



**HAL**  
open science

# Continuous time semi-Markov inference of biometric laws associated with a Long-Term Care Insurance portfolio

Guillaume Biessy

► **To cite this version:**

Guillaume Biessy. Continuous time semi-Markov inference of biometric laws associated with a Long-Term Care Insurance portfolio. 2015. hal-01220564v1

**HAL Id: hal-01220564**

**<https://hal.science/hal-01220564v1>**

Preprint submitted on 27 Oct 2015 (v1), last revised 19 May 2016 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Continuous time semi-Markov inference of biometric laws associated with a Long-Term Care Insurance portfolio

Guillaume BIESSY<sup>1</sup>

SCOR Global Life<sup>2</sup> - LaMME<sup>3</sup>

## Abstract

Unlike the mortality risk on which actuaries have been working for more than a century, the long-term care risk is young and as of today hardly mastered. Semi-Markov processes have been identified as an adequate tool to study this risk. Nevertheless, access to data is limited and the associated literature still scarce. Insurers mainly use discrete time methods directly inspired from the study of mortality in order to build experience tables. Those methods however are not perfectly suited for the study of competing risk situations.

The present paper aims at providing a theoretical framework to estimate biometric laws associated with a long-term care insurance portfolio. The presented method relies on a continuous-time semi-Markov model with three states: autonomy, dependency and death. The dependency process is defined using its transition intensities. We provide a formula to infer the mortality of autonomous people from the general population mortality, on which we ought to have more reliable knowledge. We then propose a parametric expression for the remaining intensities of the model. Incidence in dependency is described by a logistic formula. Under the assumption that the dependent population is a mixture of two populations with respect to the category of pathology that caused dependency, we show that the resulting intensity of mortality for dependent people takes a very peculiar form, which is semi-Markov. Estimation of parameters relies on the maximum likelihood method. A parametric approach eliminates issues related to segmentation in age categories, smoothing or extrapolation at higher ages. While creating model uncertainty, it proves very convenient for the practitioner. Finally, we provide an application using data from a real long-term care insurance portfolio.

**Keywords:** Long-Term Care Insurance, continuous time semi-Markov process, competing risks, maximum likelihood, mixture model, parametric model.

## 1 Introduction

Dependency among elderly people can be defined as a permanent state of inability to perform activities of daily living on one's own. It is mostly caused by diseases linked to ageing, such as dementia, neurological diseases, cardiovascular diseases and cancer. Dependent elderly people require regular care whose frequency increases with the severity of their status. If some people can rely at least partially on their family or their friends for help, others have to hire professional caregivers or join a nursing home, whose average cost exceeds 3,000 € a month. Despite public aids, this cost proves overwhelming for most pensioners. Therefore, long-term care is associated a financial risk to which most people are exposed. In France, part of this risk is transferred through private insurance contracts.

The long-term care risk is complex. Its study requires to take into account incidence in dependency as well as probabilities of death for both autonomous and dependent people, which are very different from another. This risk is directly related to ageing through pathologies, and longevity gains in the second half of the 20th century made it paramount. The first long-term care insurance products appeared for instance in France at the end of the 1980's. Average age at subscribing for those products is close to 60 when the average age at which dependency occurs is close to 85. Therefore, even on older portfolio, the number of claims remains limited, and so does the data available. Moreover, insurers and public aids use different definitions to assess the level of required care. Insurers may also change either their definition or

---

1. Contact: gbiessy@scor.com.

2. SCOR Global Life SE - 5 avenue Kléber, 75795 Paris Cedex 16, France.

3. Laboratoire de Mathématiques et Modélisation d'Évry - Université d'Évry Val d'Essonne, UMR CNRS 8071, USC INRA, Évry, France.

their underwriting and claims policy over time. All those elements make data aggregation from several sources very difficult, which may explain the difficulty of getting a better knowledge of the risk.

Markov processes are such that their transition probabilities only depend on the current state of the process. Semi-Markov process is a generalization for which transition probabilities depend on both the current state and the time spent in the current state. One can find more details about those processes in Cinlar (1969). Multi-state models based on Markov and semi-Markov processes have led to many application in the field of epidemiology. Due to the similarities between dependency and pathologies, those processes appear as natural candidates to study this phenomenon. One can refer to Haberman and Pitacco (1998) or more recently Christiansen (2012) which handle life insurance models in general and among them long-term care insurance models. Practitioners have also played a key role in the knowledge of the dependency risk. One of the very first model was presented in SCOR (1995). Relying on a parametric approach, it highlights the exponential increase in the probabilities of incidence in dependency, and defines dependent mortality (resp. autonomous mortality) as a linear function of the general population mortality, computed via an exogenous mortality table. With only 5 parameters required to model the whole dependency process, this model is remarkably simple. It is however based on the Markov assumption that dependent mortality only depends on the age of the dependent, and not on the time since the entry in dependency. This assumption is still used today by many insurers as well as in recent academic papers like Pitacco (2015), because it allows for simpler models.

In this article, we present a very different approach relying on a continuous-time semi-Markov process, which is defined using its transition intensities. Compared to a discrete-time approach, it allows to get a simpler definition of the process. It takes into account the nature of the competing risks (dependency and death) and does not require any assumption on the order in which events occur within a year. Section 2 derives an equation to express the autonomous mortality using general mortality and other intensities of the model. Benefits to use general mortality instead of autonomous mortality are discussed with more details. We then introduce the intensity for general mortality of the portfolio using a simple relational model as in Brass (1971). We propose a parametric expression for the intensity of incidence in dependency, based on the logistic form introduced by Perks (1932) for the study of mortality. We use a different expression for the intensity of mortality in dependency, inspired from a mixture model where there would be two homogeneous populations of dependent people, with different mortality intensities. Estimation of various parameters relies on the maximum likelihood method. We also introduce formulas for pricing and reserving based directly on the transition intensities. Section 3 provides an application of the model based on data from a real insurance portfolio. For each transition intensity, several models of increasing complexity are compared using the Bayesian Information Criterion, and robustness of estimation is also assessed using a non-parametric quantile bootstrap method. Finally, Section 4 summarizes the results obtained and discusses limits and potential improvements of the model.

## 2 Model

### 2.1 Notations

For  $x_0 \geq 0$ , let us consider a continuous-time process  $(Z_x)_{x \geq x_0}$  with values in the 3-state set  $E = \{A, I, D\}$  of autonomy, dependency, death. Let us assume that  $Z$  is *càd-làg* and that  $Z_{x_0} = A$ . The index variable of the process  $Z$  is called age of the process. Therefore  $x_0$  is an initial age where all individuals are assumed to be autonomous. For  $x \geq x_0$  let us denote by  $A_x$  (resp.  $I_x, D_x$ ) the probability for the process to be in the state of autonomy (resp. dependency, death) at age  $x$  or more formally

$$\begin{aligned} A_x &= P(Z_x = A), \\ I_x &= P(Z_x = I), \\ D_x &= P(Z_x = D). \end{aligned}$$

Hence  $A_{x_0} = 1$  and for all  $x \geq x_0$ ,  $A_x + I_x + D_x = 1$ .

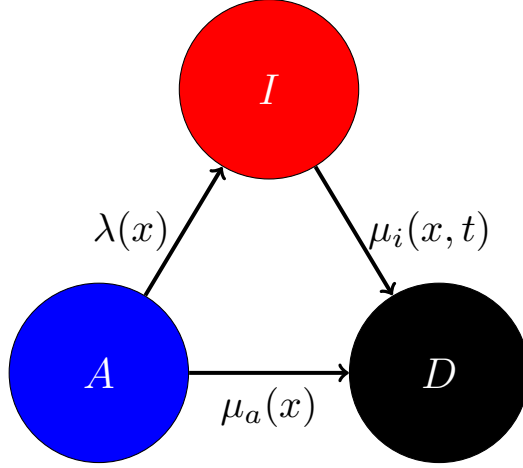


Figure 1: The 3 states continuous-time model and the associated transition intensities

We now introduce the transition intensities, also called instantaneous transition probabilities. We make the assumption that the transition between dependency and death is semi-Markov and depends on both the age of entry in the dependency state and the time spent in this state while other transitions only depend on the current age. Transition intensities allow us to fully describe the behaviour of the process

$$\begin{aligned}\mu_a(x) &= \lim_{h \rightarrow 0} \frac{1}{h} P(Z_{x+h} = D | Z_x = A), \\ \lambda(x) &= \lim_{h \rightarrow 0} \frac{1}{h} P(Z_{x+h} = I | Z_x = A), \\ \mu_i(x, t) &= \lim_{h \rightarrow 0} \frac{1}{h} P(Z_{x+t+h} = D | Z_{x-} = A, Z_x = I, Z_{x+t} = I).\end{aligned}$$

Those intensities are called respectively intensity of entry in dependency, intensity of mortality for autonomous people and intensity of mortality for dependent people. We consider that death is an absorbing state and that there is no transition allowed from dependency to autonomy. A representation of the model can be found on Figure 1. We express below  $A_x$  and  $I_x$  for  $x > x_0$  using transition intensities of the model.

