

# Forecasting of cohort fertility under a hierarchical Bayesian approach

Joanne Ellison and Erengul Dodd

*University of Southampton, UK*

and Jonathan J. Forster

*University of Warwick, Coventry, UK*

[Received March 2019. Final revision February 2020]

**Summary.** Fertility projections are a key determinant of population forecasts, which are widely used by government policy makers and planners. In keeping with the recent literature, we propose an intuitive and transparent hierarchical Bayesian model to forecast cohort fertility. Using Hamiltonian Monte Carlo methods and a data set from the human fertility database, we obtain fertility forecasts for 30 countries. We use scoring rules to assess the predictive accuracy of the forecasts quantitatively; these indicate that our model predicts with an accuracy comparable with that of the best-performing models in the current literature overall, with stronger performance for countries without a recent structural shift. Our findings support the position of hierarchical Bayesian modelling at the forefront of population forecasting methods.

**Keywords:** Cohort fertility; Forecasting; Hamiltonian Monte Carlo methods; Hierarchical Bayesian models; Human fertility database; Scoring rules

## 1. Introduction

Fertility is one of the three components of population change, together with mortality and migration—population forecasts are obtained by projecting these components forward under certain assumptions and methodologies. Government policy makers, decision makers and planners use these forecasts for numerous purposes such as planning for the future provision of basic societal needs, e.g. food, water and energy, as well as health and education services, shaping policies both locally and nationally, determining fiscal projections and informing pensions models (Office for National Statistics, 2019a; Population Reference Bureau, 2001; National Records of Scotland, 2019). Fertility forecasts specifically are required to plan maternity and child care services and to predict demand for nursery school places, as well as various other uses (Office for National Statistics, 2019b; Shang, 2012; National Records of Scotland, 2019). As a result, models that can generate plausible fertility forecasts with appropriate uncertainty are in high demand.

There is a large body of literature concerning the proposal of models to produce accurate estimates of fertility rates. Much of the early work focuses on parametric techniques, which involve choosing functions to fit closely to the bell-shaped curves of age-specific fertility rates (for example, see Peristera and Kostaki (2007)). Such functions include polynomials (Brass, 1960), the Coale–Trussell function (Coale and Trussell, 1974), beta and gamma distributions

*Address for correspondence:* Joanne Ellison, Department of Mathematical Sciences, University of Southampton, Highfield, Southampton, SO17 1BJ, UK.  
E-mail: [J.Ellison@soton.ac.uk](mailto:J.Ellison@soton.ac.uk)

(Hoem *et al.*, 1981) and the Hadwiger distribution (Hadwiger, 1940). In addition, the relational models of Brass (1974) assume a linear relationship between an observed set of fertility rates and some standard (Booth, 2006). Hoem *et al.* (1981) described a selection of these methods and gave a detailed comparison. More recently, a slight hump at younger ages has appeared in the curves for some developed countries (Chandola *et al.*, 1999), which the previous models cannot respond to (Peristera and Kostaki, 2007). Chandola *et al.* (1999) suggested a possible cause as the emerging differences between marital and non-marital fertility, proposing a mixture of Hadwiger functions to model this phenomenon.

Since the 1980s, attention has largely moved to models that treat fertility as stochastic rather than deterministic, and so can quantify the uncertainty in forecasts (Wiśniowski *et al.*, 2015). The functional model of Lee (1992) is particularly notable in this school of thought, with its use of principal component analysis and time series inspiring many approaches that involve modelling the randomness of fertility (Booth, 2006). These further functional methods include the work of Hyndman and Ullah (2007) and Shang (2012). Consistent with the change to a probabilistic viewpoint, Bayesian models are now a popular choice in the population forecasting literature as they can incorporate uncertainty naturally. Recent papers such as Wiśniowski *et al.* (2015) and Bijak and Bryant (2016) have attributed the rise in their usage to computational developments occurring as recently as the last decade. Hierarchical Bayesian models (for example, see Girosi and King (2008)), which allow borrowing of strength, are also becoming increasingly common. This strength can be borrowed from other countries, which is an approach that was used in the methodology of the first probabilistic population projections to be published by the United Nations (Ševčíková *et al.*, 2016). Alternatively, the strength can be borrowed across ages and cohorts (for example, see Czado *et al.* (2005)) for the fertility rate estimates of a single country.

In an attempt to determine whether the increasingly sophisticated and computationally expensive models that have been proposed in recent years have actually led to greater predictive accuracy, Bohk-Ewald *et al.* (2018a) performed a comprehensive comparison of 20 existing cohort fertility forecasting approaches. Taking a cohort approach, the aggregate fertility measure that they used to compare the methods is the cohort total fertility rate CFR. CFR is calculated by summing the age-specific rates for a given cohort, i.e. a group of women with the same birth year, across all reproductive ages (Jasilioniene *et al.*, 2015); as such, it can be interpreted as the average completed family size for that cohort. The equivalent measure under a period approach is the total fertility rate TFR, which sums the rates for a given calendar year—this can be interpreted as the average completed family size for a *hypothetical cohort* of women who experience the fertility rates of that one year throughout their reproductive lives (Ní Bhrolcháin, 2011).

By its definition, TFR provides a summary of fertility over a brief period of time which can be very recent, and is therefore more immediately relevant compared with CFR (Bongaarts and Feeney, 1998). However, as well as the absence of a practical interpretation, a further drawback is that changes in TFR can be due to tempo effects, i.e. changes in the average age of childbearing during the period in question, rather than just quantum effects, i.e. changes in the average number of children per woman (Bongaarts and Feeney, 1998). As a result, cohort fertility tends to be more stable across time than period fertility is (de Beer, 1985; Li and Wu, 2003), which makes it a more appealing measure for forecasting. For these reasons, combined with the frequent adoption of the cohort approach in the recent fertility forecasting literature, we also decide to take a cohort approach.

Returning to the work of Bohk-Ewald *et al.* (2018a), they found that, in terms of forecast accuracy, four methods perform better than the naive freeze rates approach, which simply freezes the age-specific rates at their most recent observed values. These superior approaches include the two simple extrapolation methods of Myrskylä *et al.* (2013a) and de Beer (1985, 1989). The

former extrapolates the age-specific trends that are exhibited over the previous 5 years for a further 5 years before freezing the rates; the latter extrapolates patterns that are exhibited by the rates across ages and cohorts jointly by fitting two interconnected auto-regressive integrated moving average time series models. The remaining two successful approaches are both Bayesian methods, namely the conjugate normal–normal model of Schmertmann *et al.* (2014a) and the aforementioned model of Ševčíková *et al.* (2016). The former constructs a prior from quadratic penalties, simultaneously penalizing potential future patterns of rates in the age and cohort dimensions that are deemed unlikely by the historical data; the latter first forecasts TFR by using a hierarchical Bayesian model, subsequently decomposing it into age-specific projections according to a particular target pattern determined by expert opinion.

