

Modelling and quantifying mortality and longevity risk

Module C.1 on Statistical and machine learning methods for portfolio data

Katrien Antonio & Michel Vellekoop

September 18-19, 2025

Actuarial Summer School, Faculty of Economic Sciences, University of Warsaw

In this module you will learn:

- different methods to analyze portfolio mortality using one or multiple risk factors
- pros and cons of the different approaches
- how information on insured amounts can be used to construct longevity assumptions for valuation purposes
- how machine learning methods can be used to do covariate ('feature') selection automatically.

Motivation

Research and industry perspectives

Risk classification approach

Relative approach

Industry approach

Wrap-up

Motivation

Population viz portfolio specific mortality

Actuarial associations (like IA|BE in Belgium, KAG in the Netherlands, Institute and Faculty of Actuaries in UK) publish **mortality projection models**:

- industry standard at national **population** level ('pop')
- insurance companies or pension funds then estimate **portfolio** correction factors ('pf'), e.g.

$$\mu_{x,t,g,\ell}^{pf} = \mu_{x,t,g}^{pop} \cdot \exp(f(\mathbf{x}_\ell)),$$

for age x , period t , gender g and risk profile ℓ (with covariates \mathbf{x}_ℓ).

Population viz portfolio specific mortality

- ▶ Read the motivation in BAJ (2015, 20, pp. 461-490), *A methodology for assessing basis risk - Abstract of the London discussion* and IFoA/LLMA (2014) *Longevity Basis Risk. A methodology for assessing basis risk by Cass Business School and Hymans Robertson LLP.*
- ▶ Longevity basis risk:
 - e.g., in a hedging instrument: longevity outcomes of **hedged portfolio** differs from **published mortality index**
 - e.g., in pricing of and reserving for life contingent risks: use of population vs insured portfolio data
 - hence, need to capture the **gap** – and evolutions over time in this gap – between mortality in the **reference population** and in the **portfolio**.

Research and industry perspectives

Focus in the actuarial literature on **two** research lines:

1. a **risk classification** approach, as in Gschlößl, Schoenmaekers & Denuit (2011, EAJ), Chapter 2 of Cass Business School and Hymans Robertson LLP (2014) and van Berkum, Antonio & Vellekoop (2021, JRSS A)
2. a **relative approach**: (1) first model reference population, (2) then model the mortality dynamics in the book given the reference model, as discussed in Chapter 4-5-6 of IFoA/LLMA (2014).

In **research**, focus is on explaining observed mortality rates as accurate as possible:

- all individual observations are given the same weight in a regression model.

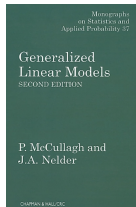
In **industry**, mortality rates are used for **valuation of the liabilities**:

- individuals with high insured amounts contribute more to the liabilities than individuals with low insured amounts.
- this has inspired actuaries and researchers to investigate the weighting of mortality rates by the insured amount.

Risk classification approach

Risk classification approach

Gschlössl, Schoenmakers & Denuit (2011, EAJ)



Key ideas:

- (1) **differentiation** based on multiple factors (e.g., age, gender, smoking habit, amount insured)
 - separate analysis per subportfolio not statistically meaningful, due to limited size of data
 - therefore: apply **regression** techniques \in statistical learning \in machine learning methods!
- (2) **Gschlössl et al. [2011]** use a **Poisson regression model** as a method for risk classification
 \Rightarrow sound **statistical framework** of **GLMs**
- (3) alternative: (at micro-level) use techniques from survival analysis (e.g., Cox proportional hazard, see **Richards [2008]**).

Risk classification approach

Gschlössl, Schoenmakers & Denuit (2011, EAJ)

9

- ▶ Mortality experience from insurance portfolio in the German marketplace over period of 5 years.

