. The present paper aims at filling this gap. Some useful applications of MSM-based calculations to derive health indices such as disease incidence risk and the number of years of life lost due to cancer targeted to this public have been illustrated.

Most studies refer to the number of years of life lost or remaining life expectancy starting from the date of diagnosis as an estimate of the disease burden [6,11,12,46–48]. This is undoubtedly useful when considering patients who have just been diagnosed, the time at which a patient is most likely to be concerned about his/her survival. Nonetheless, its relevance need not be limited to quantifying the loss of survival at the time of diagnosis. For long-term survivors, it becomes even more pertinent when considering its evolution over time [40,44]. Indeed, there are many applications where one would be interested in the loss of survival due to cancer, given that the patient already survived some years after diagnosis. This is particularly useful for cancers where the amount of time survived since diagnosis has an influence on the patient's survival. This is actually the underlying basis



behind the right to be forgotten [28,29,49,50]. Implemented since 2016 in France and since 2019 in Belgium, it states that no difference can be made, in terms of access to an insurance product and the level of its premiums, between a healthy client and a cancer patient if he/she survived at most 10 years after the end of the therapeutic protocol.  $YLL^i$  over time since diagnosis can be interpreted as a measure of how close to being cured long-term survivors can be considered [40]. A decreasing  $YLL^i$  over time since diagnosis shows some evidence that patients who are still alive are approaching the same mortality risks as the general population. In this context, Capocaccia et al. [44] proposed a cut-off of less than two years of life lost for colon cancer patients to be considered as statistically cured.

It is important to note that there have been improvements in the treatment of advanced melanoma over the last decade, leading to a positive impact on quality of life and overall patient survival [51–54]. Obviously, the bigger the improvements in treatment and overall survival, the more the duration in the ill state is underestimated and the more the number of years of life lost is overestimated. This does not, nonetheless, undermine our analyses for multiple reasons. First, a better prognosis has no impact on the incidence nor on the incidence risk (i.e. the first application of the present study). Second, the largest improvements in treatment and overall survival concern advanced melanoma, so stages III and IV. These two advanced stages represent a limited share of all tumours considered here (8.96% and 4.16% for stages III and IV, respectively). Third, improvements in treatment are quite recent, limiting the impact on the obtained results. Fourth, in the context of the right to be forgotten and from an insurer's point of view, it is more conservative if the number of years of life lost due to cancer before a certain age is overestimated than if it is underestimated.

Melanoma, thyroid and female breast cancers may include a variety of cancer sub-types and could be diagnosed at different stages of severity, leading to differences in terms of survival. It is thus undeniable that including the information on stages of severity would refine the analysis. This could be achieved, for instance, by stratifying the analyses by cancer stage. However, it has been omitted on purpose for the sake of illustration of the proposed approach.

Cancer is not one disease but a family of many diverse diseases with different outcomes. Results in the present paper focus on melanoma, thyroid, and female breast cancer patients, and cannot, at this stage, be transferred to other cancer types. A natural extension of this work would be to repeat the analyses for all major cancer types. Arik et al. [55] even showed, in a comprehensive study using UK data, that for female breast cancer there are regional differences in terms of cancer morbidity. Thus, the analysis could also be refined to a regional level instead of a national level. This is not done in the present paper as it goes beyond the scope of this study which primarily aims to advocate a new method to estimate the number of years of life lost.

Cancer patient survival has improved over the last few decades, with an increasing proportion of patients being cured for many types of cancer [56–58]. Given the increasing numbers of people being diagnosed with cancer, informing patients and involved parties with relevant risk information is crucial [11]. Providing a precise and informative estimate of the reduction in the remaining life expectancy in case cancer is diagnosed or to long-term cancer survivors is therefore of prime importance, for patients, policy-makers, and society as a whole. From the literature, it is clear that the number of years of life lost is an important addition to existing measures that give a complete picture of the impact of a cancer diagnosis. The methods proposed in this paper help to estimate this important health

indicator from a multi-state model's perspective. This will undoubtedly help to assess when the excess mortality from cancer becomes negligible in cancer survivors, in turn allowing the right to be forgotten to be developed further.