Of the eight methods that enable uncertainty quantification (this includes the methods that we have described except for that of de Beer (1985, 1989)), the Bayesian model of Schmertmann *et al.* (2014a) appears to perform strongest and therefore could be seen as the best when considering forecast accuracy and uncertainty together. Overall, the questionable dominance of the Bayesian approaches over simple extrapolation methods causes Bohk-Ewald *et al.* (2018a) to question whether such complex models requiring large amounts of data and computation time are really necessary to obtain accurate cohort fertility forecasts—this is one of the motivations of our work.

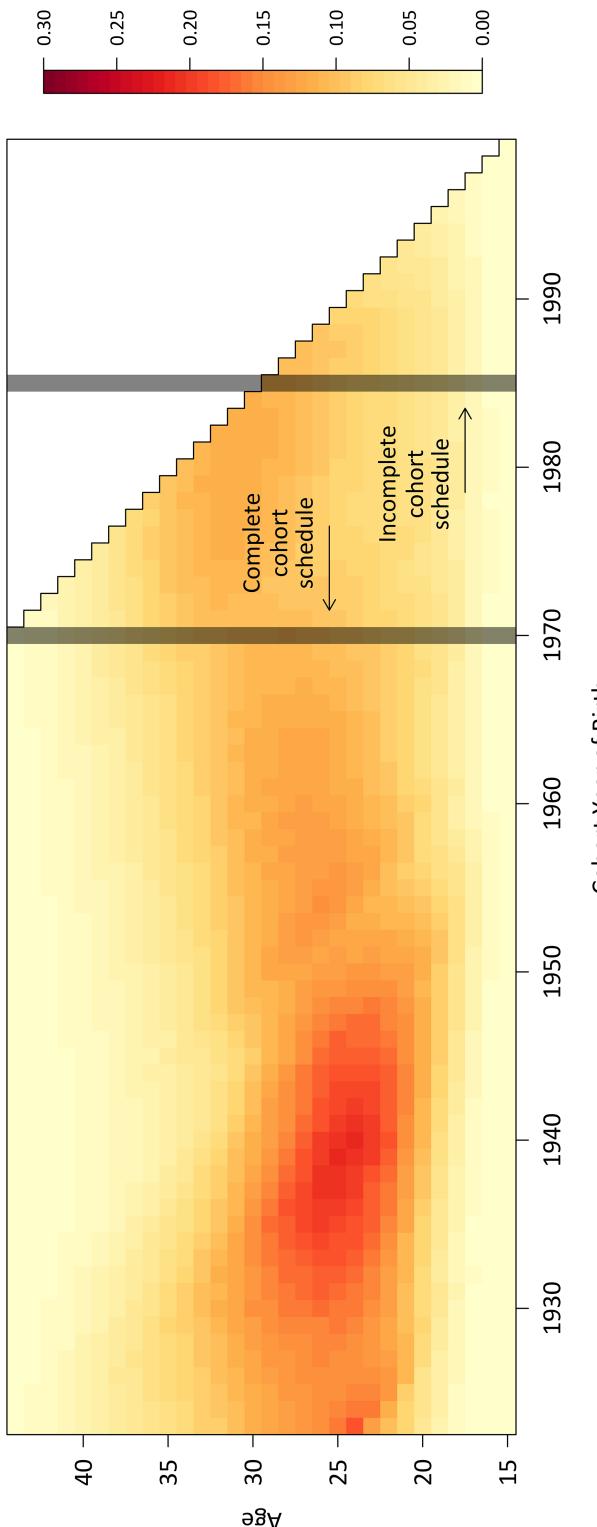
This discussion brings us to the main purpose of this paper. In the spirit of the highly successful model of Schmertmann *et al.* (2014a) and in keeping with the most recent literature, we propose a hierarchical Bayesian model for forecasting cohort fertility. By incorporating our assumptions explicitly into the model structure and then letting the data determine their precise satisfaction, we aim to construct a transparent and intuitive model with realistic levels of forecast uncertainty that can compete with the current best-performing models in the field. We fit our model by using the state of the art Hamiltonian Monte Carlo computational methodology implemented by the software RStan (Stan Development Team, 2018a).

After presenting our approach in Section 2, in Section 3 we compare the forecast performance of our model with that of Schmertmann *et al.* (2014a), Myrskylä *et al.* (2013a) and de Beer (1985, 1989). We fit to the fertility data that were available in 2014 to generate forecasts for 30 countries and perform a qualitative comparison. We also fit to the data that were available 10 years earlier in 2004 to generate forecasts for 29 countries, enabling us to use scoring rules and other summary statistics to compare the approaches quantitatively. Lastly, we discuss our findings in Section 4.

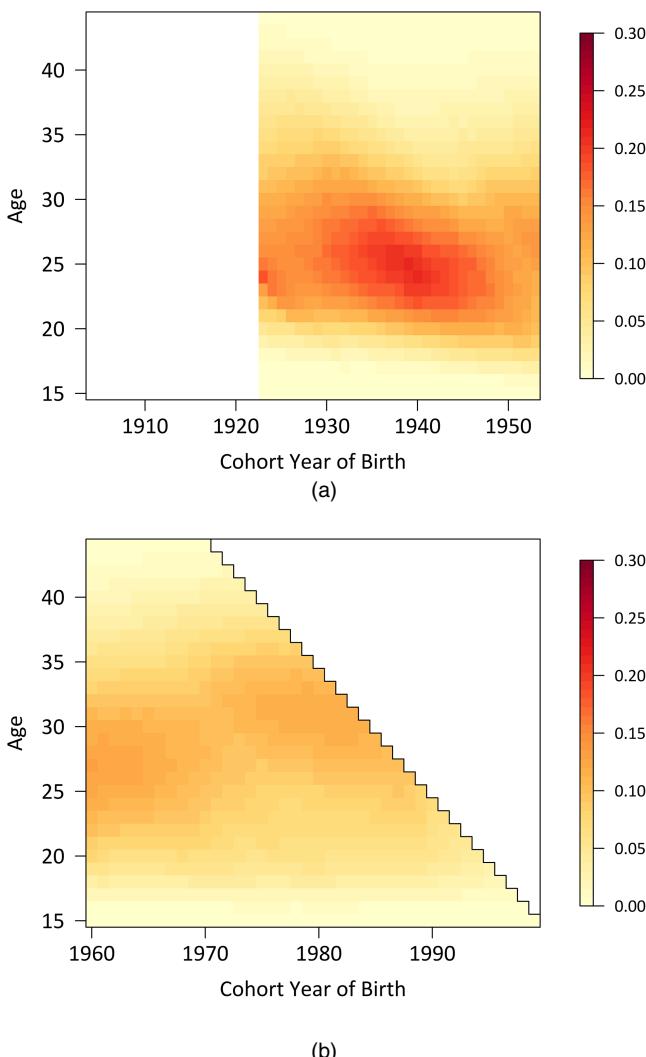
## 2. Method

### 2.1. Introduction

Consider an incomplete Lexis surface, i.e. a heat map of fertility rate estimates plotted by age against cohort year of birth. For a cohort of women at a given age, the fertility rate is estimated by dividing the number of live births by the number of women (Jasilioniene *et al.*, 2015). Fig. 1 gives an example of a Lexis surface for England and Wales data from Human Fertility Database (2019), taking the present to be 2014 by using only the data that would have been available in that year; also note that the surface is restricted to ages 15–44 years and the 1923–1999 cohorts. Features of interest include the high rates in the dark region which occurred during the 1960s and early 1970s for women in their 20s and the increase in the peak age of childbearing from the 1960s cohorts to the late 1970s cohorts (Office for National Statistics, 2017). We call the set of rates for each cohort a cohort schedule, with a cohort schedule complete if it is fully observed and incomplete otherwise. Using this terminology, we see from Fig. 1 that the cohort schedules