- ▶ Characteristics:

- male insured lives, with an individual policy, and ages 18-85
- total exposure-to-risk > 5 mio exposure years, > 12 000 deaths registered
- a set of covariates, e.g., product type, medical underwriting (yes or no), extra mortality, amount insured.

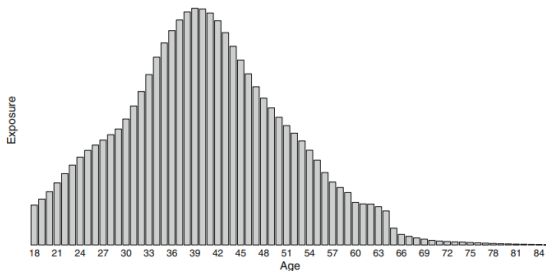


Fig. 1 Distribution of the central exposure-to-risk by attained age

► Exploratory analysis:

- (1) estimate **baseline mortality rates** by Poisson regression for death counts
disregard all other covariates at this stage

$$D_x \sim \text{Poi}(E_x \cdot \mu_x^b),$$

where D_x is **death counts** at age x , E_x corresponding **exposure** and

$$\log \mu_x^b = f(x)$$

for some smooth function $f(x)$.

Alternatively, a **population reference model** (e.g., KAG, IA|BE) can be used.

► Exploratory analysis:

(2) using the estimated baseline mortality

$$\hat{\mu}_x^b = \exp(\hat{f}(x))$$

one can explore the Standardized Mortality Rates (SMRs) for different subgroups or risk profiles of insured lives.

We construct

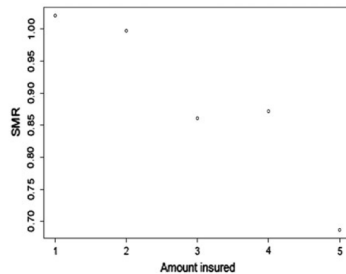
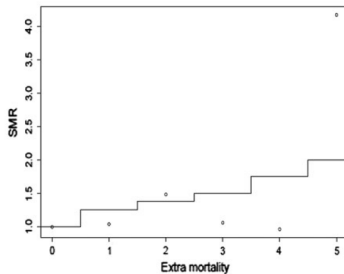
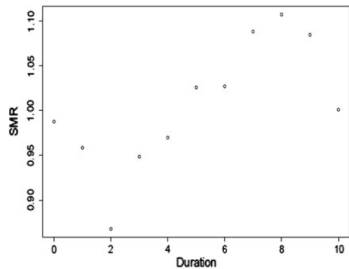
$$\text{SMR}_x = \frac{d_x}{\hat{d}_x} = \frac{d_x}{er_x \cdot \hat{\mu}_x^b} \quad \text{and} \quad \text{SMR} = \frac{\sum_{x=18}^{85} d_x}{\sum_{x=18}^{85} \hat{d}_x} = \frac{\sum_{x=18}^{85} d_x}{\sum_{x=18}^{85} er_x \cdot \hat{\mu}_x^b}.$$

Plot these SMRs against the covariates for an exploratory analysis of covariate effects.

Risk classification approach

12

Gschlössl, Schoenmakers & Denuit (2011, EAJ) - SMR



Risk classification approach

Gschlössl, Schoenmakers & Denuit (2011, EAJ) - Poisson regression

- Specify a 'risk cell' ℓ as a unique combination of covariates.
- Apply Poisson regression

$$D_{\ell} \sim \text{Poi}(E_{\ell} \cdot \mu_{\ell})$$

$$\ln \mu_{\ell} = \beta_0 + \beta_1 \ln \mu_{\ell}^b + \sum_j \beta_j x_{\ell,j},$$

where:

- E_{ℓ} is the exposure-to-risk for cell ℓ
- μ_{ℓ}^b is the baseline mortality for cell ℓ , a function of age (or, age, gender and period) only
- $x_{\ell,j}$'s are binary variables coding the covariates for cell ℓ .