For  $x \geq x_0$ ,  $h \geq 0$ , we have

$$P(Z_{x+h} = A) = [1 - P(Z_{x+h} = I | Z_x = A) - P(Z_{x+h} = D | Z_x = A)] \times P(Z_x = A)$$

and therefore

$$\frac{d}{dx} P(Z_x = A) = -[\mu_a(x) + \lambda(x)] P(Z_x = A).$$

As  $A_{x_0} = 1$ , this equation can be solved in

$$A_x = \exp \left( - \int_{x_0}^x [\lambda(u) + \mu_a(u)] du \right). \quad (1)$$

For  $x \geq x_0$ ,  $t, h \geq 0$ , we can write

$$\begin{aligned}P(Z_{x+t+h} = I | Z_{x-} = A, Z_x = I) &= P(Z_{x+t+h} = I | Z_{x-} = A, Z_x = I, Z_{x+t} = I) \\ &\quad \times P(Z_{x+t} = I | Z_{x-} = A, Z_x = I).\end{aligned}$$

which gives us

$$\frac{d}{dt} P(Z_{x+t} = I | Z_{x-} = A, Z_x = I) = -\mu_i(x, t) P(Z_{x+t} = I | Z_{x-} = A, Z_x = I).$$

As

$$P(Z_x = I | Z_{x-} = A, Z_x = I) = 1$$

we obtain

$$P(Z_{x+t} = I | Z_{x-} = A, Z_x = I) = \exp \left( - \int_0^t \mu_i(x, u) du \right).$$

Moreover, as we have

$$I_x = \int_{x_0}^x P(Z_u = A)P(Z_{u+du} = I|Z_u = A)P(Z_x = I|Z_u = A, Z_{u+du} = I)$$

we get an expression of the probability to be dependent at age  $x \geq x_0$

$$I_x = \int_{x_0}^x \lambda(u)A_u \exp\left(-\int_u^x \mu_i(u, v-u)dv\right) du. \quad (2)$$

## 2.2 Link with general mortality

Let us consider the intensity of mortality for the general population defined by

$$\mu_g(x) = \lim_{h \rightarrow 0} \frac{1}{h} P(Z_{x+h} = D | Z_x \in \{A, I\}).$$

Figure 2 represents the fourth transition in our model, a transition between life and death for the general population.

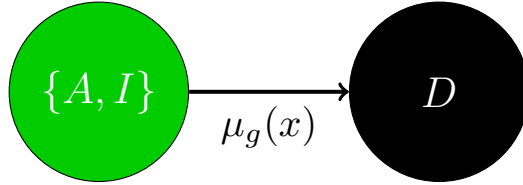


Figure 2: Intensity of transition for the general population

**Lemma.** For  $x \geq x_0$  and  $t \geq 0$ , let us denote by  $\Delta(x, t)$  the difference between the intensity of mortality for dependent people and the intensity of mortality of autonomous people with the same age, so that  $\Delta(x, t) = \mu_i(x, t) - \mu_a(x + t)$ . Then the intensity of mortality for autonomous people is solution of the following equation

$$\mu_a(x) = \mu_g(x) - \frac{\int_{x_0}^x \lambda(u)\Delta(u, x-u) \exp\left(-\int_u^x [\Delta(u, v-u) - \lambda(v)] dv\right) du}{1 + \int_{x_0}^x \lambda(u) \exp\left(-\int_u^x [\Delta(u, v-u) - \lambda(v)] dv\right) du}. \quad (3)$$

*Proof.* Differentiating (1) and (2) gives us equations (4) and (5) below which describe the evolution of the probabilities  $A_x$  and  $I_x$ . Similarly, from the definition of  $\mu_g$  we get equation (6). We obtain a system of 3 differential equations

$$\frac{d}{dx} A_x = -[\lambda(x) + \mu_a(x)]A_x \quad (4)$$

$$\frac{d}{dx} I_x = \lambda(x)A_x - \int_{x_0}^x \lambda(u)A_u \exp\left(-\int_u^x \mu_i(u, v-u)dv\right) \mu_i(u, x-u)du \quad (5)$$

$$\frac{d}{dx} (A_x + I_x) = -\mu_g(x)(A_x + I_x). \quad (6)$$

Summing the evolution equations (4) and (5), then identifying with equation (6) yields

$$\mu_g(x)(A_x + I_x) = \mu_a(x)A_x + \int_{x_0}^x \lambda(u)A_u \exp\left(-\int_u^x \mu_i(u, v-u)dv\right) \mu_i(u, x-u)du.$$

With simple algebra we get

$$\mu_a(x) = \mu_g(x) \left(1 + \frac{I_x}{A_x}\right) - \int_{x_0}^x \lambda(u) \frac{A_u}{A_x} \exp\left(-\int_u^x \mu_i(u, v-u)dv\right) \mu_i(u, x-u)du.$$

Using (2) and (1), we obtain after a few simplifications

$$\mu_a(x) = \mu_g(x) - \int_{x_0}^x \lambda(u) \exp \left( - \int_u^x [\mu_i(u, v - u) - \lambda(v) - \mu_a(v)] dv \right) [\mu_i(u, x - u) - \mu_g(x)] du.$$

Now, we replace the intensity of mortality for dependent people using the formula

$$\mu_i(x, t) = \mu_a(x + t) + \Delta(x, t)$$

which gives us

$$\mu_a(x) = \mu_g(x) - \int_{x_0}^x \lambda(u) \exp \left( - \int_u^x [\Delta(u, v - u) - \lambda(v)] dv \right) [\mu_a(x) - \mu_g(x) + \Delta(u, x - u)] du.$$

This finally leads to the result. □

Equation (3) allows us to use the general mortality instead of the autonomous mortality in the model. On one hand, mortality of autonomous people is complex to predict, because people can leave the autonomy state either by becoming dependent or dying. Furthermore, the scope of autonomous people depends directly on the definition used for dependency. Therefore predicting the autonomous mortality requires to have intensive knowledge of the dependency process beforehand. On the other hand, general mortality has been studied for a long time by actuaries, demographers, biologists and is very well documented. One can therefore rely on reference mortality table and a wide range of models.

The formula does not give an analytic expression for the intensity of autonomous mortality in the most general case. Numeric methods can however be used to compute it. With an *ad hoc* choice of model, the inner integrals in the formula take an analytical expression and numeric approximation is only required for the outer integrals. Intensity of general mortality appears directly in the equation, which is very convenient if we want to use an external reference for this intensity.

### 2.3 Data structure

In practice, data issued from insurance portfolio generally consists in two databases. The first database gathers information on contributors and the second on annuitants. From one portfolio to another, data quality and available information may vary a lot. In what follows, we assume both databases contain the variables of Table 1, listed as follows

- DoB: date of birth of the individual,
- DoS: date of start. For contributors, it is the date of subscribing. For annuitants, the date of entry in dependency.
- DoE and CoE: Respectively the date of end and cause of end for the individuals. In the case where the observation ends because of death, we use code 1 for the cause, in the case of exit because of entry in dependency, we use code 2. For individuals still autonomous when the observation stops, trajectories are right-censored. We use code 0 and the date of exit is the date at which observation ends.

DoB	DoS	DoE	CoE
23/12/1941	11/10/1992	27/09/2006	2
14/06/1926	28/03/1997	31/12/2013	0
17/04/1937	28/03/1995	04/08/2003	1

Table 1: Example of a database of contributors.

Other covariates such as gender, type of residence (home or facility), marital status, cause of dependency, amount of annuity and premium for aggravated risk may be available and bring useful additional informations. In what follows, we assume that the only covariate available is gender and we estimate a different model for each gender.

The observation period must often be limited in some way

- By setting the date at which the observation ends. With each database is associated a date of extraction, which is the date of the latest entry in the database. In practice, some claims are reported up to one year after their occurrence, which may result in some missing information during the last year of observation. It may therefore be a good idea to arbitrarily set a date for the end of the observation one year prior to the date of extraction, in order not to underestimate the number of events. It may then be required to change some date of ending accordingly and set the associated code to 0.
- By removing the first years of individual exposure. In individual long-term care insurance, insured often have to pass some sort of medical selection. While it is generally very simple, with just a few questions, it has a huge impact on incidence probabilities in dependency for the first few years after subscribing. In order not to underestimate the incidence rate, it may be a good practice to discard the first 3 years of observation for each individual.
- By limiting the length of the observation period. When we study the behaviour of a population for a specific risk, it may change over time. There are several factors involved, such as changes in the definition of dependency, underwriting and claim management policy or underlying biometric changes. Limiting the length of the observation study is therefore required. The optimal observation period would be around 5 years. However, due to limited volume of data, it may sometimes be necessary to consider a longer period.

Once the data has been processed, we can put it in a form which can be used directly to fit the model for biometric laws

- The age of entry  $x = \text{DoS} - \text{DoB}$ ,
- The age of exit  $y = \text{DoE} - \text{DoB}$ ,
- The cause of exit  $c = \text{CoE}$ .

## 2.4 Parametric modelling of the intensities

In this section, we propose to rely on a parametric expression for each of the transition intensities in the model.

### 2.4.1 Intensity of general mortality

For the intensity of general mortality of the portfolio, we use an external reference through a relational model as described in Brass (1971, 1974) or Hannerz (2001).

Let  $F_g$  be the cumulative distribution function associated with the intensity of general mortality  $\mu_g$  such that

$$F_g(x) = 1 - \exp\left(-\int_{x_0}^x \mu_g(u) du\right).$$

Then we define the cumulative distribution odds (CDO) by

$$\text{CDO}_g(x) = \frac{F_g(x)}{1 - F_g(x)}.$$

In his model, Brass makes the assumption that the logarithm of the CDO associated with the mortality of the standard population and the specific population are parallel curves. We denote by  $F_g^{ref}$  (resp.  $\mu_g^{ref}$ ) the cumulative distribution function (resp. intensity of mortality) associated with the mortality of the standard population. Under this assumption, a natural estimator for the intensity of general mortality of the portfolio is

$$\widehat{\mu}_g(x) = \frac{\widehat{\beta} \mu_g^{ref}(x)}{1 - (1 - \widehat{\beta}) F_g^{ref}(x)}$$

where  $\widehat{\beta}$  is the solution of the equation

$$\sum_x D_x = \sum_x D_x^{ref} \frac{\widehat{\beta} N_x}{N_x^{ref} (1 - (1 - \widehat{\beta}) F_g^{ref}(x))}$$

where  $D_x$  and  $N_x$  (resp.  $D_x^{ref}$  and  $N_x^{ref}$ ) are the number of death observed and the number of years lived between ages  $x$  and  $x + 1$  for the specific population (resp. for the standard population). Brass

model only requires the estimation of a single parameter  $\hat{\beta}$ . It gives a new estimator which converges smoothly towards the mortality reference at higher ages while predicting the same number of deaths as in the empirical data.