**Fig. 1.** Lexis surface of England and Wales Human Fertility Database (2019) fertility rate estimates by age against cohort year of birth, taking the present to be 2014: white cells correspond to future rates which are yet to be observed; the last complete cohort schedule (of the 1970 cohort) is indicated, as well as an incomplete cohort schedule (of the 1985 cohort)

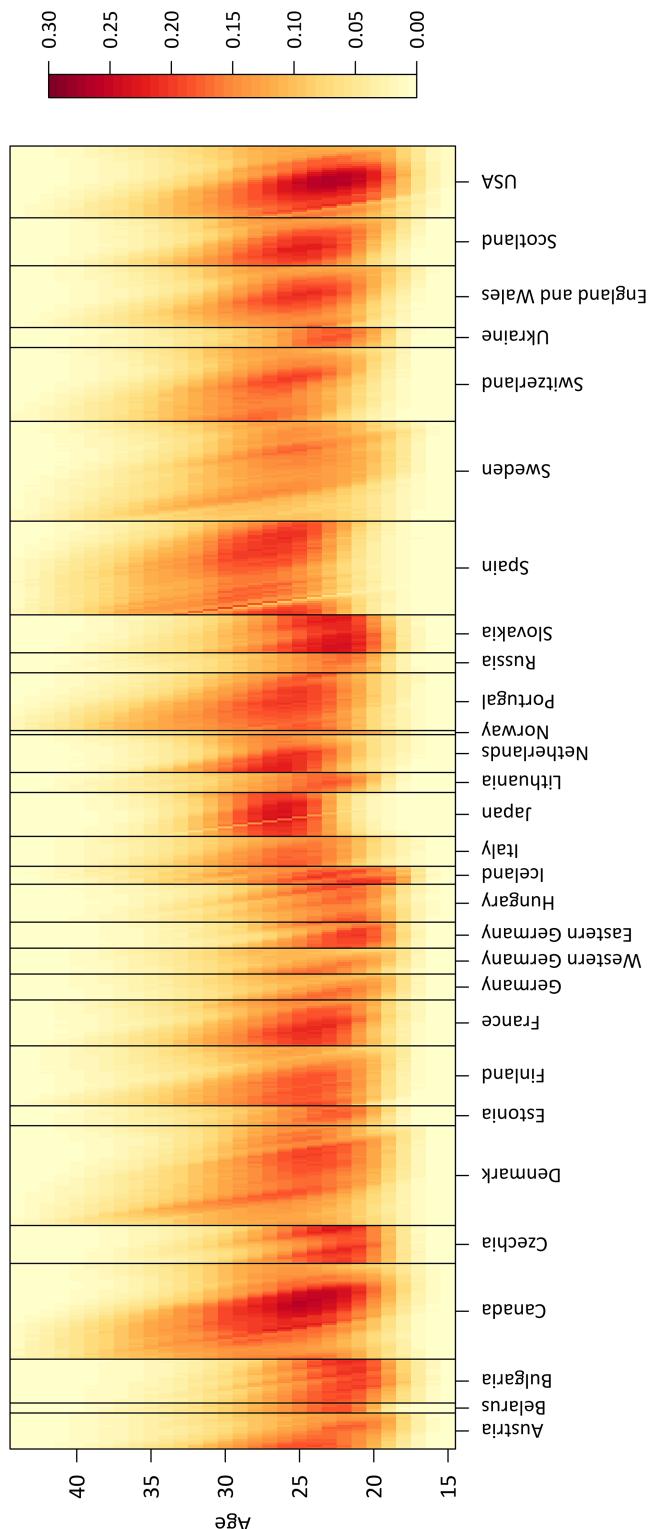


**Fig. 2.** England and Wales (a) historical and (b) contemporary Lexis surfaces (Human Fertility Database, 2019)

up to and including the 1970 cohort are complete—this is because, in 2014, the 1970 cohort would have been the youngest cohort to have an observable rate for women aged 44 years.

The data set that we consider as a whole consists of complete and incomplete cohort schedules for ages 15–44 years from countries across Europe, with several countries from North America and Asia. Following Schmertmann *et al.* (2014a), we separate the data into historical and contemporary sections. For each country, we take all the available complete cohort schedules for the 1904–1953 cohorts to form the historical part, and the 1960–1999 cohort schedules to form the contemporary part. We illustrate this in Fig. 2 for England and Wales. From Fig. 2(a), we see that England and Wales contribute only 31 of the 50 desired historical cohort schedules from the 1923 cohort onwards.

The reason for this separation is apparent on considering the model structure. For a given country, the parameters are the true fertility rates in its contemporary Lexis surface (one for



**Fig. 3.** Combined historical Lexis surfaces from all the countries in our data set (Human Fertility Database, 2019)

each of the possible cohort–age combinations in Fig. 2(b)). The historical data from all the countries (which are displayed in Fig. 3) inform the core of the model. We specify a prior for the (hyper)parameters, which is then updated by combining its information with that obtained from the country's observed contemporary rates (the rate estimates in Fig. 2(b)) in the likelihood. This gives a posterior distribution for the parameters of the entire contemporary Lexis surface. We specify our model in Section 2.2.

## 2.2. Model specification

Suppose that there are  $C$  birth cohorts ( $c = 1, \dots, C$ ) and  $A$  ages ( $a = 1, \dots, A$ ) in the contemporary Lexis surface of a particular country. For each cohort–age combination  $(c, a)$ , let  $N_{ca}$  be the observed number of births,  $W_{ca}$  be the number of women alive (i.e. the exposure) and  $\theta_{ca}$  be the true fertility rate; then  $\mu_{ca} = W_{ca}\theta_{ca}$  is the mean number of births. Because of its suitability for modelling count data, we assume a Poisson distribution for  $N_{ca}$  with mean  $\mu_{ca}$ , i.e.  $N_{ca} \sim \text{Poisson}(\mu_{ca})$ . We model  $\theta_{ca}$  on the logarithmic scale, which is the standard approach to take when modelling rates. It also ensures that our model will not generate negative forecasts, which is especially beneficial in instances where we have diminishing fertility at ages where the level is already low. Letting  $H$  be the total number of complete historical cohort schedules contributed by all the countries, we define  $\Phi$  to be the  $A \times H$  matrix of these cohort schedules (see Fig. 3 in Section 2.1 for a visual representation of such a matrix with  $A = 30$  and  $H = 653$ ). Let  $\Pi$  be the equivalent matrix on the logarithmic scale, i.e.  $\Pi_{ah} = \log(\Phi_{ah})$ ,  $a = 1, \dots, A$ ,  $h = 1, \dots, H$ . We then let  $\mathbf{X}$  be the  $A \times 3$  matrix consisting of the first three principal components of  $\Pi$ , obtained from its singular value decomposition. These components are essentially highly significant covariates that together explain a large proportion of the variation of the historical cohort schedules. The underlying assumption, which was also made in Schmertmann *et al.* (2014a), is that these components will explain a similar proportion of the variation in the contemporary cohort schedules. We now define the core of the model, which takes the following log-linear form:

$$\log(\theta_{ca}) = [\mathbf{X}\beta]_a + \varepsilon_{ca}, \quad (1)$$

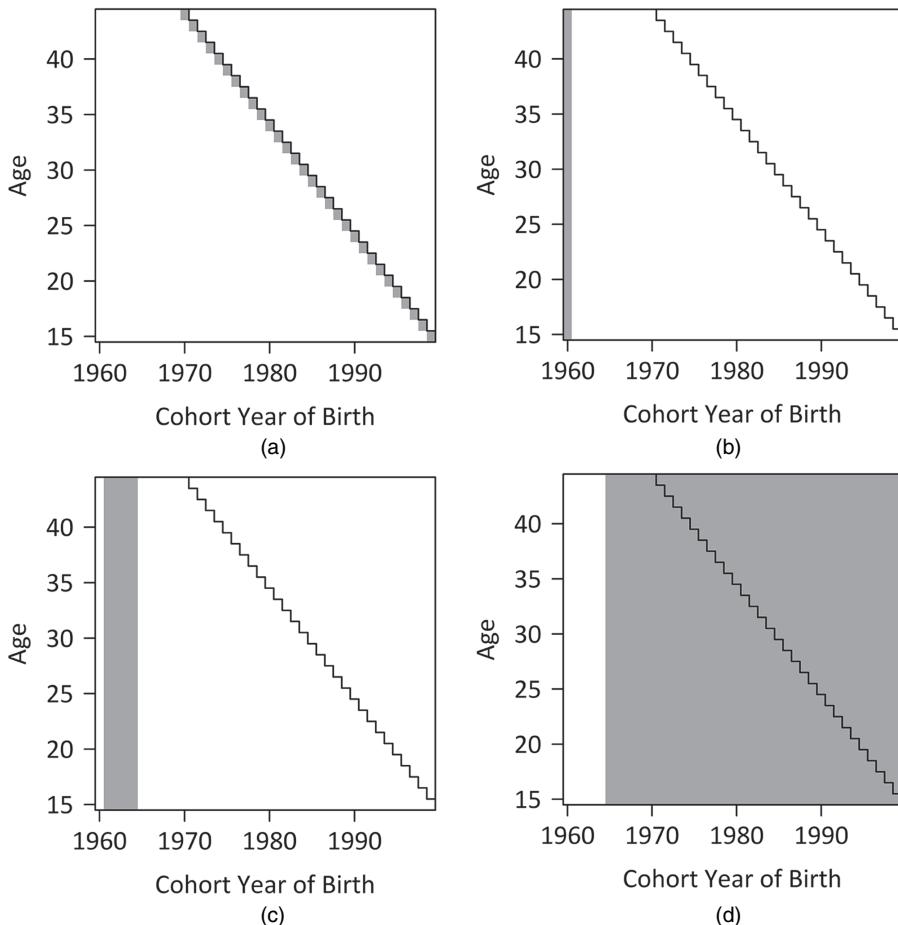
where  $\beta = (\beta_1, \beta_2, \beta_3)'$  is the vector of regression parameters, i.e. the respective weights on the three principal components—see Schmertmann *et al.* (2014a) for an interpretation of these weights when the principal components of  $\Phi$  are computed. The  $\beta$ s are allowed to vary freely, imposing no constraints on the possible shapes or levels of the incomplete cohort schedules; we ensure that this is so by giving each  $\beta_i$  a diffuse  $N(0, 30^2)$  prior. The  $\varepsilon_{cas}$  are the error terms, and we enforce our remaining model assumptions through their prior distribution.

Firstly, letting  $\mathbf{b} = (b_1, b_2, b_3)'$ , from model (1) it is clear that the model is invariant under the following parameterization:

$$\{\beta, \varepsilon_{ca}\} \rightarrow \{\beta + \mathbf{b}, \varepsilon_{ca} - [\mathbf{X}\mathbf{b}]_a\}.$$

To resolve this identification problem, we impose a constraint on the vector of ‘jump-off’  $\varepsilon_{cas}$ , i.e. the  $\varepsilon_{cas}$  whose cohort–age combinations correspond to the calendar year that we are taking to be the present—we shall discuss the reason for this choice in due course. We call this vector  $\varepsilon_{JO}$  and illustrate it in Fig. 4(a). We let  $\varepsilon_{JO} = (\mathbf{I}_A - \mathbf{P}_X)\boldsymbol{\eta}$ , where  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_A)'$  and  $\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ : the projection matrix of  $\mathbf{X}$ . We then let  $\eta_a \stackrel{\text{IID}}{\sim} N(0, \sigma_1^2)$ ,  $a = 1, \dots, A$ . Hence  $\mathbf{X}'\varepsilon_{JO} = \mathbf{0}$ . As  $\mathbf{X}$  has three columns, this imposes three linear constraints on  $\varepsilon_{JO}$ , fixing its level and therefore making the model identifiable.

Next, we divide the contemporary Lexis surface into three regions A–C, which are illustrated



**Fig. 4.** Representation of (a)  $\varepsilon_{10}$ , (b) region A, (c) region B and (d) region C on a typical contemporary Lexis surface where the present is 2014, with the included cohort–age combinations filled in grey

in Figs 4(b)–4(d). Region A consists of the  $\varepsilon_{cas}$ s with cohort–age combinations in the first cohort only, region B those in cohorts 2–5, and region C those in cohort 6 onwards. We then specify the priors on the  $\varepsilon_{cas}$ s by the region that they belong to. For region A, let

$$\varepsilon_{1a} \stackrel{\text{IID}}{\sim} N(0, \sigma_2^2), \quad a = 1, \dots, A;$$

for region B, let

$$\varepsilon_{ca} \sim N(\rho_1 \varepsilon_{(c-1)a}, \sigma_3^2), \quad c = 2, \dots, 5;$$

for region C, let

$$\varepsilon_{ca} \sim N\{\rho_2 \varepsilon_{(c-1)a} + \rho_3 (\varepsilon_{(c-1)a} + \hat{\Delta}_{(c-1)a}), \sigma_4^2\}, \quad c = 6, \dots, C,$$

where  $\hat{\Delta}_{(c-1)a} = (1/30)(10\varepsilon_{(c-1)a} - \varepsilon_{(c-2)a} - 2\varepsilon_{(c-3)a} - 3\varepsilon_{(c-4)a} - 4\varepsilon_{(c-5)a})$  is the ordinary least squares slope estimator obtained by fitting a linear regression model without an intercept to the five error terms corresponding to age  $a$  and cohorts  $c-5, \dots, c-1$ .