In this study, the assumption is made that a cancer patient cannot become healthy again (i.e. transition from the ill to the healthy state is not possible). Although this assumption is believed to be reasonable for most cancers, one may argue that it does not always hold. However, in our context, the real transition of interest is more from ill to dead than from ill to healthy, following the reasonable paradigm that staying long enough in the ill state to die from something else is, at least from a statistical point of view, equivalent to being cured (cfr. the idea of 'statistical cure' for example in Boussari et al. [59], Jakobsen et al. [60], Tralongo et al. [61]). Also, the main objective of this study is to illustrate how the concept of MSM can be applied to estimate another well-known quantity in medicine and epidemiology, which has not yet been done so far. Using more advanced MSM to estimate the number of years of life lost is undoubtedly an interesting question, but left for future research.

For cancer patients, quality of life may be considered as important as the length of life itself [62]. The number of years of life lost gives an easily interpretable measure about the survival of cancer patients. However, other indicators such as, among others, the disability-adjusted life years (DALY) should also be considered, in particular for diseases or conditions that cause significant disability or do not result in death. Another metric which could be estimated via the proposed methods is the cancer-free life expectancy. This could be estimated (i) via the mean sojourn time in the healthy state, or (ii) by subtracting, from the life expectancy in the general population, the number of years of life lost due to cancer up to the age corresponding to the life expectancy in the general population. Note that even though it is the number of years of life lost due to cancer that is estimated, the methods proposed in this paper are not limited to cancer and could be applied to several other diseases or conditions (diabetes and HIV, amongst others).

## Acknowledgments

The three first authors gratefully acknowledge the Belgian Cancer Registry for providing access to the data and for research assistance. The authors additionally thank the staff of the Belgian Cancer Registry and all physicians, pathologists, and data managers involved in Cancer Registration in Belgium for their dedicated data collection. We used the STROBE cohort checklist when writing our report [63]. The authors also thank the two anonymous referees for comments that have been very helpful for revising previous versions of the present work.

## Data availability statement

The datasets generated and/or analysed during the current study are not publicly available due to privacy reasons but are available from the corresponding author on reasonable request. The pseudonymized data can be provided within the secured environment of the Belgian Cancer Registry according to its regulations, and only upon approval by the Information Security Committee.

## Authors' contributions

All authors contributed to the study conception, design, methodology, data analysis, and interpretation. The first draft of the manuscript was written by AS, CL, and MD, and all authors commented

on previous versions of the manuscript. All authors revised the draft and approved the final version of the manuscript.

## Disclosure statement

No potential conflict of interest was reported by the author(s).

## Funding

This research project was supported by UCLouvain and funding from the FWO and F.R.S.-FNRS under the Excellence of Science (EOS) program, project EOS 40007517.

## Ethics approval and consent to participate

Ethical review and approval was not required for the study involving humans in accordance with the local legislation and institutional requirements. Written informed consent to participate in this study was not required from the participants or the participants' legal guardians/next of kin in accordance with the national legislation and the institutional requirements.

## Consent for publication

Not applicable.

