

# Mortality forecasting using variational inference

Patrik Andersson

*Uppsala University*

Mathias Lindholm

*Stockholm University*

---

## Abstract

This template helps you to create a properly formatted

*Keywords:* Non-linear state-space-models, Mortality forecasting, Hidden Markov model, Variational inference

---

## 1. Introduction

Attempts to forecast mortality go back perhaps as far as [1]. A more recent example is the Lee-Carter model [2] and its extensions, see [3, 4, 5] for a survey. Applications of mortality forecasting can be found in for example demographic predictions and in the insurance industry.

The Lee-Carter model is a log-linear multivariate Gaussian model of mortality rates. A major criticism of Lee-Carter-type models is that the model training is done as a two-step process. In the first step, point estimates of the mortality rates are obtained, for example as the maximum likelihood estimate of a Poisson distribution, and in the second step, a latent process is fitted to these estimates. This method has the advantage that it is simple and fast to implement, but it is inefficient when compared to simultaneous estimation of all unknown parameters. Also, it is not possible to distinguish between the finite population noise of the mortality estimates and the noise from the latent process. Both of these issues can potentially affect the quality of the forecasts.

Simultaneous estimation of parameters has been considered in [6] where particle filtering methods are used for a Poisson Lee-Carter model similar to [7]. However, this method has its drawbacks. It could be considered cumbersome for practitioners as it requires custom implementation and tuning, and

since the particle filter methods are computationally expensive, the number of parameters must not be too large. The complexity of the modelling, is, of course, reduced when removing the Poisson assumption, see e.g. [8] for a simpler state-space model treatment of a standard Gaussian Lee-Carter model.

Recently it has been suggested to use models from deep learning, sometimes called deep factor models, to forecast high-dimensional multivariate time series. Some examples of this can be found in [9, 10, 11, 12]. The applications presented in those articles differ from mortality forecasting in the scale of the problem. In mortality forecasting, the dimension of the time series is about 100 (the lifetime in years of a human) and the number of observations of time series is also about 100 (although some countries do have reliable data for longer than that). As a consequence, to avoid overfitting, we need to consider simpler models. This includes shallower networks for mapping latent variables to the observed time series, linear Gaussian models instead of RNNs for propagating the latent variables forward in time and fewer latent factors.

Compared to previous mortality forecasting models, the novelty of this paper is therefore to use black-box variational inference (VI) [13] to solve the inference problem. This means that after specifying how to sample from the model and the approximate posterior, the inference is done automatically without any model-specific customisation. The family of models that can be handled is also expanded. For example, one can consider neural networks for mapping the latent process to mortality rates. Also, in this case, all the parameters can be estimated simultaneously. This latter point is problematic when using particle filter techniques, where it is necessary to estimate the linear mapping from the latent process to the mortality rates in isolation before continuing with the estimation of the other parameters.

Other approaches that use machine learning techniques for forecasting mortality rates can be found in e.g. [Wüthrich & co] that consider various types of Gaussian recurrent neural network structures, [Italienarna, Lindholm & Palmborg] that consider univariate LSTM NNs, both with and without a Poisson population assumption, and see [xx] that consider tree-based techniques.

The model is implemented using the probabilistic programming language Pyro [14] and the code is available on Github, see [reference].

The rest of the paper is organised as follows: In Section 2 we describe the probabilistic model that will be used for forecasting. In Section 3 we give a

brief introduction to variational inference and in Section 4 we describe how to  
60 forecast the mortality once the model has been trained and how we validate  
the forecast. In Section 5 we demonstrate our method on a few examples.  
Section 6 concludes the paper.

## 2. Model

In this section, we describe the probabilistic model that defines the mor-  
65 tality dynamics. Uppercase letters will denote random variables and low-  
ercase letters the corresponding observed value. Greek letters will denote  
unknown parameters that are to be estimated.

Mortality data can be aggregated in different forms and clearly the choice  
of model will have to be adjusted accordingly. For example, the data could  
70 contain the population size of each age group at the beginning of the year  
and the number of deaths during the year. In this case, a binomial model  
seems natural. We however consider data on the yearly number of deaths  
and the *exposure to risk* in each age group. The exposure to risk in this  
setting is the total time that the individuals in the population were a certain  
75 age in a certain year. The number of deaths in age  $a \in \{0, 1, \dots, \bar{a}\}$ , year  
 $t \in \{0, 1, \dots, \bar{t}\}$  is denoted by  $D_{a,t}$  and the exposure by  $E_{a,t}$ .