Lastly, we assume that each  $\sigma_i \sim N^+(0, 0.25^2)$  and each  $\rho_i \sim N^+(0, 0.5^2)$ , where  $N^+(\cdot, \cdot)$  indicates a half-normal prior.

In words, our region A prior states that the error terms are simply independently and identically normally distributed with mean 0 and variance  $\sigma_2^2$ . Our region B prior states that, for each age, the  $\varepsilon_{cas}$ s follow an auto-regressive AR(1) process across cohort with coefficient  $\rho_1$  and constant error variance  $\sigma_3^2$ . In region C, the prior states that the  $\varepsilon_{cas}$ s are normally distributed with mean equal to a weighted combination of the previous age-specific error  $\varepsilon_{(c-1)a}$ , and the sum of this error term and the slope  $\hat{\Delta}_{(c-1)a}$  that it gives rise to along with the previous four age-specific errors, with weights  $\rho_2$  and  $\rho_3$  respectively and variance  $\sigma_4^2$ . The region C prior is the most important as it will determine the forecasts—this is because this region includes all the future cohort–age combinations. What we are actually doing in this prior is balancing the two most common extrapolation methods for observed fertility rates in the demographic forecasting literature, which we shall call the freeze rate and freeze slope approaches in line with Schmertmann *et al.* (2014a). The freeze rate approach assumes that the next age-specific rate will be similar to the previous rate, i.e.  $\theta_{ca} \approx \theta_{(c-1)a}$ . In contrast, the freeze slope approach assumes that the next age-specific rate will follow the recent trend of its past rates, i.e.  $\theta_{ca} \approx \theta_{(c-1)a} + \hat{\delta}_{(c-1)a}$ , where  $\hat{\delta}_{(c-1)a}$  is the recent slope—we take this to be calculated by using the last five rates in the same spirit as Myrskylä *et al.* (2013a) and Schmertmann *et al.* (2014a).

In our model we are working on a logarithmic scale and with these assumptions applied to the error terms instead of the actual rates; these are two of the key differences from the model of Schmertmann *et al.* (2014a). For the freeze rate approach the two are equivalent by using model (1), i.e.  $\varepsilon_{ca} \approx \varepsilon_{(c-1)a} \Rightarrow \theta_{ca} \approx \theta_{(c-1)a}$ . However, for the freeze slope approach they are not, as  $\varepsilon_{ca} \approx \varepsilon_{(c-1)a} + \hat{\Delta}_{(c-1)a} \Rightarrow \theta_{ca} \approx \theta_{(c-1)a} \exp(\hat{\Delta}_{(c-1)a})$ , again using model (1). This is intuitive, as a small change on the log-scale is approximately equivalent to the proportionate change on the original scale.

Returning to the region C prior, it is now clear that we are allowing the data to choose how much weight to put on the freeze rate and freeze slope assumptions through  $\rho_2$  and  $\rho_3$  respectively. We do not constrain these parameters to sum to 1, as if  $\rho_2 + \rho_3 < 1$  each age-specific process is stationary and will revert to  $[X\beta]_a$  in the long term. If  $\rho_2 + \rho_3 > 1$  then the process is non-stationary and will not exhibit long-term reversion, instead having a more explosive nature. By leaving the sum of  $\rho_2$  and  $\rho_3$  unconstrained, for each country we allow the parameters to learn from the observed contemporary fertility rate estimates, choosing whether they want to follow a stationary or non-stationary process as a result. The degree of stationarity, i.e. how close  $\rho_2 + \rho_3$  is to 1, is also significant, as it determines how strongly the reversion occurs in the forecasts. It is this reversion of a stationary process to  $[X\beta]_a$  that motivated our decision to constrain  $\varepsilon_{JO}$  to make the model identifiable earlier in this section. The rationale behind this is that, if we could choose where we would want our forecasts to revert to in the future, it would be somewhere close to where we started forecasting from, i.e. the value in the calendar year taken to be the present, as opposed to the initial value or some average across the contemporary cohort schedules. This is in line with the current fertility literature, e.g. the use of only the last 5 years of data in Myrskylä *et al.* (2013a). In constraining  $\varepsilon_{JO}$ , therefore, the desire is that for each country  $[X\beta]_a$  will be close to the jump-off value for each age  $a$ . We discuss this further in Section 3.2.3.

Here we propose a hierarchical model to borrow strength across ages and cohorts. This is particularly important in region C, where we allow  $\rho_2$ ,  $\rho_3$  and  $\sigma_4^2$  to learn from all the observed cohort–age combinations after the fifth cohort. In this way, our forecasts take as much information as possible from the contemporary Lexis surface about the relative importance of the freeze rate and freeze slope assumptions. Our model is also hierarchical in the typical sense, in that we have two levels of priors due to the presence of the hyperparameters (the  $\sigma_i$ s and the  $\rho_i$ s). We give an overview of the model fitting in Section 2.3.

### 2.3. Model fitting

The hierarchical nature of our proposed model means that it is not possible to write the posterior in closed form—instead, we need to approximate it by using Monte Carlo methods. Such methods are also required for the hierarchical Bayesian model of Ševčíková *et al.* (2016), whereas the posterior of the conjugate Bayesian model of Schmertmann *et al.* (2014a) is tractable and hence can be computed precisely. For complex hierarchical models such as ours, well-known Markov chain Monte Carlo methods like the Metropolis algorithm are less satisfactory because of their local random-walk behaviour, i.e. slow exploration of the posterior (for example, see Gelman *et al.* (2014)). The method of Hamiltonian Monte Carlo increases the efficiency of this exploration through Hamiltonian dynamics (for more details see Stan Development Team (2018b)). The variant of Hamiltonian Monte Carlo sampling that we shall use in the computation of the model proposed is the no-U-turn sampler, which determines certain algorithm parameters adaptively in each iteration to maximize the exploration distance relative to the current position. The sampler is implemented by the software RStan (Stan Development Team, 2018a), which we shall use to fit our model in Section 3.

The fitting process consists of two parts. First, we use RStan to perform  $T$  iterations of the no-U-turn sampler algorithm, following a warm-up period of  $T'$  iterations for estimation and optimization of algorithm parameters. This generates  $T$  samples of  $\beta$ , the  $\sigma_i$ s and the  $\rho_i$ s, as well as the  $\varepsilon_{ca}$ s with observed cohort–age combinations and therefore observed values of  $N_{ca}$  and  $W_{ca}$ . Second, for forecasting we need to obtain  $T$  samples of the  $\varepsilon_{ca}$ s with unobserved cohort–age combinations. We do this by simulating them from  $N\{\rho_2\varepsilon_{(c-1)a} + \rho_3(\varepsilon_{(c-1)a} + \hat{\Delta}_{(c-1)a}), \sigma_4^2\}$ : one set of samples for each of the  $T$  original samples. This gives  $T$  samples of  $\theta_{ca} = \exp([X\beta]_a + \varepsilon_{ca})$  for each cohort–age combination. We can code our model in such a way that RStan can perform this simulation within each iteration.