## ORCID

Antoine Soetewy  <http://orcid.org/0000-0001-8159-0804>

## References

- [1] Andersen PK. Decomposition of number of life years lost according to causes of death. *Stat Med*. 2013;32(30):5278–5285. doi: [10.1002/sim.v32.30](https://doi.org/10.1002/sim.v32.30)
- [2] Andersen PK, Canudas-Romo V, Keiding N. Cause-specific measures of life years lost. *Demogr Res*. 2013;29:1127–1152. doi: [10.4054/DemRes.2013.29.41](https://doi.org/10.4054/DemRes.2013.29.41)
- [3] Latouche A, Andersen PK, Rey G, et al. A note on the measurement of socioeconomic inequalities in life years lost by cause of death. *Epidemiology*. 2019;30(4):569–572. doi: [10.1097/EDE.0000000000001022](https://doi.org/10.1097/EDE.0000000000001022)
- [4] Percy C, Stanek 3rd E, Gloeckler L. Accuracy of cancer death certificates and its effect on cancer mortality statistics. *Am J Public Health*. 1981;71(3):242–250. doi: [10.2105/AJPH.71.3.242](https://doi.org/10.2105/AJPH.71.3.242)
- [5] Aragon TJ, Lichtensztajn DY, Katcher BS, et al. Calculating expected years of life lost for assessing local ethnic disparities in causes of premature death. *BMC Public Health*. 2008;8(1):116. doi: [10.1186/1471-2458-8-116](https://doi.org/10.1186/1471-2458-8-116)
- [6] Andersson TM-L, Dickman PW, Eloranta S, et al. Estimating the loss in expectation of life due to cancer using flexible parametric survival models. *Stat Med*. 2013;32(30):5286–5300. doi: [10.1002/sim.v32.30](https://doi.org/10.1002/sim.v32.30)
- [7] Belot A, Ndiaye A, Luque-Fernandez M-A, et al. Summarizing and communicating on survival data according to the audience: a tutorial on different measures illustrated with population-based cancer registry data. *Clin Epidemiol*. 2019;11:53–65. doi: [10.2147/CLEP](https://doi.org/10.2147/CLEP)
- [8] Chu P-C, Wang J-D, Hwang J-S, et al. Estimation of life expectancy and the expected years of life lost in patients with major cancers: extrapolation of survival curves under high-censored rates. *Value Health*. 2008;11(7):1102–1109. doi: [10.1111/j.1524-4733.2008.00350.x](https://doi.org/10.1111/j.1524-4733.2008.00350.x)
- [9] Manevski D, Ružić Gorenjec N, Andersen PK, et al. Expected life years compared to the general population. *Biom J*. 2023;65(4):2200070. doi: [10.1002/bimj.v65.4](https://doi.org/10.1002/bimj.v65.4)