### 2.4.2 Intensity of incidence in dependency

Let us consider the following parametric expression for the intensity of incidence in dependency

$$\lambda(x) = \frac{\exp(a_\lambda x + b_\lambda)}{1 + \exp(a_\lambda x + c_\lambda)} + d_\lambda \quad (7)$$

with  $a_\lambda > 0$ ,  $b_\lambda, c_\lambda \in \mathbb{R} \cup \{-\infty\}$  and  $d_\lambda \geq 0$ .

This is a logistic formula with 4 parameters, which was introduced by Perks (1932) under an equivalent form to explain human mortality. A few remarks can be made about this model. Parameter  $d_\lambda$  is the initial intensity, present at all ages. If we set  $d_\lambda = 0$  we obtain the logistic model with 3 parameters presented in Beard (1959, 1971). Parameter  $a_\lambda$  is related to the growth rate of the intensity, which is comprised between 0 and  $a_\lambda$  with the intensity at the tipping point of the function being  $a_\lambda/2$ . Parameter  $c_\lambda$  gives the position of the tipping point of the logistic function, reached for  $x = -c_\lambda/a_\lambda$ . Finally, parameter  $b_\lambda$  allows to set the asymptotic limit of the intensity, which is given by  $d_\lambda + \exp(b_\lambda - c_\lambda)$ .

In the case where  $c_\lambda = -\infty$ , we obtain the exponential models introduced in Gompertz (1825) and Makeham (1867) that can therefore be considered as sub-models of the logistic model.

Experience from insurers shows that the intensity of incidence in dependency increases exponentially with respect to age. At higher ages, data becomes scarcer. As dependency is linked to ageing, it is reasonable to expect that the behaviours of mortality and morbidity are quite similar and that an exponential or logistic form is suited to model incidence in dependency. The logistic formula has already been used to this extent, e.g. in Rickayzen and Walsh (2002).

For an individual  $p$  defined by his/her age of entry in the portfolio  $x_p \geq 0$ , his/her age of exit  $y_p > x_p$  and the associated exit cause  $c_p \in \{0, 1, 2\}$  the partial log-likelihood has the following expression

$$\begin{aligned} l_p(\lambda) &= \log \left[ \exp \left( - \int_{x_p}^{y_p} \lambda(u) du \right) \lambda(y_p)^{\delta_{c_p}^2} \right] \\ &= \delta_{c_p}^2 \log(\lambda(y_p)) - \int_{x_p}^{y_p} \lambda(u) du \\ &= \delta_{c_p}^2 \log \left( \frac{\exp(a_\lambda y_p + b_\lambda)}{1 + \exp(a_\lambda y_p + c_\lambda)} + d_\lambda \right) - \frac{\exp(b_\lambda - c_\lambda)}{a_\lambda} \log \left( \frac{1 + \exp(a_\lambda y_p + c_\lambda)}{1 + \exp(a_\lambda x_p + c_\lambda)} \right) - d_\lambda (y_p - x_p), \end{aligned}$$

where for  $k, l \in \mathbb{N}$ ,  $\delta_k^l = \begin{cases} 1 & \text{if } k = l, \\ 0 & \text{otherwise.} \end{cases}$

### 2.4.3 Intensity of dependent mortality

In this section, we consider a mixture of two populations and within each population we assume the mortality is the sum of a common term and a population specific mortality which only depends on the age of entry in dependency  $x$ . Under those assumptions, we show that the resulting mortality of the mixture takes a peculiar parametric form, which is semi-Markov. In what follows, we discuss how those assumptions apply to the population of dependent people and we then propose a parametric model of the previously mentioned form for the mortality of dependent people.

**Lemma.** *Let us consider a model with 2 distinct states of dependency  $I_1$  and  $I_2$ , such that the respective transition intensities from autonomy to those states are  $\lambda_1(x)$  and  $\lambda_2(x)$  and no transition is allowed between those states or back to autonomy (see Figure 3). Let us assume that the intensity of mortality for dependent in state  $I_k$  can be written as*

$$\mu_{i,k}(x, t) = \mu_0(x, t) + \Delta_k(x)$$

with  $x \geq x_0$ ,  $t \geq 0$  and  $k \in \{1, 2\}$ .



Then to ensure the consistency between this model and the 3-state model, the following relations must be satisfied

$$\lambda(x) = \lambda_1(x) + \lambda_2(x)$$

$$\mu_i(x, t) = \mu_0(x, t) + \theta_1(x) + \frac{\theta_2(x)}{1 + \theta_3(x) \exp(\theta_2(x)t)}$$

where

$$\begin{cases} \theta_1(x) = \Delta_1(x) \\ \theta_2(x) = \Delta_2(x) - \Delta_1(x) \\ \theta_3(x) = \frac{\lambda_1(x)}{\lambda_2(x)}. \end{cases}$$

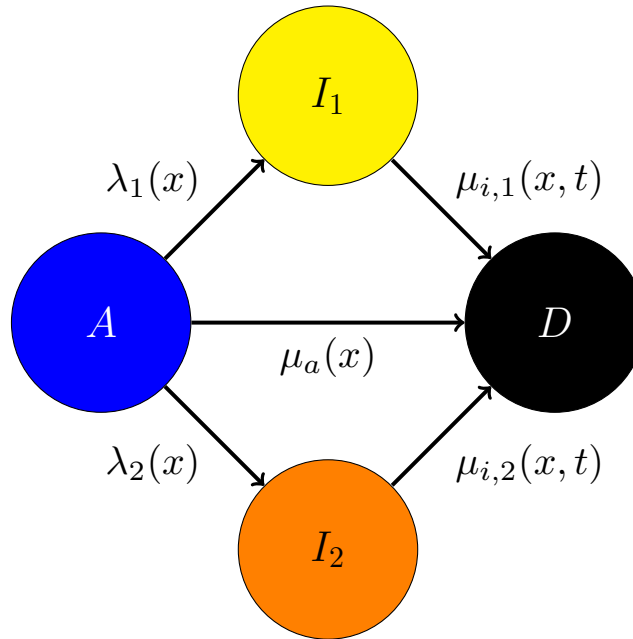


Figure 3: Model with 2 completely separate states of dependency.

Let us discuss the assumptions made in the lemma. The dependent population is heterogeneous because of the underlying pathologies that can cause dependency. Main causes of dependency are dementia, cancer, neurological and cardiovascular diseases. From an insurance perspective, we are concerned about the mortality associated with those pathologies starting from the entry in the state of functional limitation defined as dependency. On one hand people affected with terminal cancer, as well as people who had multiple strokes or a high severity infarction have a very high probability of dying within the few months following their entry in dependency. On the other hand, people affected with dementia, neurological diseases or less severe cardiovascular diseases may expect to survive for several years. It then becomes natural to see the dependent population as a mixture. We further assume that it has only two components which corresponds to the two groups of pathologies mentioned above and that within each component, the mortality of the population is homogeneous. A more realistic model would consider as much components as the number of pathologies, however it would prove very difficult to assess the mortality for each component given the limited amount of data available.

Another key assumption of the lemma is the additive mortality structure within each population. As regards dependency, the common term would be the mortality of autonomous people, and the specific term would be the additional mortality linked to the pathology responsible for dependency. The additive structure is based on the assumption that dependent people have increased mortality because of the pathology that caused dependency but are still exposed to other causes of death. The specific mortality associated with the cause of dependency is determined at the age of entry in dependency and does not change over time. This assumption is paramount if we want the structure of the mixture to remain understandable. It is a strong assumption from a medical point of view. However, given the average

duration of dependency, most of the time age of entry in dependency and current age remain within a few years of each other so the impact should be quite limited.

In the lemma we also assume that there is no transition between the two states of dependency. Of course, if someone was to be afflicted by dementia and then by a terminal cancer, its mortality level would better match the mortality of the group with cancer, while our model gives it the mortality of the group with dementia, because it is the first cause of dependency. However, such cases remain very rare, as the weight of cancer among causes of dependency greatly decreases with age.

*Proof.* The first relation on the incidences in dependency is obvious. For the second relation, let us denote by  $\eta_k(x, t)$  the proportion of dependent in state  $I_k$  among the population of people who became dependent at age  $x \geq x_0$  and then survived for a time  $t \geq 0$ .

We have for  $x \geq x_0, t \geq 0, k \in \{1, 2\}$

$$\begin{aligned} \eta_k(x, t) &= \frac{P(Z_{x+t} = I_k | Z_{x^-} = A, Z_x \in \{I_1, I_2\})}{P(Z_{x+t} \in \{I_1, I_2\} | Z_{x^-} = A, Z_x \in \{I_1, I_2\})} \\ &= \frac{P(Z_x = I_k | Z_{x^-} = A, Z_x \in \{I_1, I_2\}) \times P(Z_{x+t} = I_k | Z_{x^-} = A, Z_x = I_k)}{\sum_{l=1}^2 P(Z_x = I_l | Z_{x^-} = A, Z_x \in \{I_1, I_2\}) \times P(Z_{x+t} = I_l | Z_{x^-} = A, Z_x = I_l)} \\ &= \frac{\lambda_k(x) \exp\left(-\int_0^t \mu_{i,k}(x, u) du\right)}{\sum_{l=1}^2 \lambda_l(x) \exp\left(-\int_0^t \mu_{i,l}(x, u) du\right)}. \end{aligned}$$