The posterior distribution enables us to quantify our uncertainty about the true rates  $\theta_{ca}$  but not the empirical birth rates. We need to incorporate the additional variation to account for the fact that we are predicting an observation and not the mean. The process by which we do this for a conjugate model (e.g. Schmertmann *et al.* (2014a)) is described in Appendix A and is simple because we have the posterior distribution in closed form. The process is slightly more involved for the proposed model, however, as we have only  $T$  samples of each  $\theta_{ca}$  available to us. Our goal is to have a credible interval for each empirical rate as these are what we are trying to forecast; for this reason we need to generate empirical birth rates from our  $T$  samples of each true birth rate  $\theta_{ca}$ . For each  $(c, a)$ , we do this by first sampling a random observation from  $\text{Poisson}(\mu_{ca}^t)$  for  $t = 1, \dots, T$ , where  $\mu_{ca}^t = W_{ca}\theta_{ca}^t$  is the  $t$ th sampled value of  $\mu_{ca}$  and  $\theta_{ca}^t$  the  $t$ th sampled value of  $\theta_{ca}$ . If  $(c, a)$  is unobserved, we take  $W_{ca}$  to be its most recently observed value at age  $a$ . We then divide each of these  $T$  Poisson realizations by  $W_{ca}$  to obtain a sample of  $T$  empirical birth rates. We then compute the 90% and 50% credible intervals (CIs), which are probability intervals based on the posterior predictive distribution. We do this by extracting the (5%, 95%) and (25%, 75%) quantiles of the sample of empirical birth rates. Note that the additional uncertainty has a noticeable effect only for small countries with comparatively low exposures.

## 3. Results

### 3.1. Data and computation

To assess the forecast performance of our proposed model, we follow the advice of Bohk-Ewald *et al.* (2018a) to compare against the naive freeze rates method and the simple extrapolation

models of Myrskylä *et al.* (2013a) and de Beer (1985, 1989) at a minimum; we shall refer to these as models MGC and dB respectively, after the authors. We additionally include the model of Schmertmann *et al.* (2014a) in our comparison (denoted model SZGM), as our proposed hierarchical Bayesian model (denoted model hB) has been developed in the same spirit. We fit the models to the countries that are available in the Human Fertility Database (2019) data set. (Note that this is an update to the 2011 version of the data set that was used in Schmertmann *et al.* (2014a). It includes 12 additional countries (Belarus, Chile, Croatia, Denmark, Iceland, Italy, Japan, Norway, Poland, Spain, Taiwan and Ukraine), modifications to rate estimates available in 2011 and additional rate estimates.) We first fit the models to the data that were available in 2004, enabling us to use the more recent (or hold-out) data to perform a quantitative comparison using scoring rules and various summary statistics relating to point *and* probabilistic accuracy. This generates forecasts for 29 countries, which we call ‘2004 (fertility) forecasts’ and present in Section 3.2. We then incorporate the hold-out data directly by fitting the models to the data that were available in 2014, generating ‘2014 (fertility) forecasts’ for 30 countries which we discuss in Section 3.3. The R code and Stan files to obtain the model hB forecasts are available from <https://github.com/jvellison/hBfert>. For the model SZGM computation we use the R code that is available from the Schmertmann *et al.* (2014a) project web site (Schmertmann *et al.*, 2014b), slightly modified to account for the change in data and incorporation of additional variation (see Appendix A). For model MGC we use our own R code written according to the method that is described in Myrskylä *et al.* (2013a) and inspired by the corresponding Stata code (Myrskylä *et al.*, 2013b) as well as the R code that was used to produce the work of Bohk-Ewald *et al.* (2018a) (Bohk-Ewald *et al.*, 2018b); we also use this R code to implement model dB. We fit model hB to each country separately, with  $T' = 1000$  warm-up iterations followed by  $T = 4000$  retained iterations (see Section 2.3) thinned by a factor of 2 from an initial 8000. We examine convergence for each fit in the conventional way, and find that the samples mix well across the  $T$  iterations.

### 3.2. 2004 fertility forecasts

In this section we provide an analysis of the 2004 forecasts by using various approaches. First, in Section 3.2.1 we use scoring rules to compare the accuracy of the model hB, SZGM, MGC, dB and freeze rates forecast distributions quantitatively. In Section 3.2.2 we use typical summary statistics to compare their point accuracy and coverage (where possible). Then, in Section 3.2.3 we graphically explore both the stationarity of the model hB forecasts and the way in which the degree of stationarity affects the nature of the reversion.

#### 3.2.1. Scoring rules

To compare the models on their forecast precision and uncertainty we use scoring rules, which are measures of predictive accuracy for probabilistic prediction (Gelman *et al.*, 2014). This means that they consider the posterior distribution as a whole rather than a summary statistic of it such as the mean or median. For a probabilistic forecast  $G$  (with associated cumulative distribution function (CDF)  $G$  and probability density function  $g$ ) and observed value  $y$ , a scoring rule summarizes the suitability of  $G$  in light of  $y$  by a score. The scoring rules that we shall use are negatively oriented, which means that smaller scores are desirable (Jordan *et al.*, 2017). Following Jordan *et al.* (2017) we compute the logarithmic score LogS (Good, 1952) and the continuous ranked probability score (CRPS) (Matheson and Winkler, 1976):

$$\text{LogS}(G, y) = -\log\{g(y)\}; \quad (2)$$

$$\text{CRPS}(G, y) = \int_{\mathbb{R}} \{G(z) - \mathbb{I}(y \leq z)\}^2 dz, \quad (3)$$

where  $\mathbb{I}(\cdot)$  is the indicator function. From equation (2), it is clear that LogS penalizes forecast distributions that assign negligible probability to  $y$ . In contrast, the CRPS ‘generalises the absolute error’ (Gneiting and Raftery, 2007) by penalizing forecast distributions whose CDFs differ substantially from the empirical CDF of  $y$ , i.e. the ‘perfect forecast’ (Bröcker, 2012). In fact, equation (3) reduces to the absolute error if  $G$  is a deterministic (or point) forecast; therefore the CRPS enables us to score probabilistic and deterministic forecasts under the same metric, which is an attractive property (Gneiting and Raftery, 2007). Consequently, it is not surprising that the CRPS is more ‘sensitive to distance’ (Gneiting and Raftery, 2007) and so would score a narrow distribution with median close to  $y$ , but negligible probability assigned to  $y$ , more favourably than LogS—this is because the forecast is accurate, even though it is too precise. For further discussion of scoring rules, see Gneiting and Raftery (2007).