- [10] Antero-Jacquemin J, Pohar-Perme M, Rey G, et al. The heart of the matter: years-saved from cardiovascular and cancer deaths in an elite athlete cohort with over a century of follow-up. *Eur J Epidemiol.* **2018**;33:531–543. doi: [10.1007/s10654-018-0401-0](https://doi.org/10.1007/s10654-018-0401-0)
- [11] Baade PD, Youlten DR, Andersson TM, et al. Estimating the change in life expectancy after a diagnosis of cancer among the Australian population. *BMJ Open.* **2015**;5(4):e006740. doi: [10.1136/bmjopen-2014-006740](https://doi.org/10.1136/bmjopen-2014-006740)
- [12] Licher S, Heshmatollah A, van der Willik KD, et al. Lifetime risk and multimorbidity of non-communicable diseases and disease-free life expectancy in the general population: a population-based cohort study. *PLoS Med.* **2019**;16(2):e1002741. doi: [10.1371/journal.pmed.1002741](https://doi.org/10.1371/journal.pmed.1002741)
- [13] Andersen PK. Life years lost among patients with a given disease. *Stat Med.* **2017**;36(22):3573–3582. doi: [10.1002/sim.7357](https://doi.org/10.1002/sim.7357)
- [14] Andersen PK, Borgan O, Gill RD, et al. Statistical models based on counting processes. New York, NY: Springer Science & Business Media; **2012**.
- [15] Hougaard P. Multi-state models: a review. *Lifetime Data Anal.* **1999**;5:239–264. doi: [10.1023/A:1009672031531](https://doi.org/10.1023/A:1009672031531)
- [16] De Wreede LC, Fiocco M, Putter H. The mstate package for estimation and prediction in non- and semi-parametric multi-state and competing risks models. *Comput Methods Programs Biomed.* **2010**;99(3):261–274. doi: [10.1016/j.cmpb.2010.01.001](https://doi.org/10.1016/j.cmpb.2010.01.001)
- [17] Geskus RB. Data analysis with competing risks and intermediate states. Boca Raton, FL: Chapman and Hall/CRC; **2019**.
- [18] Putter H, Fiocco M, Geskus RB. Tutorial in biostatistics: competing risks and multi-state models. *Stat Med.* **2007**;26(11):2389–2430. doi: [10.1002/sim.v26:11](https://doi.org/10.1002/sim.v26:11)
- [19] Meira-Machado L, de Uña-Álvarez J, Cadarso-Suárez C, et al. Multi-state models for the analysis of time-to-event data. *Stat Methods Med Res.* **2009**;18(2):195–222. doi: [10.1177/0962280208092301](https://doi.org/10.1177/0962280208092301)
- [20] Andersen PK, Pohar Perme M. Inference for outcome probabilities in multi-state models. *Lifetime Data Anal.* **2008**;14(4):405–431. doi: [10.1007/s10985-008-9097-x](https://doi.org/10.1007/s10985-008-9097-x)
- [21] Touraine C, Helmer C, Joly P. Predictions in an illness-death model. *Stat Methods Med Res.* **2016**;25(4):1452–1470. doi: [10.1177/0962280213489234](https://doi.org/10.1177/0962280213489234)
- [22] Denuit M, Lucas N, Pitacco E. Pricing and reserving in LTC insurance. In: Etienne Dupourqué, Frédéric Planchet, Néfissa Sator editors. Actuarial aspects of long term care. Cham, Switzerland: Springer International Publishing, 2019, p. 129–158.
- [23] Anderson RM. Discussion: the kermack-mckendrick epidemic threshold theorem. *Bull Math Biol.* **1991**;53(1):1–32. doi: [10.1007/BF02464422](https://doi.org/10.1007/BF02464422)
- [24] Kermack WO, McKendrick AG. A contribution to the mathematical theory of epidemics. *Proc R Soc Lond. Ser A, Contain Pap Math Phys Character.* **1927**;115(772):700–721.
- [25] CRUK. Melanoma skin cancer survival statistics. Available from: <https://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/melanoma-skin-cancer>, 2023.
- [26] CRUK. Thyroid cancer survival statistics. Available from: <https://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/thyroid-cancer>, 2023.
- [27] NHS Digital. Cancer survival in England, cancers diagnosed 2016 to 2020, followed up to 2021. Available from: <https://digital.nhs.uk/data-and-information/publications/statistical/cancer-survival-in-england/cancers-diagnosed-2016-to-2020-followed-up-to-2021>, 2023.
- [28] Soetewey A, Legrand C, Denuit M, et al. Waiting period from diagnosis for mortgage insurance issued to cancer survivors. *Eur Actuar J.* **2021**;11(1):135–160. doi: [10.1007/s13385-020-00254-x](https://doi.org/10.1007/s13385-020-00254-x)
- [29] Soetewey A, Legrand C, Denuit M, et al. Right to be forgotten for mortgage insurance issued to cancer survivors: critical assessment and new proposal. *Eur Actuar J.* **2024**;15:1–29.
- [30] CRUK. Breast cancer survival statistics. Available from: <https://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/breast-cancer>, 2023.
- [31] Esteve J, Benhamou E, Raymond L, et al. Statistical methods in cancer research. Volume IV. Descriptive epidemiology. *IARC Sci Publ.* **1994**;128(1):302.