The intensity of the mortality for the population of dependent people is

$$\begin{aligned} \mu_i(x, t) &= \sum_{k=1}^2 \eta_k(x, t) \mu_{i,k}(x, t) \\ &= \sum_{k=1}^2 \frac{\lambda_k(x) \exp\left(-\int_0^t \mu_{i,k}(x, u) du\right)}{\sum_{l=0}^2 \lambda_l(x) \exp\left(-\int_0^t \mu_{i,l}(x, u) du\right)} \mu_{i,k}(x, t) \\ &= \sum_{k=1}^2 \frac{\lambda_k(x) \exp(-\Delta_k(x)t)}{\sum_{l=0}^2 \lambda_l(x) \exp(-\Delta_l(x)t)} \mu_{i,k}(x, t) \\ &= \mu_0(x, t) + \sum_{k=1}^2 \frac{\lambda_k(x) \exp(-\Delta_k(x)t)}{\sum_{l=0}^2 \lambda_l(x) \exp(-\Delta_l(x)t)} \Delta_k(x) \\ &= \mu_0(x, t) + \Delta_1(x) + \frac{\lambda_2(x) \exp(-\Delta_2(x)t)}{\lambda_1(x) \exp(-\Delta_1(x)t) + \lambda_2(x) \exp(-\Delta_2(x)t)} [\Delta_2(x) - \Delta_1(x)] \\ &= \mu_0(x, t) + \Delta_1(x) + \frac{\Delta_2(x) - \Delta_1(x)}{1 + \frac{\lambda_1(x)}{\lambda_2(x)} \exp([\Delta_2(x) - \Delta_1(x)]t)} \end{aligned}$$

which proves the lemma. □

Let us consider a mixture model as described in the lemma with the following parametrization

$$\begin{cases} \mu_0(x, t) = \mu_a(x + t) \\ \Delta_1(x) = \phi_1 + \exp(\phi_2 x + \phi_3) \\ \Delta_2(x) = \Delta_1(x) + \phi_4 \\ \lambda_1(x) = \frac{\exp(\phi_5 x + \phi_6)}{1 + \exp(\phi_5 x + \phi_6)} \lambda(x) \\ \lambda_2(x) = \lambda(x) - \lambda_1(x) \end{cases}$$

with  $\phi_3, \phi_5, \phi_6 \in \mathbb{R}$  and  $\phi_1, \phi_2, \phi_4 \geq 0$ .

In this model, intensities of dependent mortality for specific populations increase with age of entry in dependency according to two parallel Makeham's models. For a given pathology, associated mortality is expected to increase with age, as older people prove more fragile. The choice of Makeham's model, which includes both a constant term and a term which increases exponentially with age, seems adapted. We assume the exponential term is common for the two components mainly for computational purposes. Should the difference  $\Delta_2(x) - \Delta_1(x)$  depend on  $x$ , then estimation of parameters would become hazardous. To set the level of the intensity of incidence in dependency for the two components, we choose to model directly the weight of each component in the mixture at entry in dependency, using a logistic law so that weights remain between 0 and 1.

The resulting intensity of mortality for the aggregated population takes the form presented in the lemma with

$$\begin{cases} \theta_1(x) = \phi_1 + \exp(\phi_2x + \phi_3), \\ \theta_2(x) = \phi_4, \\ \theta_3(x) = \exp(\phi_5x + \phi_6), \end{cases}$$

which leads to

$$\mu_i(x, t) = \underbrace{\mu_a(x + t) + \phi_1 + \exp(\phi_2x + \phi_3) + \frac{\phi_4}{1 + \exp(\phi_4t + \phi_5x + \phi_6)}}_{\Delta(x, t)}.$$

The associated partial log-likelihood for an individual  $p$  with an age of entry in dependency  $x_p \geq 0$ , an age of exit  $y_p > x_p$  and the associated cause of exit  $c_p \in \{0, 1\}$  is

$$\begin{aligned} l_p(\mu_a, \Delta) &= \log \left[ \exp \left( - \int_{x_p}^{y_p} \mu_i(x_p, u - x_p) du \right) \mu_i(x_p, y_p - x_p)^{\delta_{c_p}^1} \right] \\ &= \delta_{c_p}^1 \log(\mu_i(x_p, y_p - x_p)) - \int_{x_p}^{y_p} \mu_i(x_p, u - x_p) du \\ &= \delta_{c_p}^1 \log \left( \mu_a(y_p) + \phi_1 + \exp(\phi_2x_p + \phi_3) + \frac{\phi_4}{1 + \exp(\phi_4[y_p - x_p] + \phi_5x_p + \phi_6)} \right) - \int_{x_p}^{y_p} \mu_a(u) du \\ &\quad - (\phi_1 + \exp(\phi_2x_p + \phi_3) + \phi_4) [y_p - x_p] + \log \left( \frac{1 + \exp(\phi_5x_p + \phi_6)}{1 + \exp(\phi_4[y_p - x_p] + \phi_5x_p + \phi_6)} \right). \end{aligned}$$

## 2.5 Parameters estimation procedure

To estimate the parameters, we use the following procedure

1. We estimate the parameters for the intensity of incidence in dependency  $\hat{\lambda}$  (resp. the autonomous mortality  $\hat{\mu}_a^{(1)}$ ), using the contributors database and Perks logistic model. More precisely  $\hat{\lambda}$  is the maximum likelihood estimator (MLE) constructed by summing the partial log-likelihoods given in section 2.4.2 and  $\hat{\mu}_a^{(1)}$ , the MLE for the intensity of autonomous mortality, takes a similar form.
2. We estimate the parameters for the intensity of general mortality  $\hat{\mu}_g$  by using the individuals of both databases and Brass relational model in order to get a robust estimate of the intensity of general mortality with forced convergence towards a reference mortality table at higher ages.
3. We estimate the parameters for the intensity of dependent mortality  $\hat{\Delta}$ , thanks to the parametric semi-Markov mixture model introduced previously, to  $\hat{\mu}_a^{(1)}$  and to the annuitant database. To do so, we use the MLE constructed by summing the partial log-likelihoods given in section 2.4.3.
4. Thanks to equation (3), we compute the value of a second estimator for the intensity for autonomous mortality  $\hat{\mu}_a^{(2)}$ , relying on  $\hat{\lambda}$ ,  $\hat{\Delta}$ ,  $\hat{\mu}_g$  and using numerical methods to approximate the outer integrals in (3).

A summary of the procedure is provided in Figure 4.

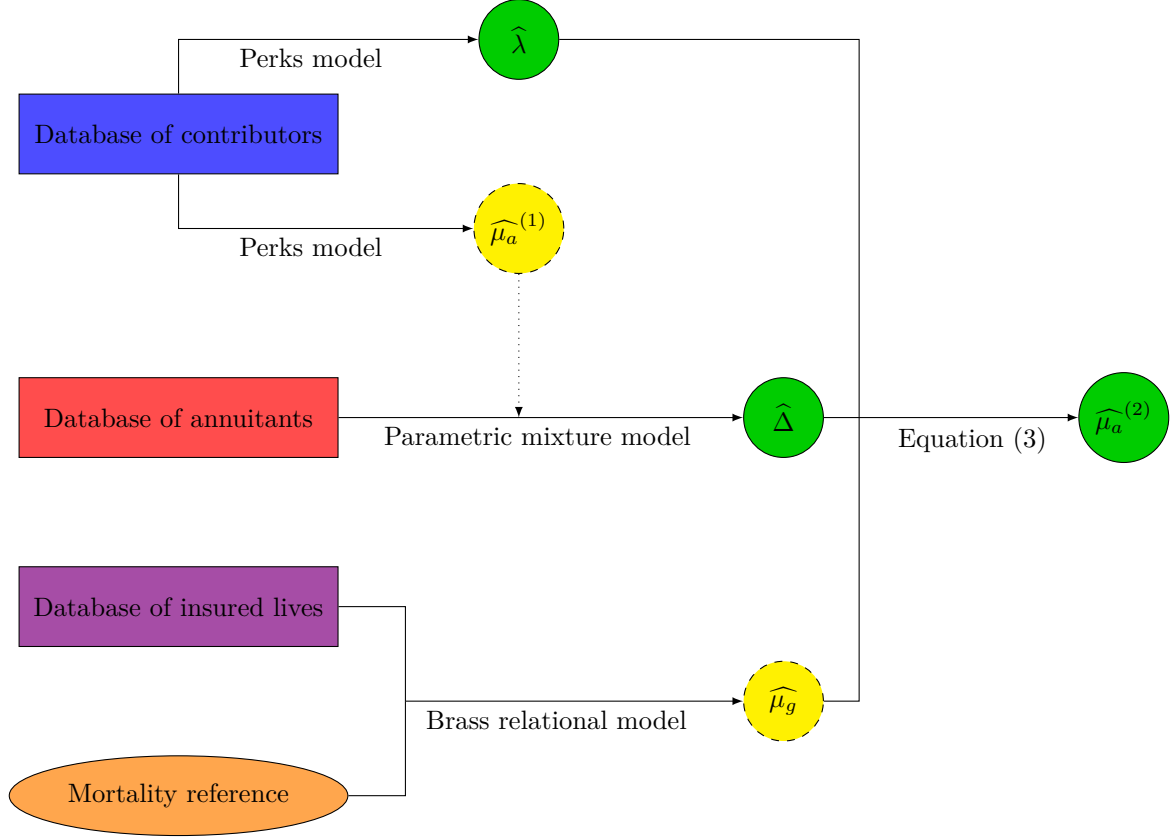


Figure 4: Procedure for the estimation of biometric laws. Dashed (resp. plain) circles represent intermediary (resp. final) estimates of biometric laws.

## 2.6 Link with annual probabilities

Once we have proceeded to the estimation of parameters for transition intensities, annual probabilities of transition can easily be derived. Those probabilities are generally used by insurers for pricing and reserving purpose. In the following section, we introduce pricing and reserving formula based on transition intensities directly. However, expressions of annual probabilities of transition are still interesting for comparison purposes.

