

# Forecasting vital rates from demographic summary measures

Carlo G. Camarda<sup>1</sup>

<sup>1</sup> Institut national d'études démographiques (INED), Paris, France

E-mail for correspondence: `carlo-giovanni.camarda@ined.fr`

**Abstract:** In population and actuarial sciences, time-trends of summary measures (such as life expectancy or the average number of children per woman) are easy to interpret and predict. Most summary measures are nonlinear functions of the vital rates, the key variable we usually want to estimate and forecast. Furthermore smooth outcomes of future age-specific vital rates are desirable. Therefore, optimization with nonlinear constraints in a smoothing setting is necessary. We propose a methodology that combines Sequential Quadratic Programming and a *P*-spline approach, allowing to forecast age-specific vital rates when future values of demographic summary measures are provided. We provide an application of the model on Italian mortality and Spanish fertility data.

**Keywords:** Vital rates forecast; Smoothing; Constrained nonlinear optimization; Summary measures.

## 1 Introduction

Future mortality and fertility levels can be predicted either by modelling and extrapolating rates over age and time, or by forecasting summary measures, later converted into age-specific rates. The latter approach takes advantage of the prior knowledge that demographers and actuaries have on possible future values of measures such as life expectancy at birth and total fertility rate. Among others, this methodology has been lately adopted by the United Nations (Ševčíková et al., 2016). In this paper, we propose a model to derive future mortality and fertility age-patterns complying with projected summary measures. Unlike comparable approaches, we assume only smoothness of future vital rates, which is achieved by a two-dimensional *P*-spline approach as in Currie et al. (2004), and we allow constraints to multiple series of summary measures. Since these measures are

---

This paper was published as a part of the proceedings of the 34th International Workshop on Statistical Modelling (IWSM), Guimarães, Portugal, 7–12 July 2019. The copyright remains with the author(s). Permission to reproduce or extract any parts of this abstract should be requested from the author(s).

commonly nonlinear functions of the estimated penalized coefficients, Lagrangian multipliers cannot be directly implemented. We hence opted for a Sequential Quadratic Programming (SQP) procedure (Nocedal & Wright, 2006) to perform the associated constrained nonlinear optimization. We illustrate our approach with two data sets. We forecast mortality of Italian females, based on future life expectancy predicted by UN World Population Prospects (2017) and a future trend of a lifespan disparity measure obtained by time-series analysis. We also forecast Spanish fertility constrained to future values of total fertility rates, mean and variance of age at childbearing, derived by time-series analysis.

## 2 Model on Italian mortality data

For ease of presentation, we formulate the model on mortality data. We suppose that we have deaths, and exposures to risk, arranged in two matrices,  $\mathbf{Y} = (y_{ij})$  and  $\mathbf{E} = (e_{ij})$ , each  $m \times n_1$ , whose rows and columns are classified by age at death,  $\mathbf{a}$ ,  $m \times 1$ , and year of death,  $\mathbf{t}_1$ ,  $n_1 \times 1$ , respectively. We assume that the number of deaths  $y_{ij}$  at age  $i$  in year  $j$  is Poisson distributed with mean  $\mu_{ij} e_{ij}$ . Forecasting aims to reconstruct trends in  $\mu_{ij}$  for  $n_2$  future years,  $\mathbf{y}_2$ ,  $n_2 \times 1$ .

It is common practice to summarize mortality age-patterns by computing measures such as life expectancy at birth ( $e_0$ ) and lifespan disparity measures. Time-trends of these summary measures are often regular and well-understood. Forecasting these time-series is therefore an easier task. Figure ?? (top-left panel) presents observed  $e_0$  for Italian females from 1960 to 2016 along with the medium variant up to 2050 as computed by the UN. A second constraint is given by future values of  $e^\dagger$ , a lifespan disparity measure defined as the average years of life lost in a population attributable to death (Vaupel & Canudas Romo, 2003). Future values of this measure are obtained by conventional time-series models and portrayed in the top-right panel of Figure ?. Future mortality patterns, both by age and over time, must adhere to these predicted trends.