Regarding computation for the models with probabilistic forecasts (models hB, SZGM and MGC), LogS and the CRPS can be calculated exactly under models SZGM and MGC because of their Gaussian forecast distributions; the same is not true for model hB because of its intractable posterior (see Section 2.3), and as a result approximations are required. We use a Gaussian approximation for  $g$  to compute LogS and an empirical CDF-based approximation for  $G$  to compute the CRPS (see Krüger *et al.* (2019) for details).

We fit the models over ages 15–44 years, using the 1950–1989 cohorts as our contemporary data and maintaining the 1904–1953 cohorts as our historical data as in Section 2.1. We fix the number of complete contemporary cohorts at 11 (meaning that the 1950–1960 cohorts are complete and the 1961–1989 cohorts are increasingly incomplete), and focus solely on CFR for simplicity and consistency with the recent literature. For each model we compute the CRPS for the CFR-forecasts with a corresponding observed value—the number of such forecasts varies by country because of data availability, ranging from 5 to 13 with a modal value of 12 observed CFR-values available after 2004 (for the 1961–1972 cohorts). We then plot the average CRPS by country in Fig. 5, noting that this reduces to the mean absolute error for freeze rates and model dB as their forecasts are deterministic. We order the countries by decreasing average CRPS (increasing predictive accuracy) under model hB, for ease of comparison against the other methods.

Fig. 5 provides strong support for model hB being competitive with the current best cohort fertility forecasting methods—its average score is significantly better than that for freeze rates for 27 of the 29 countries, and only marginally worse for the remaining two countries (Iceland and Bulgaria). In terms of its performance among the models that were identified as the most accurate in Bohk-Ewald *et al.* (2018a), it is fair to say from inspecting Fig. 5 that, for all countries from Denmark to the right (excluding England and Wales, and Slovakia), the difference between the model hB average score and the lowest for that particular country (achieved by model hB for Portugal and the USA) is negligible. For the countries left of Denmark we see larger differences from the minimum, in particular for Switzerland, Iceland and Estonia; model SZGM exhibits the opposite trend, performing very strongly for these countries overall but poorly for Denmark, Sweden, Finland, the USA, the Netherlands and Spain. To facilitate a fair comparison, we present histograms of the average scores for each model in Fig. 6—the model hB and SZGM plots indicate that, excluding the Netherlands, their average CRPS distributions across countries are quite similar.

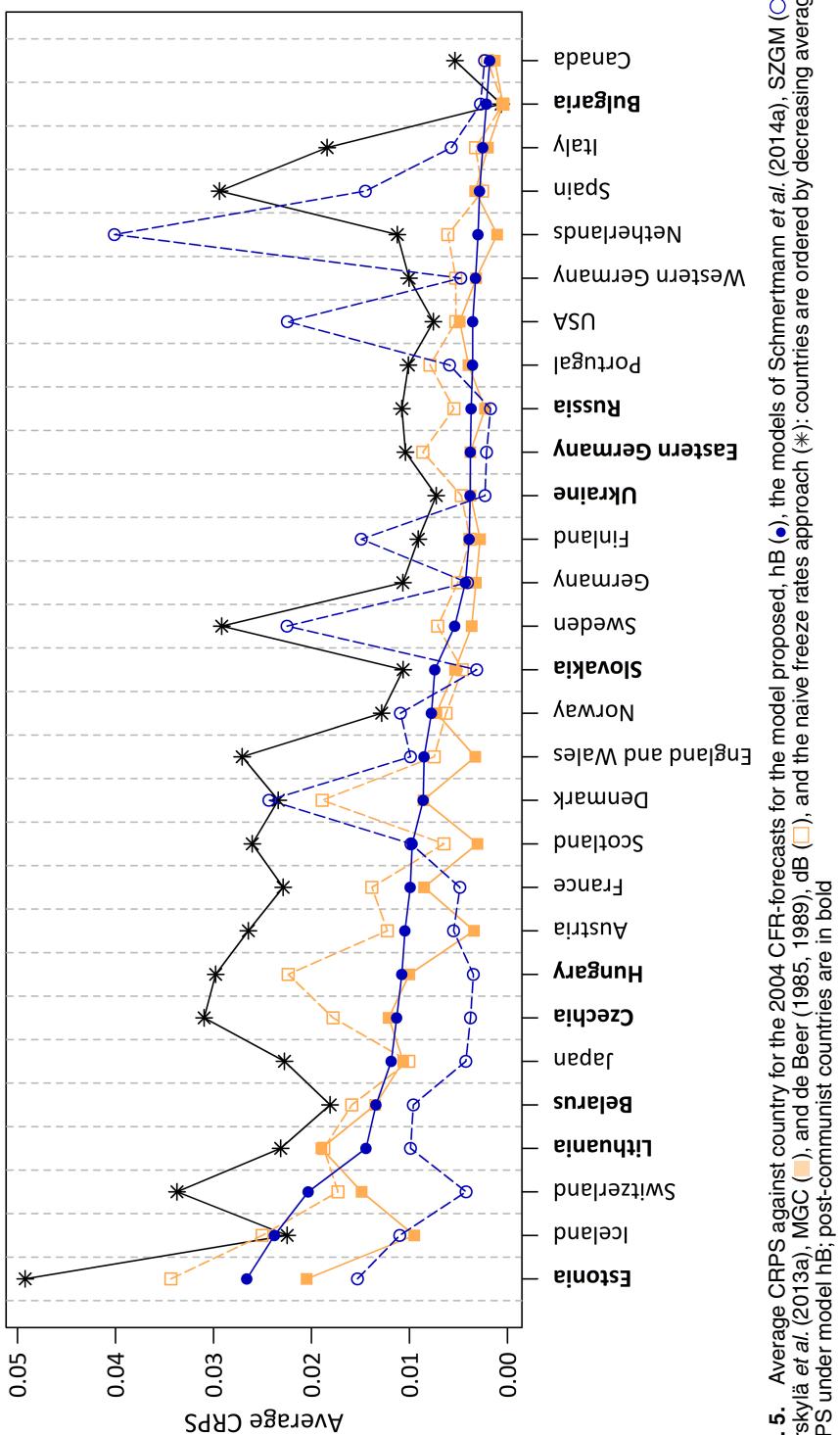
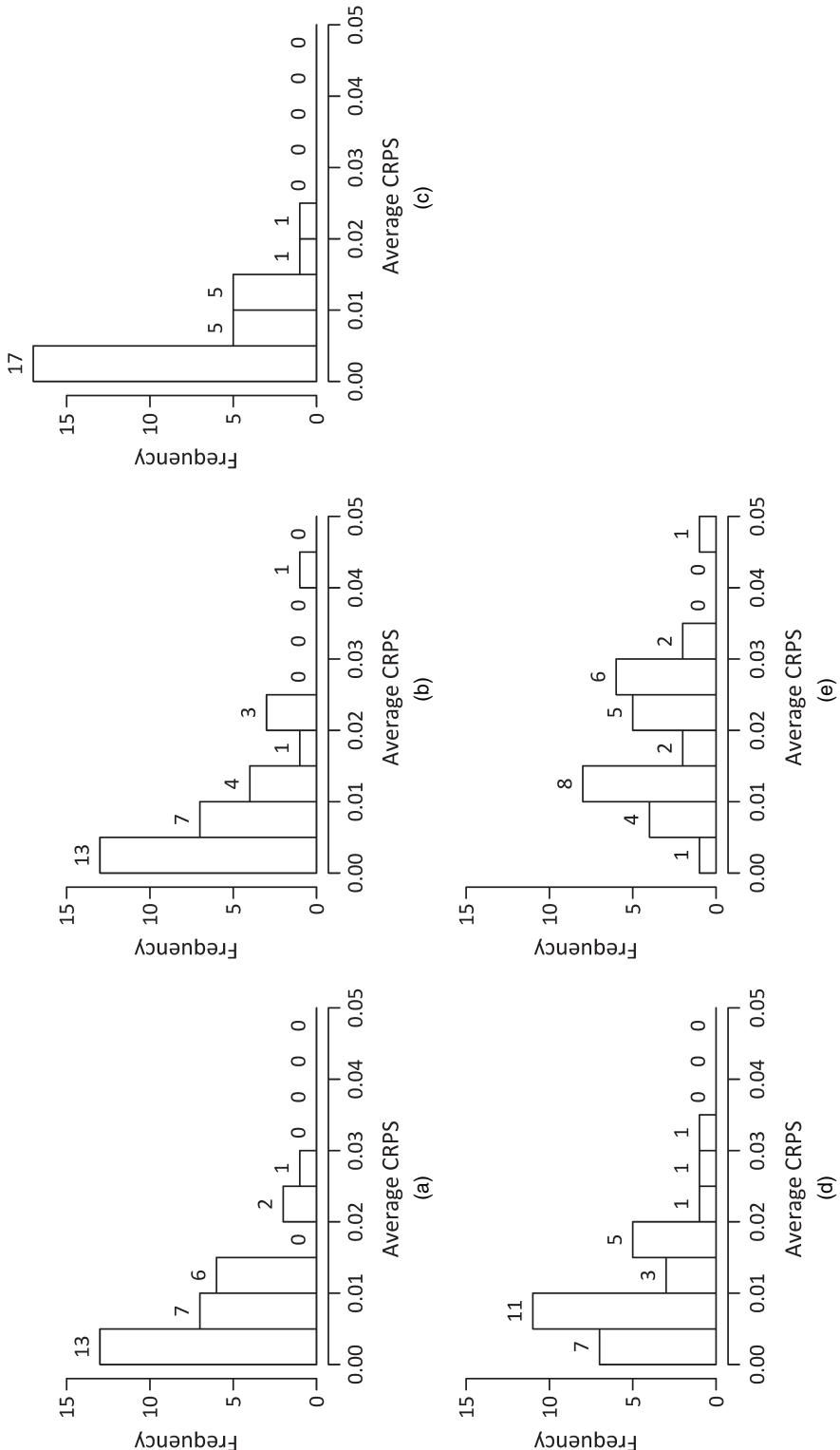