We denote, for  $x \geq x_0$  and  $t \geq 0$

$$\begin{aligned}
 q_g(x) &= P(Z_{x+1} = D | Z_x \in \{A, I\}), \\
 q_{aa}(x) &= P(Z_{x+1} = D, \forall y \in [x; x+1], Z_y \neq I | Z_x = A), \\
 i(x) &= P(\exists y \in [x; x+1], Z_y = I | Z_x = A), \\
 q_i(x, t) &= P(Z_{x+t+1} = D | Z_{x-} = A, Z_x = I, Z_{x+t} = I).
 \end{aligned}$$

For  $x \geq x_0$  and  $t \geq 0$

$$\begin{aligned}
 q_g(x) &= 1 - \exp\left(-\int_x^{x+1} \mu_g(u) du\right), \\
 q_{aa}(x) &= \int_x^{x+1} \mu_a(u) \exp\left(-\int_x^u [\lambda(u) + \mu_a(u)] du\right) du, \\
 i(x) &= \int_x^{x+1} \lambda(u) \exp\left(-\int_x^u [\lambda(u) + \mu_a(u)] du\right) du, \\
 q_i(x, t) &= 1 - \exp\left(-\int_t^{t+1} \mu_i(x, u) du\right).
 \end{aligned}$$

## 2.7 Pricing and reserving

We note  $\tau$  the continuous time actuarial rate used to compute discounted cash flows. We consider a product where the autonomous insured life pays a fixed amount of premium at the start of every period of duration  $1/f_1$  (in year). Should he/she become dependent, he/she is entitled to an annuity  $R/f_2$  paid at the end of every period of duration  $1/f_2$ . On the French long-term care insurance market most products rely on monthly premium and monthly benefit, which means  $f_1 = f_2 = 12$ . Furthermore, let us consider there is an age  $\omega$  such that  $D_\omega \simeq 1$ . In practice,  $\omega = 120$  will be used for applications.

Let us introduce additional notation

$$A(x, y) = P(Z_y = A | Z_x = A) = \frac{A_y}{A_x} = \exp\left(-\int_x^y [\mu_a(u) + \lambda(u)] du\right),$$

$$I_x(t, s) = P(Z_{x+s} = I | Z_{x-} = A, Z_x = I, Z_{x+t} = I) = \exp\left(-\int_t^s \mu_i(x, u) du\right)$$

and

$$\bar{A}(x, y) = e^{-\tau(y-x)} A(x, y) = \exp\left(-\int_x^y [\mu_a(u) + \lambda(u) + \tau] du\right),$$

$$\bar{I}_x(t, s) = e^{-\tau(s-t)} I_x(t, s) = \exp\left(-\int_t^s [\mu_i(x, u) + \tau] du\right)$$

for  $x_0 \leq x \leq y$  and  $0 \leq t \leq s$ .

Let us introduce

- $P(x_s, x, f_1)$  the expected value of insured liabilities for an autonomous insured life with age at subscribing  $x_s$  and current age  $x$  for a 1 € premium

$$P(x_s, x, f_1) = \frac{1}{f_1} \sum_{k \in \mathbb{N} \cap [f_1(x-x_s); f_1(\omega-x_s)]} \bar{A}(x, x_s + \frac{k}{f_1}),$$

- $RFC(x_i, t, R, f_2)$  the expected value of insurer liabilities for a dependent insured life, also called reserve for claim. For an amount of benefit  $R$ , an age  $x_i$  at entry in dependency and a time  $t$  spent in dependency, the corresponding amount of reserve is

$$RFC(x_i, t, R, f_2) = \frac{1}{f_2} \sum_{k \in \mathbb{N} \cap (f_2 t; f_2(\omega-x_i)]} \bar{I}_{x_i}(t, \frac{k}{f_2}).$$

- $B(x, R, f_2)$  the expected value of insurer liabilities for an autonomous insured life with current age  $x$

$$B(x, R, f_2) = \int_x^\omega \lambda(u) \bar{A}(x, u) RFC(u, 0, R, f_2) du$$

- The stability premium  $p^*(R, x_s)$ . It is the value of premium that matches insurer and insured liabilities at the time of subscribing. For an age  $x_s$  at subscribing and an amount of benefit  $R$  we have

$$p^*(x_s, R, f_1, f_2) = \frac{B(x_s, R, f_2)}{P(x_s, x_s, f_1)}.$$

- The reserve for premium (RFP). This reserve is constituted for autonomous people. Its amount is equal to the expectancy of future discounted cash flows of benefit minus discounted cash flows of premium. For an insured of age at subscribing  $x_s$ , current age  $x$ , an amount of premium  $p$  and an amount of benefit  $R$ , the associated amount of reserve is

$$RFP(x_s, x, p, R, f_1, f_2) = B(x, R, f_2) - pP(x_s, x, f_1).$$

## 2.8 Model selection

When we deal with complex models, it can be very interesting to compare sub-models to see if the use of many parameters is justified. To this extent, we can rely on the *Bayesian Information Criterion (BIC)*. For a model  $m_i$  characterized by a number of parameters  $p_i$  and a log-likelihood function  $l_i$  maximized at  $\theta_i$ , the expression of the criterion is as follows

$$BIC_i = -2l_i(\theta_i) + p_i \log(n)$$

where  $n$  represents the number of observed transitions in the expression of the likelihood. The choice of the coefficient in front of the number of parameters  $\log(n)$  differs from the one made in the original Akaike's Information Criterion (AIC) where this coefficient is 1. Also, let us note that in the present version of the criterion we consider the number of observed transition and not the number of individuals as in the original criterion. The interest in introducing this modification in the case of censored data is discussed in Volinsky and Raftery (2000). Using the BIC, we are able to compare models, the model with the lower BIC being the best model.

## 3 Results

In this section, we provide an example using data from a long-term care insurance portfolio. The definition used for dependency is "3ADL4" which means that an insured life is considered dependent if he/she has permanently lost the ability to do on his/her own at least 3 among the 4 activities of daily living defined by the contract (functional mobility, dressing, bathing, eating). The portfolio we consider contains a large number of policies and covers a relatively long period. We have information about policyholders until the age of 95 and trajectories in dependency lasting up to 13 years. The date of extraction is 11/31/2013 for both contributors and annuitants. We consider a 10 year observation period between 1/1/2002 and 12/31/2011 for contributors, and remove the first 3 years spent in the portfolio. For annuitants, we consider the observation period between 1/1/1994 and 12/31/2012, and keep the full exposure over this period. Database of contributors contains 160,669 individuals after the data has been processed, for a total of 1,325,578 years of exposure with 75.9 % of the lines being right censored. Database of annuitants contains a total of 17,632 individuals for a total of 43,010 years of exposure and 31.4 % of right censored individuals. Women account for 65.3 % of contributors and 65.9 % of annuitants. Separate models are estimated for men and women.

### 3.1 General mortality

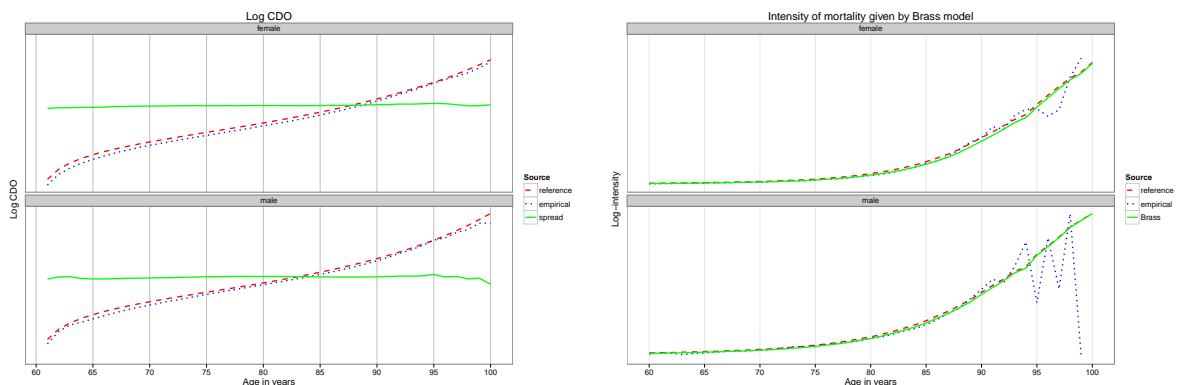


Figure 5: Left: Logarithm of cumulative distribution odds for empirical mortality (dotted) and reference mortality (dashed). Their difference (plain) should be a straight line. Right: Intensity of mortality estimated from the data (dotted), from the mortality reference (dashed) and resulting from Brass model (plain). The  $y$ -scale has been removed to preserve confidentiality of results.

We use Brass relational model with a reference mortality table. Figure 5 (Left) displays the logarithm of the cumulative distribution odds for empirical probabilities and the mortality reference. Figure 5 represents the empirical and reference intensities of mortality, as well as the intensity fitted by Brass model. The intensity fitted by the model is close to the empirical value but converges toward the reference at higher ages, when no empirical data is available.

### 3.2 Incidence in dependency

The results of the estimation of incidence in dependency can be found in Table 2. According to the BIC, Gompertz's model is the best model for men and Beard's model the best for women. The choice of Beard's model however leads to an asymptotic value for the incidence of 0,155. As far as we know, this slowing down in the growth of incidence has not been observed in any other study. In addition, by looking at figure 6, it does not strike us that Beard's model is a better fit, and the data is really scarce at higher ages. In what follows, in order to take a conservative estimate, we use Gompertz model, which gives higher incidence at higher ages, for both men and women.