- ▶ Perform **model building** and **variable selection** using tools for GLMs.
- ▶ End product:
 - identification and quantification of statistically significant risk factors
 - confidence intervals for estimated effects
 - mortality rates for specific **risk profiles**.

Risk classification approach

15

Gschlössl, Schoenmakers & Denuit (2011, EAJ) - Poisson regression end product

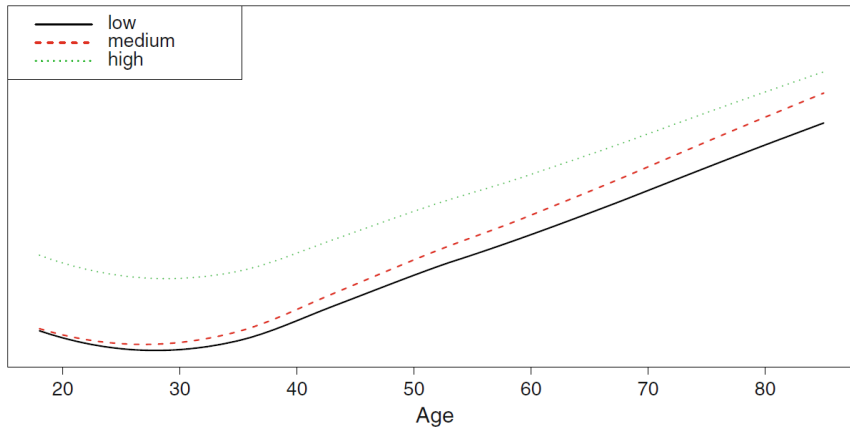


Fig. 4 Mortality rates according by attained aged and risk profile

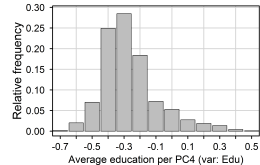
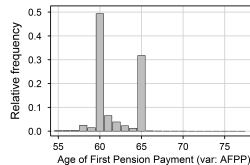
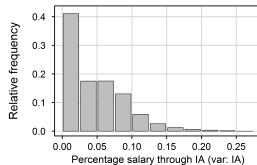
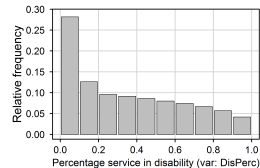
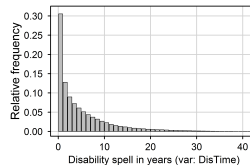
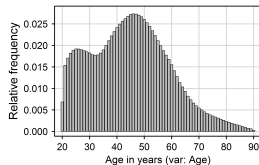
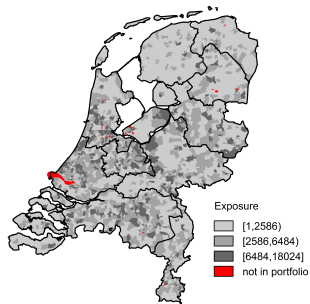
van Berkum, Antonio & Vellekoop (2021, JRSS A)

- ▶ We discuss the van Berkum, Antonio & Vellekoop (2021, JRSS A) paper on **Quantifying longevity gaps using micro-level lifetime data**:
 - data from a Dutch pension fund that follows individuals during the period 2006 to 2011
 - 11 325 511 individual observations on 2 162 899 unique individuals resulting in a total of 11 304 448 years lived, and during the observed period 41 622 deaths were recorded
 - risk factors available, e.g., salary and disability information
 - GAMs with KAG model for $\mu_{x,t}$ (\sim Li & Lee model) as baseline.