Our model is a dynamic factor model that can be written as:

$$D_{a,t} \mid X_t, E_{a,t} \sim \text{Poisson} \left( E_{a,t} \exp \left( f_a^\psi(X_t) \right) \right), \quad (1)$$

$$X_{i,t} = X_{i,t-1} + K_{i,t-1} + U_{i,t}, \quad U_{i,t} \text{ iid } \mathbf{N}(0, \sigma_{X,i}^2), \quad (2)$$

$$K_{i,t} = \mu_i + \varphi_i(K_{i,t-1} - \mu_i) + V_{i,t}, \quad V_{i,t} \text{ iid } \mathbf{N}(0, \sigma_{K,i}^2). \quad (3)$$

Here  $i = 1, 2, \dots, d$ , where  $d$  is the dimension of the latent variables. We  
also require  $0 \leq \varphi_i \leq 1$ . The function  $f_a^\psi$  is the  $a$ :th component of  $f^\psi : \mathbb{R}^d \rightarrow \mathbb{R}^{\bar{a}+1}$ . In our examples in Section 5  $f^\psi$  will be given by either an  
affine transformation or a sum of radial basis functions. That is either,  
 $f^\psi(x) = Ax + b$ , where  $A$  and  $b$  are trainable or  $f^\psi$  is a sum of radial basis  
functions, that is, the  $a$ th component is

$$f_a^\psi(x) = x^T \sum_{i=1}^p w_i e^{-\tau^2(a-\mu_i)^2} + b.$$

We will fix  $\tau = 10$  and therefore  $\{w_i, \mu_i, b\}_{i=1}^p$  are the trainable parameters.  
Compared to the more general affine transformation, radial basis functions

have the advantage of guaranteeing a certain smoothness of  $f_a$  as a function of  $a$ , encoding a prior that similar ages should have similar mortality.

We remark here that the exact specification of the above model is not critical for the continuation. For example, the exponential link function in the Poisson distribution could be changed to some other positive differentiable function without complication. We are assuming that the components of the latent process are independent, instead, we let any dependence be captured by  $f$ . However, this latent process could be replaced with some other Markov process.

### 3. Variational inference

Here we explain the main ideas of variational inference (VI) in a general setting. At the end of the section, we connect this to our specific model. For more on VI in general we refer to [13] and for the application to state space models, see [15].

We are observing  $y$ , whose distribution depends on a latent variable  $x$  and an unknown parameter  $\psi$ . This is modelled by the joint distribution

$$p_\psi(y, x) = p_\psi(y \mid x)p_\psi(x). \quad (4)$$

The likelihood,

$$L(\psi) = p_\psi(y) = \int p_\psi(y, x)dx, \quad (5)$$

is in general not tractable and therefore approximations are needed in order to be able to estimate  $\psi$ . Consider a parametrised distribution, the approximate posterior,  $q_\theta(x)$ . Then observe that, due to Jensen's inequality, the log-likelihood is

$$\begin{aligned} l(\psi) &:= \log L(\psi) = \log \int \frac{p_\psi(y, x)}{q_\theta(x)} q_\theta(x) dx \\ &\geq \int (\log p_\psi(y, x) - \log q_\theta(x)) q_\theta(x) dx =: \mathcal{L}(\psi, \theta). \end{aligned}$$

The right-hand side is known as the evidence lower bound (ELBO). The idea of VI is to instead of maximising the log-likelihood, maximise the ELBO.

Towards this, we calculate the gradients

$$\partial_\psi \mathcal{L}(\varphi, \theta) = \int \partial_\psi \log p_\psi(y, x) q_\theta(x) dx, \quad (6)$$

$$\partial_\theta \mathcal{L}(\psi, \theta) = \int (\log p_\psi(y, x) - \log q_\theta(x)) \partial_\theta \log q_\theta(x) q_\theta(x) dx. \quad (7)$$

We can then proceed to obtain unbiased estimates of the gradients by sampling from  $q_\theta$  and maximise  $\mathcal{L}$  using stochastic optimisation algorithms. 95 Once converged,  $q_\theta(x)$  can be used as an approximation of the posterior distribution of the latent variables  $p_\psi(x | y)$ .