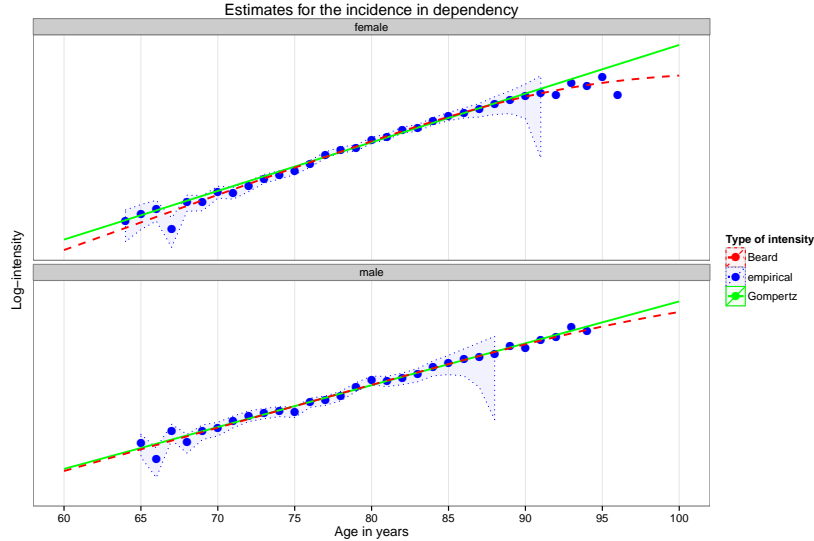


Figure 6: Estimates of the incidence in dependency. Dots and ribbon represent empirical estimates with confidence intervals, when data is sufficient. Plain (resp. dashed) line represent the Gompertz (resp. Beard) model fitted to the data. The  $y$ -scale has been removed to preserve confidentiality of results.

Model	Intensity	$l$ (men)	BIC (men)	$l$ (women)	BIC (women)
<b>Gompertz</b>	$\lambda(x) = e^{a_\lambda x + b_\lambda}$	- 21,383.00	<b>42,782.61</b>	- 42,286.02	84,590.07
Makeham	$\lambda(x) = e^{a_\lambda x + b_\lambda} + d_\lambda$	- 21,383.00	42,790.91	- 42,286.02	84,599.09
<b>Beard</b>	$\lambda(x) = \frac{e^{a_\lambda x + b_\lambda}}{1 + e^{a_\lambda x + c_\lambda}}$	- 21,382.40	42,789.71	- 42,263.72	<b>84,554.49</b>
Perks	$\lambda(x) = \frac{e^{a_\lambda x + b_\lambda}}{1 + e^{a_\lambda x + c_\lambda}} + d_\lambda$	- 21,382.05	42,797.33	- 42,262.72	84,561.51

Table 2: Value of log-likelihood  $l$  and value of  $BIC$  for men and women, for several models of intensity for the incidence in dependency.

### 3.3 Dependent mortality

For the intensity of dependent mortality, we consider 5 different models. The first model is an additive model. The dependent population is assumed to be homogeneous and the difference between the intensity of mortality for autonomous and dependent people is flat

$$\mu_i(x, t) = \mu_a(x, t) + \phi_1. \quad (8)$$

The second model is a time-homogeneous mixture model. It relies on the idea that there are two populations of dependent people, but does not consider that the incidence and specific mortality associated with each population change with respect to age of entry in dependency

$$\mu_i(x, t) = \mu_a(x + t) + \phi_1 + \frac{\phi_2}{1 + \phi_3 \exp(\phi_2 t)}. \quad (9)$$

The third model is the full 6 parameters model presented in Section 2.4.3, with

$$\mu_i(x, t) = \mu_a(x + t) + \phi_1 + \exp(\phi_2 x + \phi_3) + \frac{\phi_4}{1 + \exp(\phi_4 t + \phi_5 x + \phi_6)} \quad (10)$$

with  $\phi_1, \phi_2, \phi_4 \geq 0$  and  $\phi_3, \phi_5, \phi_6 \in \mathbb{R}$ .

For comparison purposes we also consider the historical Markov model introduced in SCOR (1995) which is a linear function of the autonomous mortality with intercept

$$\mu_i(x, t) = \alpha \mu_a(x + t) + \beta \quad (11)$$

with  $\alpha, \beta \geq 0$ .

We could also use a generalized linear model (GLM). In a Markov approach, we only consider the current age  $x + t$  as a predictor

$$\log \mu_i(x, t) = \log \mu_a(x + t) + \beta_0 + \beta_1(x + t) \quad (12)$$

but with a semi-Markov approach, we can consider separate coefficients to take into account both age at entry in dependency  $x$  and time spent in dependency  $t$

$$\log \mu_i(x, t) = \log \mu_a(x + t) + \beta_0 + \beta_1 x + \beta_2 t \quad (13)$$

with  $\beta_0, \beta_1, \beta_2 \in \mathbb{R}$ .

Model equation	l (men)	BIC (men)	l (women)	BIC (women)
(8)	- 8,991.16	17,999.26	- 17,935.40	35,888.59
(9)	- 8,345.68	16,725.23	- 17,459.71	34,955.01
<b>(10)</b>	- 8,277.13	<b>16,613.52</b>	- 17,383.60	<b>34,829.51</b>
(11)	- 8,991.15	18,007.70	- 17,933.31	35,893.32
(12)	- 8,985.00	17,995.39	- 17,911.98	35,850.66
(13)	- 8,760.70	17,555.26	- 17,823.84	35,683.29

Table 3: Value of log-likelihood  $l$  and value of  $BIC$  for men and women of previously listed models for the intensity of dependent mortality.

Table 3 gives the log-likelihood and the BIC associated with the different models. The first 3 models are embedded models of increasing complexity. As the value of the BIC diminishes significantly when we go from the first to the second and from the second to the third, the use of more parameters is validated by the criterion. Model (11) is a submodel of model (8) which has a higher value for the BIC. The criterion does not validate the use of the additional parameter  $\alpha$  in this case. Models (12) and (13) are not submodels of the other models. However, looking at the likelihood, it appears that they are better than model (8) but a lot worse than (9) and (10). Model (12) is a submodel of model (13) and the BIC shows that taking into account separate terms for age of entry and time spent in dependency clearly improves the quality of the model. In what follows, we only consider model (10) for applications.

Figure 7 represents the annual probabilities associated with the model. We compute empirical annual probabilities by grouping dependent people, according to their age of entry in dependency, in three categories: 60 to 70, 70 to 80 and 80 to 90. We then consider that individuals in one category have the same annual death probabilities which only depend on time spent in dependency. Under the normal approximation, we represent 95 % confidence intervals. Methods for the estimation of empirical probabilities and construction of confidence intervals can be found in Planchet and Thérond (2006). Even with a large volume of data and 10 years subdivisions, confidence intervals are still very wide, especially for men, at higher age and/or high time spent in dependency. Nevertheless, the annual probabilities computed using the model appear to match the empirical probabilities very well.



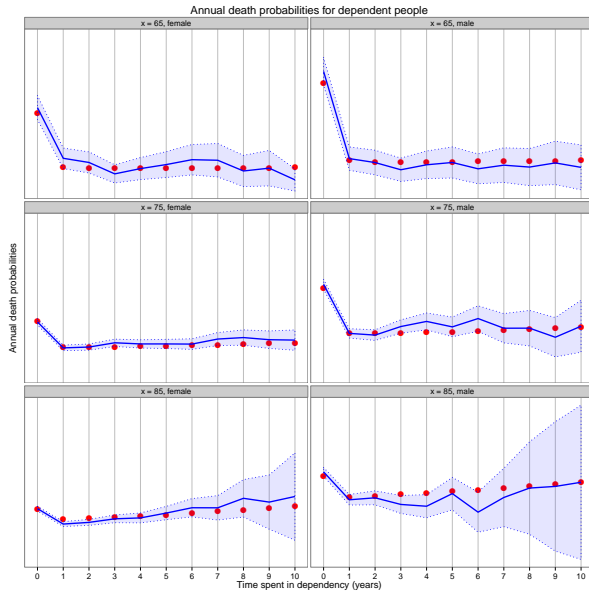


Figure 7: Consecutive death probabilities for dependent people according to the model (dots) with empirical probabilities (plain line) and 95 % confidence intervals (dashed lines). The  $y$ -scale has been removed to preserve confidentiality of results.

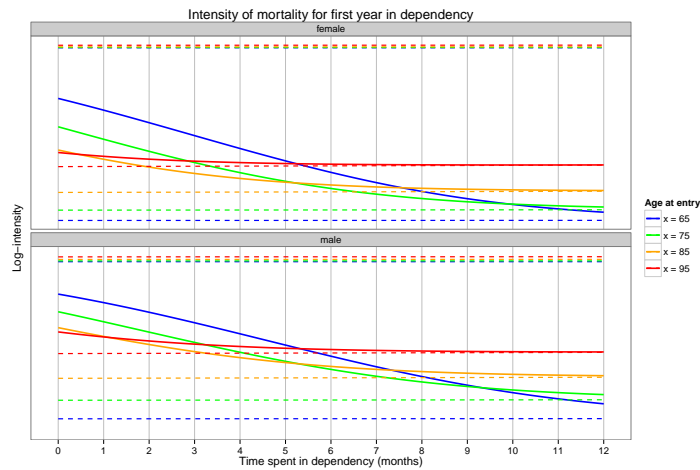


Figure 8: Intensity of dependent mortality over the first 12 months spent for several ages at entry in dependency (plain lines) with intensities of mortality among the two theoretical populations of the mixture appearing as dashed lines. The  $y$ -scale has been removed to preserve confidentiality of results.