Risk classification approach

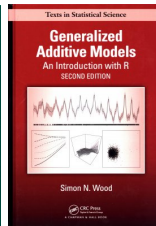
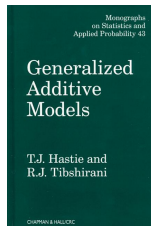
van Berkum, Antonio & Vellekoop (2021, JRSS A)

17



Risk classification approach

van Berkum, Antonio & Vellekoop (2021, JRSS A)



► Generalized Linear Models (GLMs):

- transformation of the mean modelled by a **linear predictor** (say $\mathbf{x}_i' \boldsymbol{\beta}$)
- not well suited for continuous risk factors that relate to the response in a **non-linear** way.

► Generalized Additive Models (GAMs):

- allow for **smooth effects** of continuous risk factors, including a spatial smoother, in the predictor.

van Berkum, Antonio & Vellekoop (2021, JRSS A)

- Generalized Additive Model with predictor:

$$\eta_{\ell} = \ln(\mu_{\ell}) = \beta_0 + \ln \mu_{\ell}^b + \sum_{j=1}^p \beta_j x_{\ell j}^d + \sum_{j=1}^q f_j(x_{\ell j}^c) + \sum_{j=1}^r f_j(x_{\ell j}^s, y_{\ell j}^s),$$

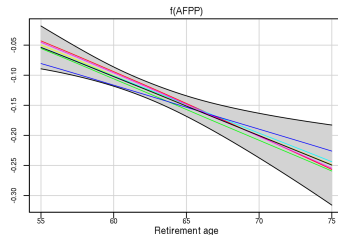
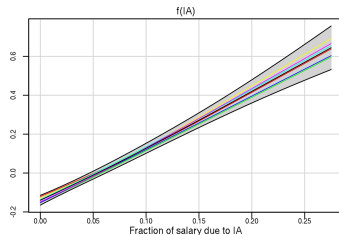
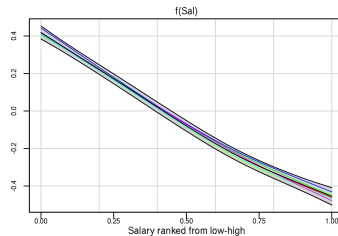
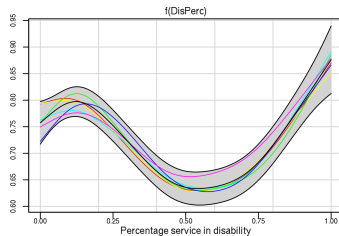
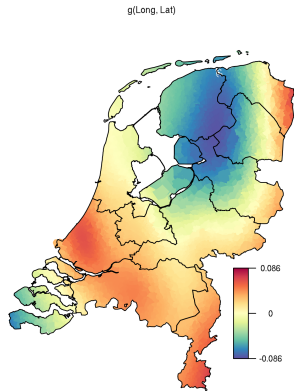
where μ_{ℓ} is the mean death count per unit of exposure for risk cell ℓ .

- The predictor includes:
- dummy-coded factor variables
 - smooth function $f_j(\cdot)$ of a one-dimensional continuous covariate X_j^c
 - smooth function $f_j(\cdot, \cdot)$ of a two-dimensional continuous variables (X_j^s, Y_j^s) (e.g., spatial effect or interaction effect of two continuous covariates).

Risk classification approach

van Berkum, Antonio & Vellekoop (2021, JRSS A)

20



Risk classification approach

van Berkum, Antonio & Vellekoop (2021, JRSS A)

- ▶ Pension funds want accurate estimates of their liabilities; we developed a **financial backtest** to test the appropriateness of the outcomes
- ▶ Define b_j as the insured amount and a_j the annuity factor for participant j
- ▶ We compared EoY expected liabilities:

$$\Gamma = \sum_{j=1}^{L_t} (Y_{t,j} \cdot b_j a_j + (1 - Y_{t,j}) \cdot 0)$$

vs realized liabilities $\tilde{\Gamma} = \sum_{j=1}^{L_t} I_j \cdot b_j a_j$.