We arrange data as a column vector, that is,  $\mathbf{y} = \text{vec}(\mathbf{Y})$  and  $\mathbf{e} = \text{vec}(\mathbf{E})$  and we model our Poisson death counts as follows:  $\ln(E(\mathbf{y})) = \ln(\mathbf{e}) + \boldsymbol{\eta} = \ln(\mathbf{e}) + \mathbf{B}\boldsymbol{\alpha}$ , where  $\mathbf{B}$  is the regression matrix over the two dimensions:  $\mathbf{B} = \mathbf{I}_{n_1} \otimes \mathbf{B}_a$ , with  $\mathbf{B}_a \in \mathbb{R}^{m \times k_a}$ . Over time, we employ an identity matrix of dimension  $n_1$  because we will incorporate a constraint for each year. Over age,  $\mathbf{B}_a$  includes a specialized coefficient for dealing with mortality at age 0. In order to forecast, data and bases are augmented as follows:

$$\check{\mathbf{E}} = [\mathbf{E} : \mathbf{E}_2], \quad \check{\mathbf{Y}} = [\mathbf{Y} : \mathbf{Y}_2], \quad \check{\mathbf{B}} = \mathbf{I}_{n_1+n_2} \otimes \mathbf{B}_a, \quad (1)$$

where  $\mathbf{E}_2$  and  $\mathbf{Y}_2$  are filled with arbitrary future values. If we define a weight matrix  $\mathbf{V} = \text{diag}(\text{vec}(\mathbf{1}_{m \times n_1} : \mathbf{0}_{m \times n_2}))$ , the coefficients vector  $\boldsymbol{\alpha}$

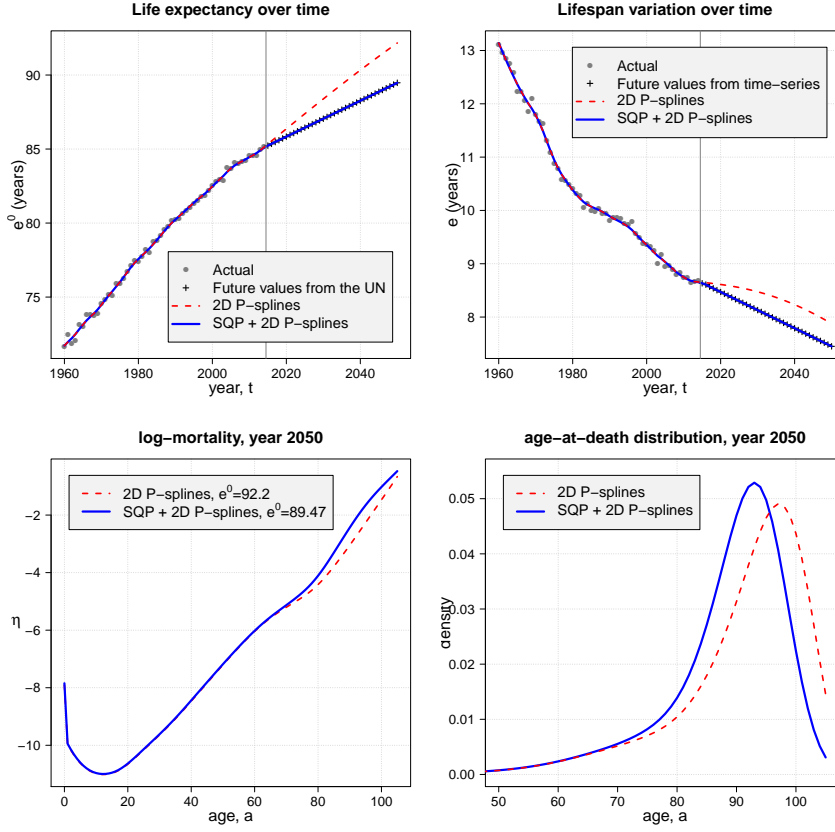


FIGURE 1. Top panels: Actual, estimated and forecast life expectancy at birth and lifespan disparity measure by United Nations and time-series, 2D  $P$ -splines and the SQP+2D  $P$ -splines. Bottom panels: Mortality in 2050 described by log-hazards and associated densities (ages 50+) by 2D  $P$ -splines and the SQP+2D  $P$ -splines. Italian females, ages 0-105, years 1960-2014, forecast up to 2050.

can be estimated by a penalised version of the iteratively reweighted least squares algorithm:

$$(\check{B}^T V \check{W} \check{B} + P) \check{\alpha} = \check{B}^T V \check{W} \check{z}, \quad (2)$$

where a difference penalty  $P$  enforces smoothness behaviour of mortality both over age and time. Outcomes from this approach in terms of life expectancy and  $e^\dagger$  are depicted with a dashed line in Figure ?? (top panels), and departures from the UN and time-series projected values are evident. Both life expectancy and average years of life lost are nonlinear function of the coefficients vector  $\alpha$ . For a year  $j$  and associated  $k_a$  coefficients  $\alpha_j$ , we denote mortality by  $\mu_j = \exp(B_a \alpha_j)$ . We can write our summary measures

as follows

$$\begin{aligned} e^0(\boldsymbol{\alpha}_j) &= \mathbf{1}_m^\top \exp[\mathbf{C} \boldsymbol{\mu}_j] + 0.5 \\ e^\dagger(\boldsymbol{\alpha}_j) &= -\exp[\mathbf{C} \boldsymbol{\mu}_j]^\top \mathbf{C} \boldsymbol{\mu}_j \end{aligned} \quad (3)$$