- [32] Oksanen H. Modelling the survival of prostate cancer patients. Tampere, Finland: University of Tampere; 1998.
- [33] Soetewey A, Legrand C, Denuit M, et al. Semi-Markov modeling for cancer insurance. *Eur Actuar J*. 2022;12(2):813–837. doi: [10.1007/s13385-022-00308-2](https://doi.org/10.1007/s13385-022-00308-2)
- [34] Chiang CL. Life table and its applications. Malabar, Florida: R.E.Krieger Pub. Co.; 1984, pp. 316.
- [35] Keyfitz N, Caswell H. Applied mathematical demography. Vol. New York, NY: Springer; 2005.
- [36] Preston S, Heuveline P, Guillot M, et al. Measuring and modeling population processes. Malden, MA: Wiley-Blackwell; 2001.
- [37] Kaplan EL, Meier P. Nonparametric estimation from incomplete observations. *J Am Stat Assoc*. 1958;53(282):457–481. doi: [10.1080/01621459.1958.10501452](https://doi.org/10.1080/01621459.1958.10501452)
- [38] Jackson C. Multi-state modelling with R: the msm package. Cambridge, UK: MRC Biostatistics Unit, University of Cambridge. 2007, p. 1–53.
- [39] Noordzij M, Dekker FW, Zoccali C, et al. Measures of disease frequency: prevalence and incidence. *Nephron Clin Pract*. 2010;115(1):c17–c20. doi: [10.1159/000286345](https://doi.org/10.1159/000286345)
- [40] Botta L, Dal Maso L, Guzzinati S, et al. Changes in life expectancy for cancer patients over time since diagnosis. *J Adv Res*. 2019;20:153–159. doi: [10.1016/j.jare.2019.07.002](https://doi.org/10.1016/j.jare.2019.07.002)
- [41] Centers for Disease Control P. Years of potential life lost before age 65–united states, 1990 and 1991. *MMWR Morb Mortal Wkly Rep*. 1993;42(13):251–253.
- [42] Gardner JW, Sanborn JS. Years of potential life lost (YPLL)—what does it measure? *Epidemiology*. 1990;1(4):322–329. doi: [10.1097/00001648-199007000-00012](https://doi.org/10.1097/00001648-199007000-00012)
- [43] Silversmit G, Vaes E, van Eycken L. Estimation of population-based cancer-specific potential years of life lost in belgium. *Eur J Cancer Prev*. 2017;26:157–163. doi: [10.1097/CEJ.00000000000000385](https://doi.org/10.1097/CEJ.00000000000000385)
- [44] Capocaccia R, Gatta G, Dal Maso L. Life expectancy of colon, breast, and testicular cancer patients: an analysis of us-seer population-based data. *Ann Oncol*. 2015;26(6):1263–1268. doi: [10.1093/annonc/mdv131](https://doi.org/10.1093/annonc/mdv131)
- [45] Arik A, Cairns AJG, Dodd E, et al. Estimating the impact of the covid-19 pandemic on breast cancer deaths among older women. In: 2023 Living to 100 Research Symposium-Asia; Kowloon, Hong Kong, on February 16, 2023, organized by the Society of Actuaries; 2023.
- [46] Andersson TM, Dickman PW, Eloranta S, et al. The loss in expectation of life after colon cancer: a population-based study. *BMC Cancer*. 2015;15(1):1–10. doi: [10.1186/s12885-015-1427-2](https://doi.org/10.1186/s12885-015-1427-2)
- [47] Baade PD, Youlten DR, Andersson TM, et al. Temporal changes in loss of life expectancy due to cancer in australia: a flexible parametric approach. *Cancer Causes Control*. 2016;27(8):955–964. doi: [10.1007/s10552-016-0762-1](https://doi.org/10.1007/s10552-016-0762-1)
- [48] Syriopoulou E, Bower H, Andersson TM, et al. Estimating the impact of a cancer diagnosis on life expectancy by socio-economic group for a range of cancer types in England. *Br J Cancer*. 2017;117(9):1419–1426. doi: [10.1038/bjc.2017.300](https://doi.org/10.1038/bjc.2017.300)
- [49] Mesnil M. What do we mean by the right to be forgotten? an analysis of the french case study from a lawyer’s perspective. *J Cancer Policy*. 2018;15:122–127. doi: [10.1016/j.jcpo.2018.01.001](https://doi.org/10.1016/j.jcpo.2018.01.001)
- [50] Scocca G, Meunier F. A right to be forgotten for cancer survivors: a legal development expected to reflect the medical progress in the fight against cancer. *J Cancer Policy*. 2020;25:100246. doi: [10.1016/j.jcpo.2020.100246](https://doi.org/10.1016/j.jcpo.2020.100246)
- [51] Pasquali S, Hadjinicolaou AV, Sileni VC, et al. Systemic treatments for metastatic cutaneous melanoma. *Cochrane Database Syst Rev*. 2018(2).
- [52] Pedersen S, Holmstroem RB, von Heymann A, et al. Quality of life and mental health in real-world patients with resected stage III/IV melanoma receiving adjuvant immunotherapy. *Acta Oncol (Madr)*. 2023;62(1):62–69. doi: [10.1080/0284186X.2023.2165449](https://doi.org/10.1080/0284186X.2023.2165449)
- [53] Tichanek F, Försti A, Hemminki A, et al. Survival in melanoma in the nordic countries into the era of targeted and immunological therapies. *Eur J Cancer*. 2023;186:133–141. doi: [10.1016/j.ejca.2023.03.019](https://doi.org/10.1016/j.ejca.2023.03.019)
- [54] Tromme I, Legrand C, Devleeschauwer B, et al. Melanoma burden by melanoma stage: assessment through a disease transition model. *Eur J Cancer*. 2016;53:33–41. doi: [10.1016/j.ejca.2015.09.016](https://doi.org/10.1016/j.ejca.2015.09.016)