- ▶ We compute mean and variance for the random liabilities at YE:

$$\mathbb{E}(\Gamma|\eta_j) = \sum_{j=1}^{L_t} p_{t,j} \cdot b_j a_j$$

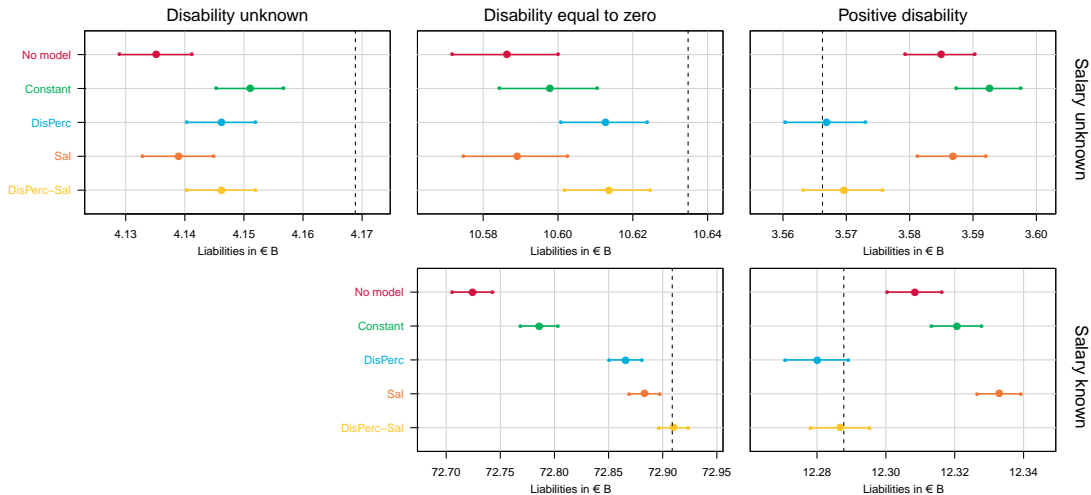
$$\text{Var}(\Gamma|\eta_j) = \sum_{j=1}^{L_t} (b_j a_j)^2 \cdot p_{t,j} \cdot (1 - p_{t,j})$$

- ▶ Assuming normality for the liabilities, we construct prediction intervals for $\Gamma \rightarrow$

Risk classification approach

van Berkum, Antonio & Vellekoop (2021, JRSS A)

22



So far, we discussed GLMs and GAMs to build a model for

$$D_{x,t,g,\ell} \sim \text{POI}(E_{x,t,g,\ell} \cdot \mu_{x,t,g}^{\text{pop}} \cdot \exp(f(\mathbf{x}_\ell)),$$

where ℓ refers to a risk cell, a unique combination of covariates, and $\mu_{x,t,g}^{\text{pop}}$ captures the age-period-gender baseline mortality rate.

In non-life insurance, the use of **machine learning methods** is (nowadays) gaining popularity as an alternative for POI regression models.

Why?

- avoid manual model building, helps with covariate selection and covariate engineering (e.g., binning)
- more automatic \Rightarrow useful with lots of covariates!
- scales better to large data sets
- useful as a first screening, then build GLM as (global) surrogate model
- ...

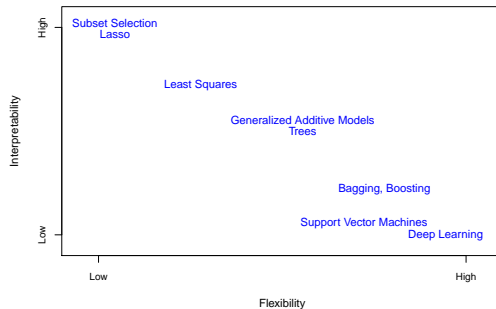
Risk classification approach

From statistical to machine learning

24

A few pointers:

- Henckaerts et al. (2021, NAAJ) on tree-based learners for claim count and severity data
- Schelldorfer & Wüthrich (2019) on Combined Actuarial Neural Networks (CANNs) + papers and lecture notes by Wüthrich
- Henckaerts et al. (2022) on GLM as a global surrogate for, e.g., a Gradient Boosting Machine (GBM).