Further, to obtain faster convergence, various variance reduction techniques are often used. Here we only mention the so-called reparametrisation trick. Suppose that we can find functions  $x_\theta$  such that

$$\int f(x) q_\theta(x) dx = \int f(x_\theta(z)) q(z) dz, \quad (8)$$

which makes the sampling distribution independent of  $\theta$ . In particular, the gradient satisfies

$$\partial_\theta \int f(x) q_\theta(x) dx = \partial_\theta \int f(x_\theta(z)) q(z) dz = \int \partial_\theta f(x_\theta(z)) q(z) dz, \quad (9)$$

which usually improves the sampling variance compared to differentiating the density directly. An important example of a distribution that allows for reparametrisation according to (8) is the Gaussian, since if  $Z \sim \mathbf{N}(0, 1)$  then 100  $\mu + \sigma Z \sim \mathbf{N}(\mu, \sigma^2)$ .

In the numerical illustrations in Section 5, the approximate posterior is modelled as a Gaussian distribution with an autoregressive covariance. That is, the distribution of the process is given by

$$\tilde{X}_{i,t} = \tilde{\mu}_{i,t}^X + \alpha_{i,t} \tilde{X}_{i,t-1} + \tilde{e}_{i,t}^X, \quad \tilde{e}_{i,t}^X \text{ iid } \mathbf{N}(0, \tilde{\sigma}_{X,i}^2), \quad (10)$$

$$\tilde{K}_{i,t} = \tilde{\mu}_{i,t}^K + \beta_{i,t} \tilde{K}_{i,t-1} + \rho_{i,t} \tilde{X}_{i,t-1} + \tilde{e}_{i,t}^K, \quad \tilde{e}_{i,t}^K \text{ iid } \mathbf{N}(0, \tilde{\sigma}_{K,i}^2). \quad (11)$$

#### 4. Forecasting and validation

Here we discuss how to forecast the mortality after maximising the ELBO and thus obtaining estimates of  $\varphi$ ,  $\theta$  and, in particular, the joint distribution of  $(\tilde{X}_{i,\bar{t}}, \tilde{K}_{i,\bar{t}})$ .

Since both the approximate posterior and the latent process is Gaussian, the forecasting distribution of the latent process is also Gaussian. That is, for  $t > \bar{t}$ ,

$$\begin{pmatrix} \hat{X}_{i,t} \\ \hat{K}_{i,t} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \hat{\mu}_{i,t}^X \\ \hat{\mu}_{i,t}^K \end{pmatrix}, \begin{pmatrix} \hat{\sigma}_{X,t}^2 & \hat{\rho}_{i,t} \hat{\sigma}_{X,t} \hat{\sigma}_{K,t} \\ \hat{\rho}_{i,t} \hat{\sigma}_{X,t} \hat{\sigma}_{K,t} & \hat{\sigma}_{K,t}^2 \end{pmatrix} \right), \quad (12)$$

where the parameters can be calculated iteratively from (2) and (3), by using the initial value

$$\begin{pmatrix} \hat{X}_{i,\bar{t}} \\ \hat{K}_{i,\bar{t}} \end{pmatrix} := \begin{pmatrix} \tilde{X}_{i,\bar{t}} \\ \tilde{K}_{i,\bar{t}} \end{pmatrix}, \quad (13)$$

105 and the forecast of mortality rates is given by  $\exp(f^{\hat{\psi}}(\hat{X}_t))$ .

If one wants to forecast the actual number of deaths, a forecast of the number of living at the beginning of the year is also needed, together with some assumption on the distribution of when in the year people are born. For a longer discussion on how this can be done, we refer to [6].

110 The forecast is validated by calculating the logarithmic score of the forecast on the validation data set, see for example [16]. That is, the mortality rates are forecasted and multiplied by the observed exposure-to-risk in each age group, giving the intensity of the Poisson distributed number of deaths in each age group. The logarithmic score of each age group is summed to  
115 give the total logarithmic score for a given year.