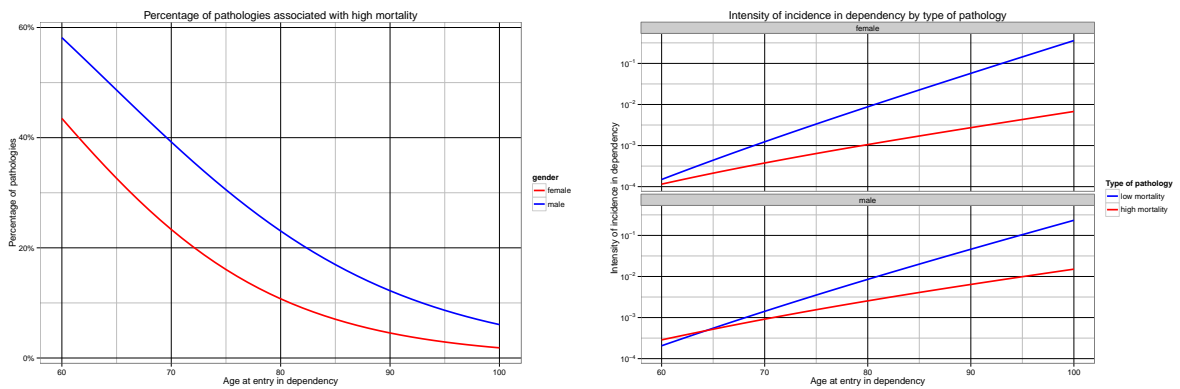


Figure 9: Left: prevalence of high mortality pathologies in the population of new dependents inferred by the model; right: Intensity of incidence in dependency by category of pathologies inferred by the model.

Figure 8 represents the intensity of dependent mortality for the first 12 months spent in dependency. Younger people have higher initial intensity of mortality but lower ultimate values. Figure 9 provides information regarding the breakdown of the population of new dependent and the incidence for each population of dependent, according to the model. One should keep in mind that pathologies are not observed in the data and this figure only represents the underlying distribution inferred by the model. The higher prevalence of high mortality pathologies among younger people explains the results of Figure 8.

### 3.4 Autonomous mortality

Figure 10 represents the intensity of autonomous mortality we get by trying to fit directly a logistic model to the data and the intensity we obtain at the end of the procedure, by using equation (3). We observe a significant difference at higher ages, where fitting directly the logistic model leads to overestimating the autonomous mortality, especially for women.

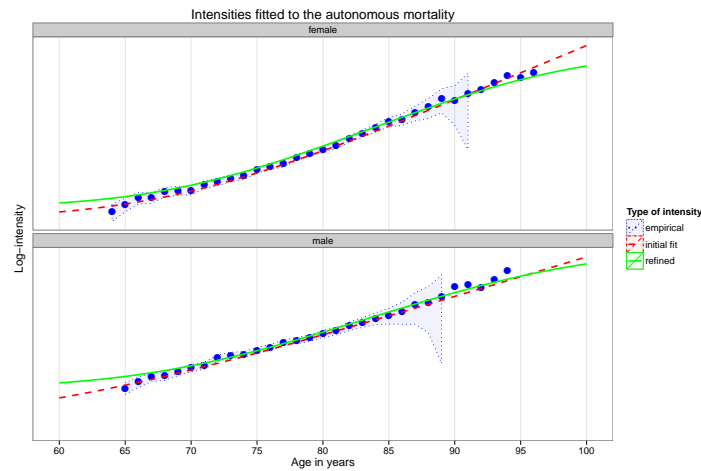


Figure 10: Intensity of autonomous mortality. Dots represents the empirical incidence rates estimated from the data. The dashed line represents the result of a direct fit of Perks logistic model and the plain line represents the intensity refined by using equation (3). The  $y$ -scale has been removed to preserve confidentiality of results.

### 3.5 Annual probabilities and prevalence of dependency

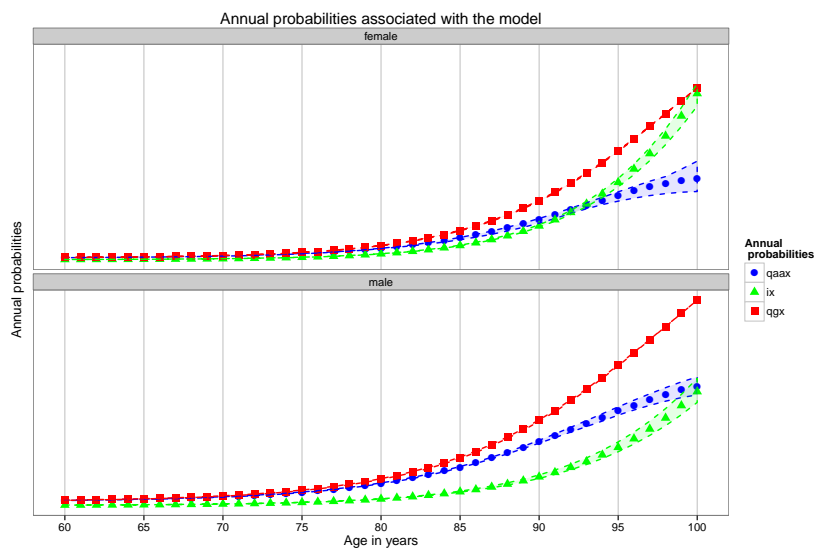


Figure 11: Annual probabilities for autonomous mortality in blue, general mortality in red and incidence in dependency in green, with 95 % confidence intervals obtained by bootstrap. The  $y$ -scale has been removed to preserve confidentiality of results.

In order to assess the robustness of the estimation performed, we use a classic non-parametric quantile bootstrap method. From the initial observation database, we build 100 new samples by drawing, with replacement, as many individuals as in the initial observation database. For each sample, we then run all the steps of the estimation procedures, and compute the associated annual probabilities as well as the prevalence of dependency, represented on Figure 11. We observe that general mortality and autonomous mortality diverge at higher ages. Figure 12 represents the prevalence of dependency and the number of dependent people as a fraction of the initial population at age 18. The peak of dependency is reached at age 89 for men and age 93 for women. The prevalence increases exponentially with respect to age, and is higher for men than women. Overall the confidence intervals are relatively tight, and the estimation method proves quite robust.

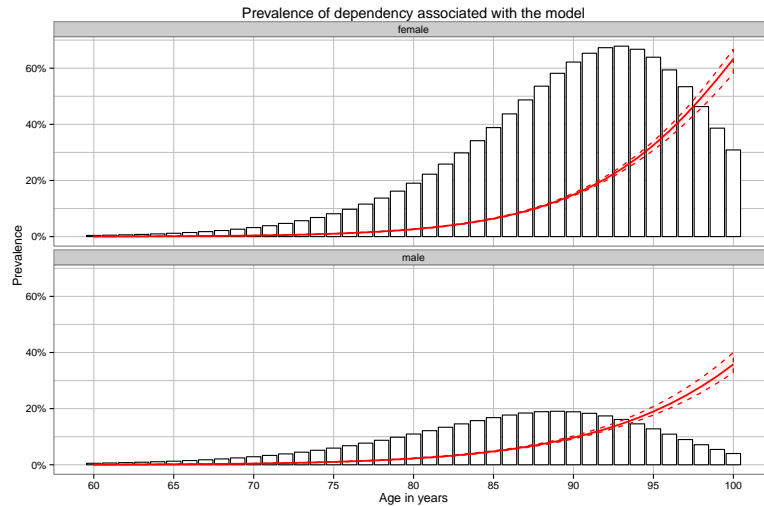


Figure 12: Prevalence of dependency by age (plain line), with 95 % confidence intervals obtained by bootstrap. The histogram represents 10 times the number of living dependent people as a percentage of the initial population at age 18.

### 3.6 Results of pricing and reserving

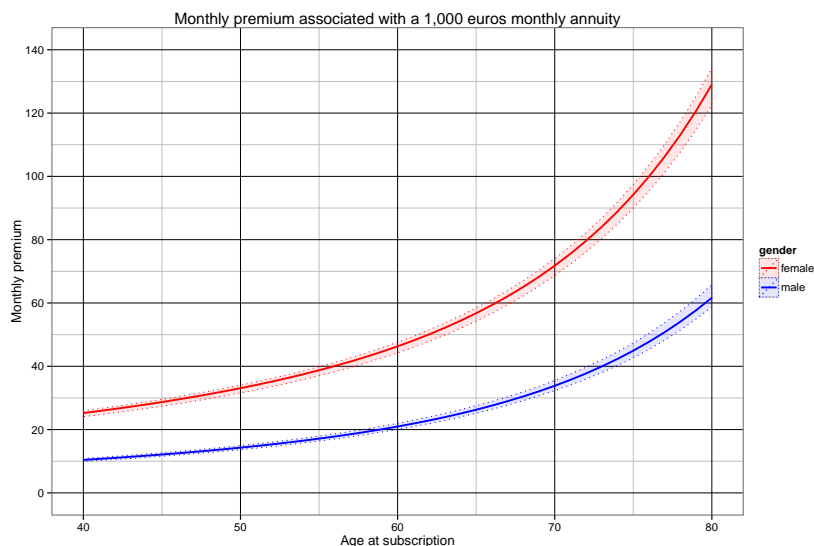


Figure 13: Amount of monthly premium required according to the model, with 95 % confidence intervals obtained by bootstrap. The  $y$ -scale has been removed to preserve confidentiality of results.

We consider a long-term care insurance product where policyholders pay a level premium, whose amount is fixed at subscribing, at the beginning of every month following subscribing while they are autonomous. Should they become dependent, they would stop paying the premium and instead receive

an annuity of 1,000 € at the end of every month following the entry in dependency. Most products in France work this way. We use a technical rate of 1 % for the pricing of the product. Figure 13 provides the required premium according to the model for ages at subscribing between 40 and 90, with confidence intervals obtained by bootstrap. The premium is higher for women than men. Confidence intervals are very tight, and the uncertainty assessed using bootstrap is very limited compared to other possible causes of uncertainty on the risk.