Picture from James et al., 2021, An introduction to statistical learning

Relative approach

- ▶ The **relative approach** proposed in **BAJ [2015]** and (IFoA/LLMA, 2014) models
 - the reference population
 - then the book **given the reference**.
- ▶ Advantages? (BAJ, 2015) and (IFoA/LLMA, 2014)
 - allows data mismatch reference vs. book
 - reference population models widely available and extensively studied
 - reference typically much larger than book
 - consistency when modelling several books using same reference.

► Step 1: fit a model for the reference

general model for the reference population (\sim Binomial distributional assumption used in IFoA/LLMA, 2014)

$$D_{xt}^R \sim \text{BIN}(E_{xt}^R, q_{xt}^R)$$
$$\text{logit}(q_{xt}^R) = \log\left(\frac{q_{xt}^R}{1 - q_{xt}^R}\right) = \alpha_x^R + \sum_{j=1}^N \beta_x^{(j,R)} \cdot \kappa_t^{(j,R)} + \gamma_c^R$$

where

- each $\kappa_t^{(j,R)}$ contributes to the reference mortality trend
- γ_c^R is a cohort effect in the reference, with $c = t - x$.

- Step 2: fit the book given the reference

$$D_{xt}^B \sim \text{BIN}(E_{xt}^B, q_{xt}^B)$$
$$\text{logit}(q_{xt}^B) - \text{logit}(q_{xt}^R) = \alpha_x^B + \sum_{j=1}^M \beta_x^{(j,B)} \kappa_t^{(j,B)} + \gamma_c^B$$

where

- α_x^B captures mortality level differences between reference and book
- M components contribute to trend in differences in mortality and γ_c^B captures differences in cohort effect.

- ▶ **Step 3:** specify and calibrate **time series dynamics** using ARIMA toolbox (\sim our discussions in Modules 1-3)

(suggestions from the IFoA/LLMA, 2014 paper, section 5.2.3)

- **multivariate random walk with drift** for the reference time indices κ_t^R
 - **integrated AR(1)** for cohort effect γ_c^R in reference
 - **vector AR(1)** process for time indices in the book
 - **AR(1) process** for cohort γ_c^B .
- ▶ Select appropriate models for reference and book, examine goodness-of-fit, robustness, ...

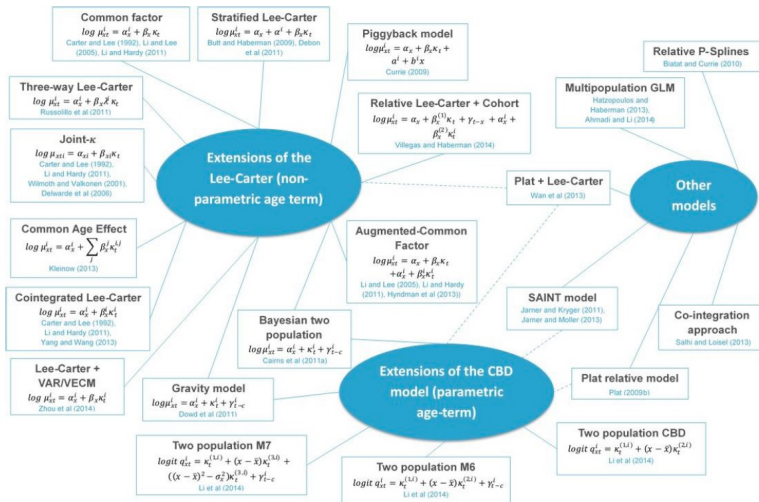


Figure 5.1: Universe of multi-population models

- ▶ The authors reflect on 3 sources of uncertainty when producing fan charts:
 1. process risk: from the time series
 2. parameter uncertainty: from estimation of parameters, including those in the time series; bootstrapping is used!
 3. sampling risk: random sampling the number of deaths under the assumed binomial model.
- ▶ Particularly relevant with smaller books!