## 5. Results

In this section, we illustrate the performance of the model by fitting it using a dataset on the mortality of Swedish males. The dataset is collected from [17]. We evaluate the models using a rolling window; using 60 years  
120 to fit the model and 10 years to evaluate the forecast against the actual outcome. The training and evaluation windows are then rolled forward one year and the process repeats. For comparison, we also fit the Lee-Carter and the model by Plat [18].

The model is as in (1) - (3) where  $f$  is either affine or a sum of radial  
125 basis functions. We begin by selecting hyperparameters for each model, e.g. dimension of the affine transformation or the number of radial bases. We then illustrate the fitted model and the out-of-sample forecasting performance. Finally, we compare our model with some alternatives.

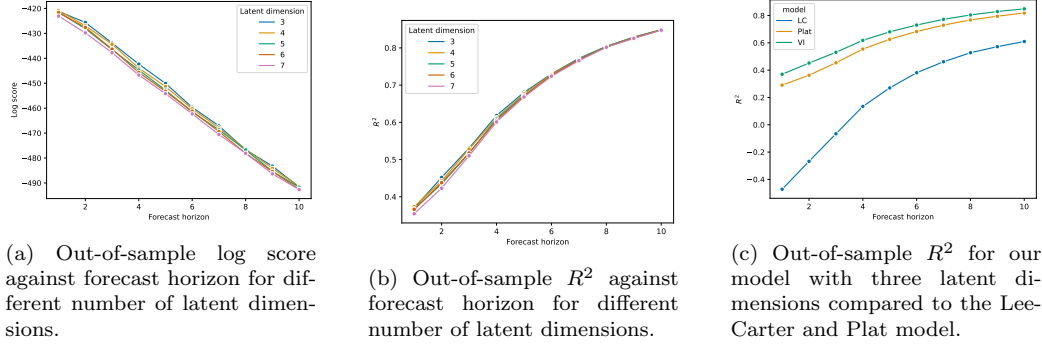


Figure 1: Model evaluation when fitted on Swedish male mortality data with the first year from 1920 to 1952. Each model is fitted using 60 years of data and evaluated on the following 10 years. The figures show that three latent dimensions give the best out-of-sample performance in our model and that our model outperforms the comparison models.

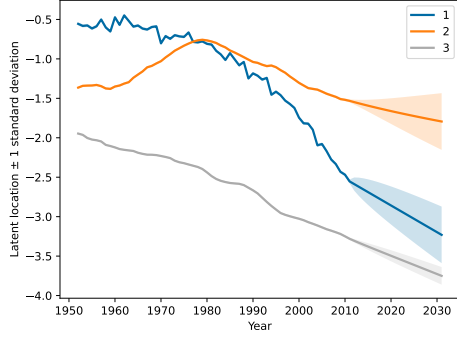
### 5.1. Affine

For model selection, we compare the out-of-sample log-score for varying dimensions of the latent process. In Figure 1a and 1b we see that dimension 3 overall performs best. We have performed experiments also for dimensions 1 and 2, but the performance was considerably worse, and they are therefore excluded from the picture. In Figure 1c we compare the  $R^2$  metric of our model with that of Lee-Carter and Plat. We see that the Lee-Carter model is not competitive, while the model by Plat is slightly worse than ours. Figure 2 illustrates the fitted model. In particular, we note that the in-sample fit of the mortality rates in 2d are quite good. Figure 3 shows the smoothed mortality rates and forecasted for four age groups from 2012 to 2023. The shaded regions represent  $\pm 1$  standard deviation and the grey dots are the observed mortality rates. That is, we should expect that around 7 out of 10 observations are within the shaded region. The pictures seem to confirm this.