- [55] Arık A, Dodd E, Streftaris G. Cancer morbidity trends and regional differences in england—a Bayesian analysis. *PLoS ONE*. 2020;15(5):e0232844. doi: [10.1371/journal.pone.0232844](https://doi.org/10.1371/journal.pone.0232844)
- [56] Andersson TM, Dickman PW, Eloranta S, et al. Estimating and modelling cure in population-based cancer studies within the framework of flexible parametric survival models. *BMC Med Res Methodol*. 2011;11(1):96. doi: [10.1186/1471-2288-11-96](https://doi.org/10.1186/1471-2288-11-96)
- [57] Lambert PC, Thompson JR, Weston CL, et al. Estimating and modeling the cure fraction in population-based cancer survival analysis. *Biostatistics*. 2006;8(3):576–594. doi: [10.1093/biostatistics/kxl030](https://doi.org/10.1093/biostatistics/kxl030)
- [58] Silversmit G, Jegou D, Vaes E, et al. Cure of cancer for seven cancer sites in the Flemish region. *Int J Cancer*. 2017;140(5):1102–1110. doi: [10.1002/ijc.v140.5](https://doi.org/10.1002/ijc.v140.5)
- [59] Boussari O, Romain G, Remontet L, et al. A new approach to estimate time-to-cure from cancer registries data. *Cancer Epidemiol*. 2018;53:72–80. doi: [10.1016/j.canep.2018.01.013](https://doi.org/10.1016/j.canep.2018.01.013)
- [60] Jakobsen LH, Andersson TM-L, Bicler JL, et al. On estimating the time to statistical cure. *BMC Med Res Methodol*. 2020;20:1–13. doi: [10.1186/s12874-020-00946-8](https://doi.org/10.1186/s12874-020-00946-8)
- [61] Tralongo P, McCabe MS, Surbone A. Challenge for cancer survivorship: improving care through categorization by risk. *J Clin Oncol*. 2017;35(30):3516–7. doi: [10.1200/JCO.2017.74.3450](https://doi.org/10.1200/JCO.2017.74.3450)
- [62] Shrestha A, Martin C, Burton M, et al. Quality of life versus length of life considerations in cancer patients: a systematic literature review. *Psychooncology*. 2019;28(7):1367–1380. doi: [10.1002/pon.v28.7](https://doi.org/10.1002/pon.v28.7)
- [63] Von Elm E, Altman DG, Egger M, et al. The strengthening the reporting of observational studies in epidemiology (strobe) statement: guidelines for reporting observational studies. *The Lancet*. 2007;370(9596):1453–1457. doi: [10.1016/S0140-6736\(07\)61602-X](https://doi.org/10.1016/S0140-6736(07)61602-X)