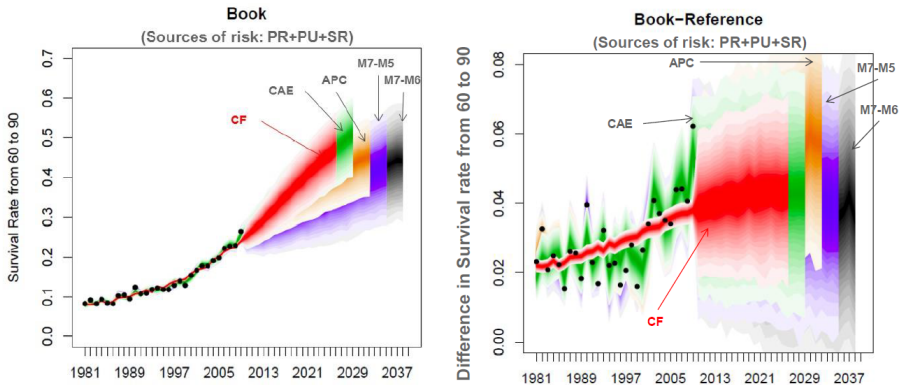


Figure 6.10: Fan charts of 30 year period survival probabilities at age 60 for the “Extreme Wealthy” test book using different mortality models and different sources of risk (PR=process risk; PU=parameter uncertainty; SR=sampling risk). Left panes present results for the book population and right panes results for the difference in survival probabilities in the book and the reference population.

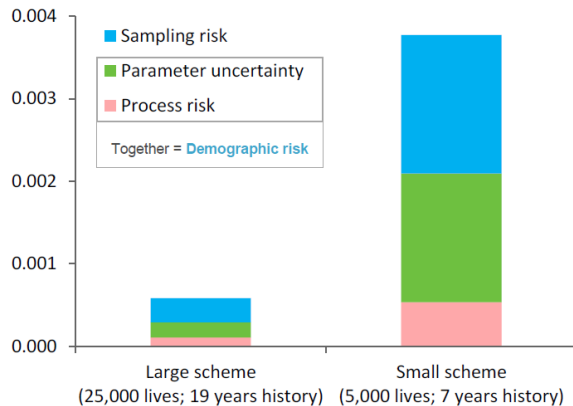


Figure 9.1 Comparison of variance of difference between book and reference population under CAE+cohort for two different book populations.

Industry approach

Industry approach

Use of mortality rates for valuation purposes

- ▶ For cash flow projection systems, it is convenient if the dimensions of the mortality rates are stable over time, so only few risk factors are used to differentiate between policyholders.
- ▶ Insurance companies typically set longevity assumptions / mortality rates by **Homogeneous Risk Group** (HRG); a group of policies with similar policy characteristics.
- ▶ For each HRG, mortality rates are then defined by age, sex and calendar year (sometimes also more granular, e.g., by policy duration or smoking status).
- ▶ **Gschlössl et al. [2011]** discuss that some actuaries use the insured amounts as weight, but the authors consider that an 'ad-hoc' approach and prefer to use insured amounts as a risk factor.

Industry approach

Two definitions of death rates

Dependence on sex is omitted, but the analyses are typically performed for the two sexes separately.

Define $d_{x,t}^{\text{pf,num}}$ and $E_{x,t}^{\text{pf,num}}$ as the observed number of deaths and corresponding exposure-to-risk in year t and at age x :

- These observations are obtained by simply counting **numbers of individuals**.
- We define the crude death rate $m_{x,t}^{\text{pf,num}} = d_{x,t}^{\text{pf,num}} / E_{x,t}^{\text{pf,num}}$.
- All individuals are given the same weight!