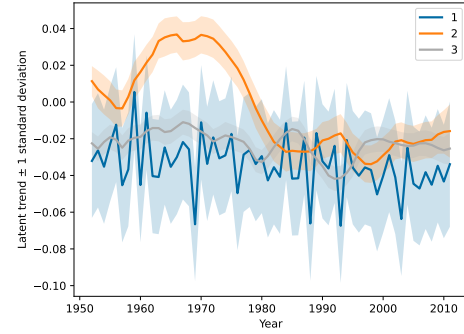
### 5.2. Radial basis functions

### 5.3. Comparison

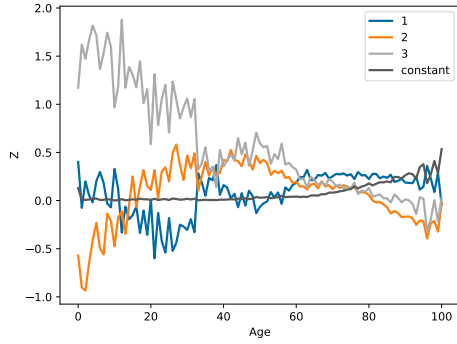
The model performance is summarized in Table 1 where the log-score and  $R^2$  are shown for each model, averaged over all the forecast horizons and rolling windows.



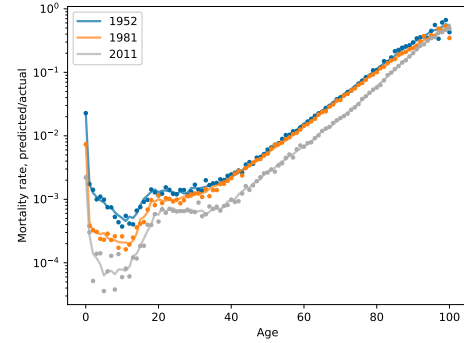
(a) Level of latent process in-sample and its forecast. Shaded region is  $\pm 1$  standard deviation. From 1950 to 2011 shows smoothed values, from 2012 shows forecast.



(b) Trend of latent process. The shaded region is  $\pm 1$  standard deviation.



(c) Factor loadings for the model with 3 latent dimensions



(d) In-sample fit of model mortality rates. Dots are observed rates and solid lines are the fitted model. This model is fitted on data from 1952 to 2011.

Figure 2: Illustration of the model fitted to Swedish data from 1952 to 2011 with three latent dimensions.



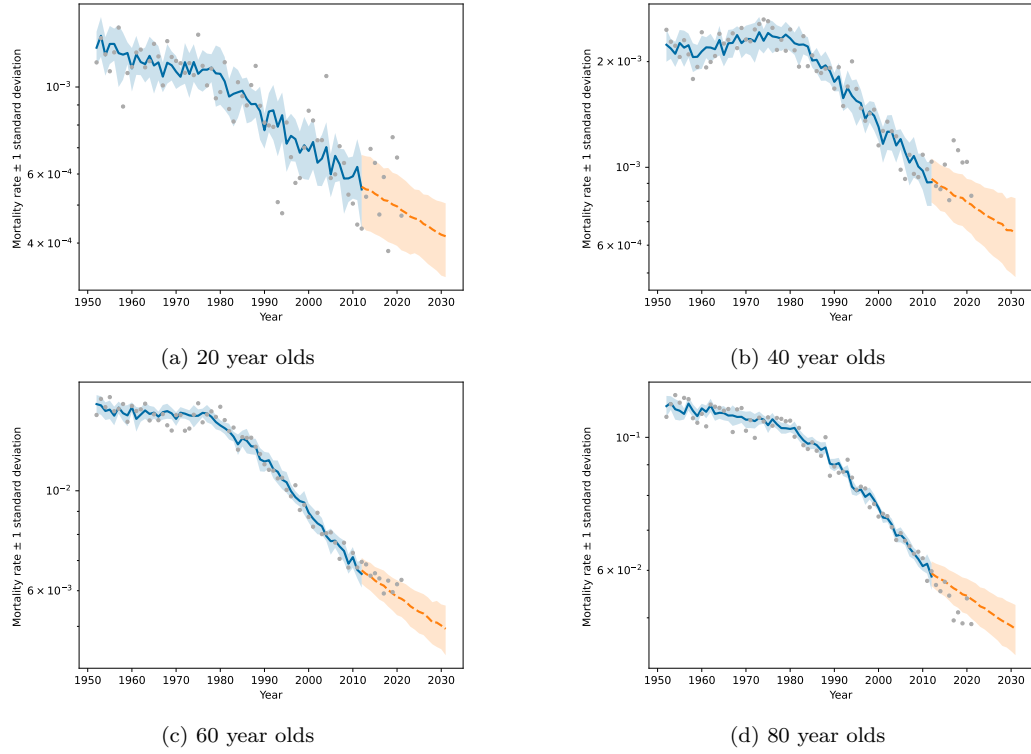


Figure 3: Mortality forecasts for four different age groups. Solid line is the smoothed mortality rate, dashed line is the mean forecasted mortality rate and shaded region is  $\pm 1$  standard deviation of the forecasted observed mortality rate. Dots indicate observed mortality rates.