where  $\mathbf{C}$  is a  $(m \times m)$  lower triangular matrix filled only with -1. Constrained nonlinear optimization is therefore necessary and a SQP approach is implemented. Let denote with  $\mathbf{N}^0$  and  $\mathbf{N}^\dagger$  the  $(k_a n_2 \times n_2)$  matrices with block-diagonal structures containing derivatives of (??) with respect to  $\boldsymbol{\alpha}_j$  for  $j = n_1 + 1, \dots, n_1 + n_2$ :

$$\begin{aligned} \frac{\partial e^0(\boldsymbol{\alpha}_j)}{\partial \boldsymbol{\alpha}_j} &= \mathbf{1}_m^\top \text{diag}[\exp(\mathbf{C} \boldsymbol{\mu}_j)] \mathbf{C} \text{diag}(\boldsymbol{\mu}_j) \mathbf{B}_a \\ \frac{\partial e^\dagger(\boldsymbol{\alpha}_j)}{\partial \boldsymbol{\alpha}_j} &= -\mathbf{B}_a^\top \{ \mathbf{C}^\top [\mathbf{C} \boldsymbol{\mu}_j \circ \exp(\mathbf{C} \boldsymbol{\mu}_j)] \circ \boldsymbol{\mu}_j \} + \\ &\quad -\mathbf{B}_a^\top \{ [\mathbf{C}^\top \exp(\mathbf{C} \boldsymbol{\mu}_j)] \circ \boldsymbol{\mu}_j \}, \end{aligned} \quad (4)$$

where  $\circ$  represents element-wise multiplication. Target life expectancy and lifespan disparity for future years are given by  $n_2$ -vectors  $\mathbf{e}_T^0$  and  $\mathbf{e}_T^\dagger$ . Solution of the associated system of equations at the step  $\nu + 1$  is given by

$$\begin{bmatrix} \boldsymbol{\alpha}_{\nu+1} \\ \boldsymbol{\omega}_{\nu+1} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_\nu & : & \mathbf{H}_\nu^0 & : & \mathbf{H}_\nu^\dagger \\ \mathbf{H}_\nu^{0T} & : & \mathbf{0}_{n_2 \times n_2} & : & \mathbf{0}_{n_2 \times n_2} \\ \mathbf{H}_\nu^{\dagger T} & : & \mathbf{0}_{n_2 \times n_2} & : & \mathbf{0}_{n_2 \times n_2} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{r}_\nu - \mathbf{L}_\nu \boldsymbol{\alpha}_\nu \\ \mathbf{e}_T^0 - \mathbf{e}^0(\boldsymbol{\alpha}_\nu) \\ \mathbf{e}_T^\dagger - \mathbf{e}^\dagger(\boldsymbol{\alpha}_\nu) \end{bmatrix}, \quad (5)$$

where  $\mathbf{L}$  and  $\mathbf{r}$  are left- and right-hand-side of the system in (??), and matrices  $\mathbf{H}^0 = [\mathbf{0}_{k_a n_1 \times n_2} : \mathbf{N}^0]^\top$  and  $\mathbf{H}^\dagger = [\mathbf{0}_{k_a n_1 \times n_2} : \mathbf{N}^\dagger]^\top$ . Vector of  $\boldsymbol{\omega}$  denotes the current solution of the associated Lagrangian multipliers for both set of constraints.

Future values for  $e^0$  and  $e^\dagger$  forecast by the proposed method are exactly equal to the UN and time-series values (Figure ??, top panels). The bottom panels show the forecast mortality age-pattern in 2050: the shape obtained by the suggested approach is not a simple linear function of the plain  $P$ -splines outcome, and differences are evident by looking at the associated age-at-death distributions.

### 3 Spanish Fertility Data

We forecast Spanish fertility using three commonly-used summary measures: Total Fertility Rate describing average number of children per women in a given year, and mean and variance of childbearing age which measure fertility shape over age. In formulas:

$$\begin{aligned} TFR(\boldsymbol{\alpha}_j) &= \mathbf{1}_m^\top \boldsymbol{\mu}_j \\ MAB(\boldsymbol{\alpha}_j) &= \boldsymbol{\mu}_j^T (\mathbf{a} + 0.5) / TFR(\boldsymbol{\alpha}_j) \\ VAB(\boldsymbol{\alpha}_j) &= \boldsymbol{\mu}_j^T (\mathbf{a} + 0.5)^2 / TFR(\boldsymbol{\alpha}_j) - MAB(\boldsymbol{\alpha}_j)^2. \end{aligned} \quad (6)$$

We forecast trends of these measures by time-series analysis. We then smooth and constrain future fertility age-patterns to comply forecast values of (??) as in (??). Summary measures as well as fertility rates in 2050 are presented in Figure ???. Differences between proposed approach and plain 2D  $P$ -splines are clear. Whereas  $P$ -splines blindly extrapolate previous trends mainly accounting for the last observed years, the proposed approach enforces future age-patterns to adhere combinations of summary measures, guiding future fertility toward demographic meaningful trends.