Define $d_{x,t}^{\text{pf,ia}}$ and $E_{x,t}^{\text{pf,ia}}$ as the release in insured amounts and total insured amount

- These observations are obtained by counting the **insured amounts** of the individuals.
- We define the crude death rate $m_{x,t}^{\text{pf,ia}} = d_{x,t}^{\text{pf,ia}} / E_{x,t}^{\text{pf,ia}}$.
- People with higher insured amount contribute more to $m_{x,t}^{\text{pf,ia}}$!

Industry approach

Two definitions of death rates

- ▶ Numbers weighted mortality can easily be analyzed through Poisson or (Negative) Binomial regression on the observed death counts $d_{x,t}^{\text{pf,num}}$, see the [Risk classification approach](#)
- ▶ For insured amounts weighted mortality, there is no obvious probability distribution. [Plat \[2009\]](#) shows (among other things) the approach often taken in (Dutch) industry:
 - define 'observed' experience factors $\theta_{x,t} = m_{x,t}^{\text{pf,ia}} / m_{x,t}^{\text{pop}}$
 - WLS is applied to these experience factors, for example as:

$$\theta_{x,t} = \beta_0 + \beta_1 x + \varepsilon_{x,t} \quad \varepsilon_{x,t} \sim \text{Normal}\left(0, \frac{\sigma^2}{w_{x,t}}\right),$$

where $w_{x,t}$ is a weight that is introduced to capture the heterogeneity in the observations.

- ▶ HRG-specific assumptions are defined as $q_{x,t}^{\text{pf}} = q_{x,t}^{\text{pop}} \cdot \hat{\theta}_x$ with $\hat{\theta}_x = \hat{\beta}_1 + \hat{\beta}_1 x$.

- ▶ **Richards [2008]** includes the insured amount as a weight in a regression approach and finds that this approach better captures the release in provision / sum insured.
- ▶ This approach is what **Gschlössl et al. [2011]** argue to be 'ad-hoc', but it enables including more risk factors.
- ▶ Recent work on the topic performed at **RCLR** will be available soon.

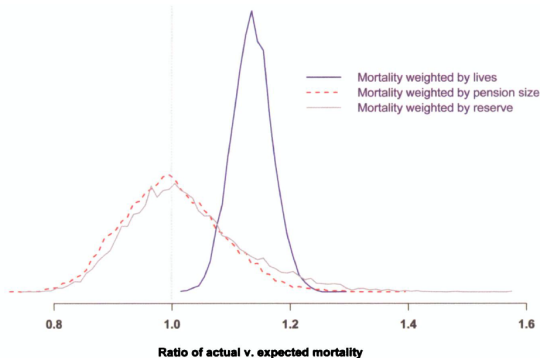


Figure 8. Frequency plot of ratio of actual v. expected mortality for 10,000 bootstrapped portfolios of 50,000 lives (amounts-weighted fit)

Source: **Richards [2008]**

Wrap-up

By now, you should be able to:

- discuss the main characteristics of the relative approach and the risk classification approach to model portfolio mortality.
- understand how different types of risk may affect future mortality rates.
- explain how portfolio mortality can be backtested, also in financial terms.
- understand why in industry sometimes more ad-hoc approaches are taken.

- A methodology for assessing basis risk - abstract of the london discussion. *British Actuarial Journal*, 20(3):461–490, 2015.
- Susanne Gschlössl, Pascal Schoenmaekers, and Michel Denuit. Risk classification in life insurance: methodology and case study. *European Actuarial Journal*, 1:23–41, 2011.
- Richard Plat. Stochastic portfolio specific mortality and the quantification of mortality basis risk. *Insurance: Mathematics and Economics*, 45(1):123–132, 2009.
- Stephen J Richards. Applying survival models to pensioner mortality data. *British Actuarial Journal*, 14(2):257–303, 2008.