Model	Log-score	$R^2$
Affine	-455.1	0.664
Lee-Carter	-608.5	0.215
Plat	-474.9	0.608

Table 1: Out-of-sample model evaluation metrics for Sweden. Numbers are averaged over all forecast horizons and rolling windows.

## 6. Conclusions

## References

- [1] B. Gompertz, 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, Philosophical transactions of the Royal Society of London (1825) 513–583.
- [2] R. D. Lee, L. R. Carter, Modeling and forecasting us mortality, Journal of the American statistical association 87 (419) (1992) 659–671.
- [3] H. Booth, L. Tickle, Mortality modelling and forecasting: A review of methods, Annals of actuarial science 3 (1-2) (2008) 3–43.
- [4] S. Haberman, A. Renshaw, A comparative study of parametric mortality projection models, Insurance: Mathematics and Economics 48 (1) (2011) 35–55.
- [5] M. F. Carfora, L. Cutillo, A. Orlando, A quantitative comparison of stochastic mortality models on italian population data, Computational Statistics & Data Analysis 112 (2017) 198–214.
- [6] P. Andersson, M. Lindholm, Mortality forecasting using a lexis based state space model (2020).
- [7] N. Brouhns, M. Denuit, J. K. Vermunt, A poisson log-bilinear regression approach to the construction of projected lifetables, Insurance: Mathematics and Economics 31 (3) (2002) 373–393.
- [8] P. De Jong, L. Tickle, Extending lee–carter mortality forecasting, Mathematical Population Studies 13 (1) (2006) 1–18.

- [9] N. Nguyen, B. Quanz, Temporal latent auto-encoder: A method for probabilistic multivariate time series forecasting, arXiv preprint arXiv:2101.10460 (2021).
- 175 [10] Y. Wang, A. Smola, D. Maddix, J. Gasthaus, D. Foster, T. Januschowski, Deep factors for forecasting, in: International Conference on Machine Learning, PMLR, 2019, pp. 6607–6617.
- [11] D. Salinas, V. Flunkert, J. Gasthaus, T. Januschowski, Deepar: Probabilistic forecasting with autoregressive recurrent networks, International  
180 Journal of Forecasting 36 (3) (2020) 1181–1191.
- [12] S. S. Rangapuram, M. W. Seeger, J. Gasthaus, L. Stella, Y. Wang, T. Januschowski, Deep state space models for time series forecasting, Advances in neural information processing systems 31 (2018) 7785–7794.
- [13] R. Ranganath, S. Gerrish, D. M. Blei, Black box variational inference.,  
185 in: Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics, 2014.
- [14] E. Bingham, J. P. Chen, M. Jankowiak, F. Obermeyer, N. Pradhan, T. Karaletsos, R. Singh, P. Szerlip, P. Horsfall, N. D. Goodman, Pyro: Deep Universal Probabilistic Programming, Journal of Machine Learning  
190 Research (2018).
- [15] E. Archer, I. M. Park, L. Buesing, J. Cunningham, L. Paninski, Black box variational inference for state space models, arXiv preprint arXiv:1511.07367 (2015).
- [16] T. Gneiting, A. E. Raftery, Strictly proper scoring rules, prediction, and  
195 estimation, Journal of the American statistical Association 102 (477) (2007) 359–378.
- [17] Human Mortality Database. University of California, Berkeley (USA), and Max Planck Institute for Demographic Research (Germany), Available at <http://www.mortality.org> (downloaded on January 22, 2022)  
200 (2022).
- [18] R. Plat, On stochastic mortality modeling, Insurance: Mathematics and Economics 45 (3) (2009) 393–404.