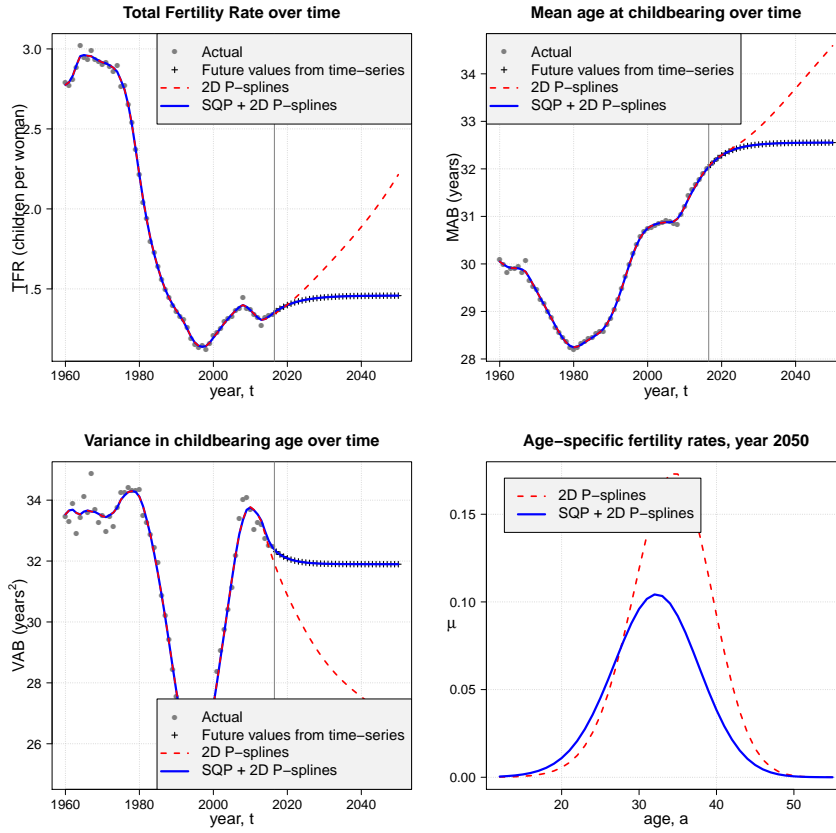


FIGURE 2. Top and left-bottom panels: Actual, estimated and forecast Total Fertility Rate, Mean and Variance in childbearing age by time-series analysis, 2D  $P$ -splines and the SQP+2D  $P$ -splines. Right-bottom panel: Age-specific fertility rate in 2050 by 2D  $P$ -splines and the SQP+2D  $P$ -splines. Spain, ages 12-55, years 1960-2016, forecast up to 2050.

## 4 Concluding remarks

In this paper, we combine smoothing models ( $P$ -splines) and optimization with nonlinear constraints (Sequential Quadratic Programming) to forecast vital rates when future values of demographic summary measures are provided.

We envisage further applications. Forecast of vital rates for partially completed cohorts is often relevant in population studies. For instance, final fertility history of a given cohort may be hypothesized though age-pattern is not yet observed and its estimation will be necessary. We also plan to adopt our approach to reconstruct demographic scenarios which are conventionally based on summary measures.

From a methodological perspective, future work will be realized to incorporate uncertainty and to objectively select the amount of smoothness in future mortality and fertility age-patterns.

## References

- Currie, I. D. et al. (2004). Smoothing and Forecasting Mortality Rates. *Statistical Modelling*, **4**, 279-298.
- Nocedal, J. & Wright, S. J. (2006). *Numerical Optimization*. Springer.
- Ševčíková, H. et al. (2016). Age-Specific Mortality and Fertility Rates for Probabilistic Population Projections. In R. Schoen (Ed.), *Dynamic demographic analysis*, 285–310. Springer.
- United Nations, Population Division (2017). *World Population Prospects: The 2017 Revision, Volume II*. ST/ESA/SER.A/400.
- Vaupel & Canudas Romo (2003). Decomposing change in life expectancy: A bouquet of formulas in honor of Nathan Keyfitz's 90th birthday. *Demography*, **40**, 201-216.