## Appendix. Development of $e_{11}^{\tau}(t; z)$

In this section, we show how Equation (8) is obtained. Assuming  $\tau > t$ , and  $\tau$  and  $t$  are integers, we have

$$e_{11}^{\tau}(t; z) = \int_t^{\tau} p_{11}(t, u; z) du \quad (\text{A1})$$

$$= \sum_{k=0}^{\tau-t-1} \int_{t+k}^{t+k+1} p_{11}(t, u; z) du \quad (\text{A2})$$

$$= \sum_{k=0}^{\tau-t-1} \int_{t+k}^{t+k+1} p_{11}(t, t+k; z) p_{11}(t+k, u; z+k) du \quad (\text{A3})$$

$$= \sum_{k=0}^{\tau-t-1} \underbrace{p_{11}(t, t+k; z)}_{(1)} \underbrace{\int_{t+k}^{t+k+1} p_{11}(t+k, u; z+k) du}_{(2)} \quad (\text{A4})$$

The terms (1) and (2) in Equation (A4) are developed below.

$$(1) p_{11}(t, t+k; z) = \exp\left(-\int_t^{t+k} \alpha_{12}(u; z+u-t) du\right) \quad (\text{A5})$$

$$= \exp\left(-\sum_{l=0}^{k-1} \int_{t+l}^{t+l+1} \alpha_{12}(u; z+u-t) du\right) \quad (\text{A6})$$

$$= \exp\left(-\sum_{l=0}^{k-1} \int_{t+l}^{t+l+1} \alpha_{12}(t+l; z+l) du\right) \quad (\text{A7})$$

$$= \exp\left(-\sum_{l=0}^{k-1} \alpha_{12}(t+l; z+l)\right) \quad (\text{A8})$$

$$(2) \int_{t+k}^{t+k+1} p_{11}(t+k, u; z+k) du = \int_{t+k}^{t+k+1} \exp\left(-\int_{t+k}^u \alpha_{12}(s; z+s-t) ds\right) du \quad (\text{A9})$$

$$= \int_{t+k}^{t+k+1} \exp\left(-\int_{t+k}^u \alpha_{12}(t+k; z+k) ds\right) du \quad (\text{A10})$$

$$= \int_{t+k}^{t+k+1} \exp(-\alpha_{12}(t+k; z+k)(u-t-k)) du \quad (\text{A11})$$

$$= \left[ \frac{\exp\left(-\alpha_{12}(t+k; z+k)(u-t-k)\right)}{-\alpha_{12}(t+k; z+k)} \right]_{t+k}^{t+k+1} \quad (\text{A12})$$

$$= \frac{1 - \exp\left(-\alpha_{12}(t+k; z+k)\right)}{\alpha_{12}(t+k; z+k)} \quad (\text{A13})$$

Hence,

$$e_{11}^{\tau}(t; z) = \sum_{k=0}^{\tau-t-1} \underbrace{\exp\left(-\sum_{l=0}^{k-1} \alpha_{12}(t+l; z+l)\right)}_{(1)} \underbrace{\frac{1 - \exp\left(-\alpha_{12}(t+k; z+k)\right)}{\alpha_{12}(t+k; z+k)}}_{(2)}. \quad (\text{A14})$$