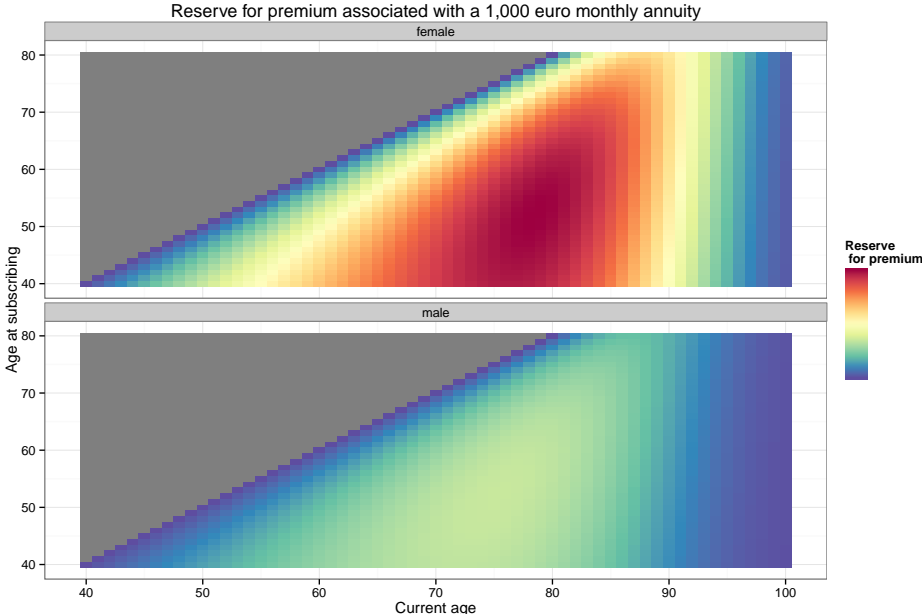


Figure 14: Expected value of reserve for premium by age at subscribing and current age, assessed at subscribing. Warm colors correspond to higher amount of reserve.

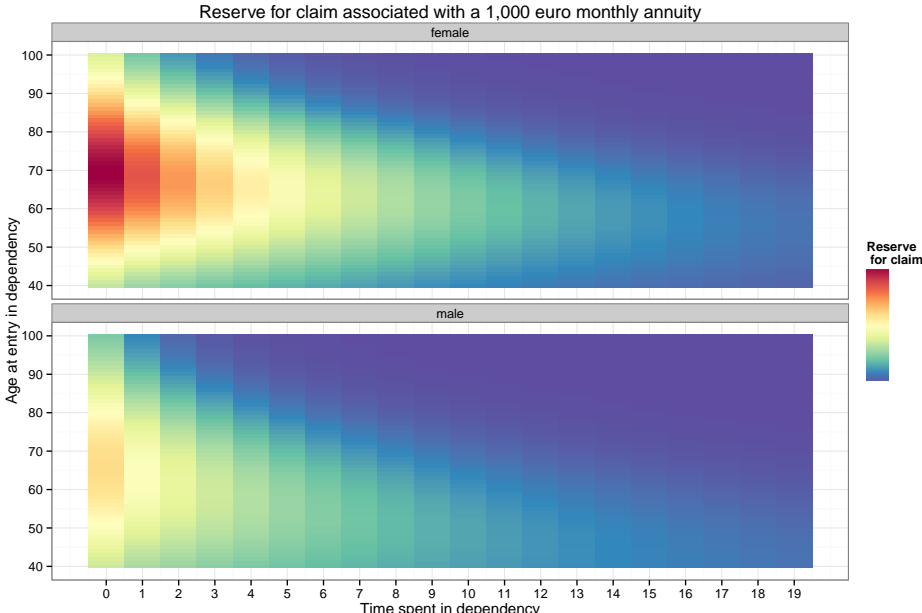


Figure 15: Expected value of reserve for claim by age at entry in dependency and time spent in dependency, assessed at claim inception. Warm colors correspond to higher amount of reserve.

We compute the expected value assessed at subscribing of the reserve for premium on Figure 14. This is the product between the amount of reserve computed in Section 2 and the probability for an individual to remain autonomous until the current age considered, as the reserve for premium only apply to autonomous insured lives. This reserve reaches a maximum between ages 75 and 88, depending on

the age at subscribing and then decreases when the cost associated with the claims starts to outweigh the amount of premium paid by the remaining autonomous insured lives. We also compute the expected value assessed at claim inception of the reserve for claim on Figure 15. This is the product between the probability to survive in dependency for the given duration and the reserve for claim computed in Section 2 associated with this same duration. This reserve decreases by duration as the survival probability gets lower. The initial amount of reserve reaches its maximum around 70 for both men and women. Indeed, younger claimants have very high death probabilities for the first year spent in dependency, and older claimants have higher death probabilities for the subsequent years. Claims coming from insured lives around 70 therefore prove the most expensive claims on average.

## 4 Discussion

In this paper, we introduce a method to estimate biometrics laws associated with a long-term care insurance portfolio. This method relies on a continuous time semi-Markov model, as opposed to discrete-time methods used by practitioners. We provide a formula to include general mortality in the model instead of autonomous mortality. We then suggest parametric models for all transition intensities. A logistic model is used for incidence in dependency and Brass relational model is used for the intensity of general mortality, while for dependent mortality we introduce a more complex model based on assumptions on a representation of the dependent population as a mixture of two groups associated with different pathologies. Estimation of parameters is performed using the maximum likelihood method. We also provide adequate formulas for pricing and reserving.

We then apply our methodology to a real long-term care insurance portfolio. Our model proves consistent with empirical estimations of annual probabilities, whenever those probabilities can be computed from the data available. Empirical probabilities for dependent mortality highlight that mortality during first year in dependency is way higher than for the subsequent years. Therefore a semi-Markov model which takes into account both the age and the time spent in dependency is required in order to explain this phenomenon.

For an insurance company, ignoring the complexity of the dependency process can be very damaging. Let us consider a company which launches a new insurance product and wants to review the level of reserves based on the first few years of experience from the portfolio. The exposure and death associated with the first year in dependency would be overrepresented in the data as many trajectories would be censored after just a few years. A Markov model only considers age to be relevant to assess transition probabilities. Such a model would therefore consider that dependent people have a very high mortality in dependency, not only for the first year but also for every subsequent year. Hence, it would greatly underestimate the time spent in dependency and the required amount of reserve, which would lead to heavy losses in the future.

In the present article, we try to assess several sources of error. As we use a parametric approach, there is a significant risk of model error. We therefore compare the results of the model with the empirical annual probabilities we obtain using a non-parametric approach. We consider several sub-models and try to remain parsimonious in the number of parameters by using the Bayesian Information Criterion to compare them. The robustness of estimation is also assessed using a non-parametric quantile bootstrap method.

The parametric form we introduce for dependent mortality stems from the assumption that pathologies can be sorted in two homogeneous groups. This assumption is unrealistic and it might further be improved by focusing on the study of the pathologies causing dependency. Data about pathologies however is extremely scarce and kept private by most insurers. Besides, our estimation approach is periodic and does not consider that intensities are changing over time. The estimation of drifts would indeed prove very difficult because of the limited observation period, and lack of consistency in definition of dependency, underwriting and claim management policies over time. Also, most products in France allow the insurer to increase the level premium in order to account for such drifts of the underlying risk. A sensitivity approach where we would consider several scenarios and look at the impact of the drifts on the premium could nevertheless prove very useful. However, as of today, to the best of our knowledge, neither the data nor the theoretical framework associated with this issue do exist. Finally, the model only covers one level of dependency, when most insurers provide products with several levels of guarantee. Extending the model to consider several levels of dependency, as in Lepez et al. (2013) or Biessy (2015) would therefore be very useful. Once again, finding real data to perform estimation of parameters proves very challenging.

## References

- Beard, R. (1959). Note on some mathematical mortality models. In: G.E.W. Wolstenholme and M. O'Connor (eds.). *The Lifespan of Animals. Ciba Foundation Colloquium on Ageing*. pp. 302-311. Little, Brown, Boston.
- Beard, R. (1971). Some aspects of theories of mortality, cause of death analysis, forecasting and stochastic processes. In: *Biological Aspects of Demography* (ed. W. Brass). London: Taylor and Francis.
- Biessy, G. (2015). Long-term care insurance: A multi-state semi-Markov model to describe the dependency process in elderly people. *Bulletin Français d'Actuariat* 15(29), 41–73.
- Brass, W. (1971). Mortality models and their uses in demography. *Transactions of the Faculty of Actuaries* 33, 123–142.
- Brass, W. (1974). Perspectives in population prediction: Illustrated by the statistics of England and Wales. *Journal of the Royal Statistical Society* 137(4), 532–583.
- Christiansen, M. C. (2012). Multistate models in health insurance. *Advances in Statistical Analysis* 96(2), 155–186.
- Cinlar, E. (1969). Markov renewal theory. *Advances in Applied Probability* 1, 123–187.
- Gompertz, B. (1825). On the nature of the function expressive of the law of human mortality, and on a new mode of determining the value of life contingencies. *Phil. Trans. R. Soc. Lond.* 115, 513–583.
- Haberman, S. and E. Pitacco (1998). *Actuarial Models for Disability Insurance*. Chapman and Hall/CRC, 1st edition.
- Hannerz, H. (2001). An extension of relational methods in mortality estimation. *Demographic Research* 4(10), 337–368.
- Lepez, V., S. Roganova, and A. Flahault (2013). A semi-Markov model to investigate the different transitions between states of dependency in elderly people. In: Colloquium of the International Actuarial Association, Lyon.
- Makeham, W. M. (1867). On the law of mortality. *Journal of the Institute of Actuaries* 13, 325–358.
- Perks, W. (1932). On some experiments in the graduation of mortality statistics. *Journal of the Institute of Actuaries* 63(1), 12–57.
- Pitacco, E. (2015). Actuarial values for long-term care insurance products. a sensitivity analysis. *Working Paper*, ARC Centre of Excellence in Population Ageing Research.
- Planchet, F. and P. Thérond (2006). *Modèles de durée*. Economica.
- Rickayzen, B. D. and D. E. P. Walsh (2002). A multi-state model of disability for the United Kingdom: Implications for future need for long-term care for the elderly. *British Actuarial Journal* 8, 341–393.
- SCOR (1995). L'assurance dépendance. Dossiers SCOR, SCOR Tech.
- Volinsky, C. T. and A. E. Raftery (2000). Bayesian information criterion for censored survival models. *Biometrics* 56(1), 256–262.