Fig. 5. Average CRPS against country for the 2004 CFR-forecasts for the model proposed, hB (●), the models of Schmertmann *et al.* (2014a), SZGM (○), Myskylä *et al.* (2013a), and dB (□), and the naive freeze rates approach (\*); countries are ordered by decreasing average CRPS under model dB, post-communist countries are in bold



**Fig. 6.** Histograms of the average CRPS across countries for the 2004 CFFR-forecasts for (a) the model proposed, hB, (b) the models of Schmertmann et al. (2013a), SZGM, (c) Myrskylä et al. (2013a), MGC, (d) dB and (e) the naive freeze rates approach

Of the two simple extrapolation models, model MGC clearly outperforms model dB as well as the two Bayesian models, with a highly competitive average score across nearly all the countries; this conclusion is also supported by the relevant histograms in Fig. 6. This strong performance is consistent with the assessment of Bohk-Ewald *et al.* (2018a) in terms of forecast accuracy; however, they did find its forecast uncertainty to be substantially weaker in comparison. From the discussion of scoring rules at the start of this section, we know that the CRPS scores accurate forecast distributions with poor coverage more favourably compared with LogS, which explicitly penalizes forecasts which assign small probabilities to the true value. It is therefore also important to consider how models hB, SZGM and MGC perform under LogS, which we do by presenting Fig. 7 (the equivalent of Fig. 5 for LogS).

The most notable difference between the trends that are exhibited in Figs 5 and 7 is the highly erratic and unpredictable nature of the model MGC scores in Fig. 7—indeed, we cannot display some of the model MGC average LogS-values as they are too large. This dramatically poorer performance for model MGC under LogS compared with the CRPS indicates that its forecast distributions grossly underperform in terms of forecast uncertainty, agreeing with Bohk-Ewald *et al.* (2018a). Regarding the two Bayesian models, again we see evidence of their complementary behaviour in that, where model SZGM performs badly, model hB tends to perform well and vice versa. Under this scoring rule the models appear to be more balanced in terms of the relative magnitudes of their differences from the smallest average score, excluding the poor performance for the Netherlands under model SZGM; this improvement for model SZGM suggests that, overall, model hB may perform slightly worse in terms of coverage but better in terms of forecast accuracy. In Section 3.2.2 we shall investigate whether the summary statistics support these conclusions.

To obtain some idea of what the forecast distributions actually look like, we present the model hB and SZGM forecast distributions graphically for five countries in Fig. 8; note that we do not present the model MGC forecast distributions in these plots for simplicity and because of their pre-established poorer coverage. Figs 8(a) and 8(b) represent countries that perform significantly better under model hB compared with model SZGM according to the scoring rules (the Netherlands and the USA). This is evident from the way that the hold-out CFR-values fall in the centre of the model hB 50% CIs, whereas they drift from the model SZGM intervals after the first few forecast years. Figs 8(c) and 8(d) are where the opposite is true, i.e. model SZGM outperforms model hB in the scoring rules. The first of these, Czechia, is a post-communist (PC) country, and the effect of the declining rates that were experienced across the PC region (Billingsley, 2010) on the forecast CIs is clear. Both models seemingly project the downturn unrealistically into the future, with the wider CIs for model SZGM being the sole reason why it outperforms model hB. The second case for Switzerland is more convincingly in favour of model SZGM, looking like the reverse of Figs 8(a) and 8(b). Lastly, Fig. 8(e) for Scotland gives an instance where there is little to choose between the models under both scoring rules; the CIs overlap for the first few forecast years before diverging, with the subsequent hold-out CFR-values falling roughly in between them. We also observe that the model hB CIs are consistently narrower than those for model SZGM—we shall see whether this difference leads to a reduction in coverage compared with model SZGM when we examine the summary statistics in Section 3.2.2.

Overall, through the use of scoring rules to assess the predictive performance of the model hB cohort fertility forecast distributions relative to some of the current best models in the field, we have found strong evidence to support model hB being competitive in terms of forecast accuracy *and* uncertainty. Fig. 13 in the on-line supporting information gives the CFR-plots for the countries that are not represented in Fig. 8.

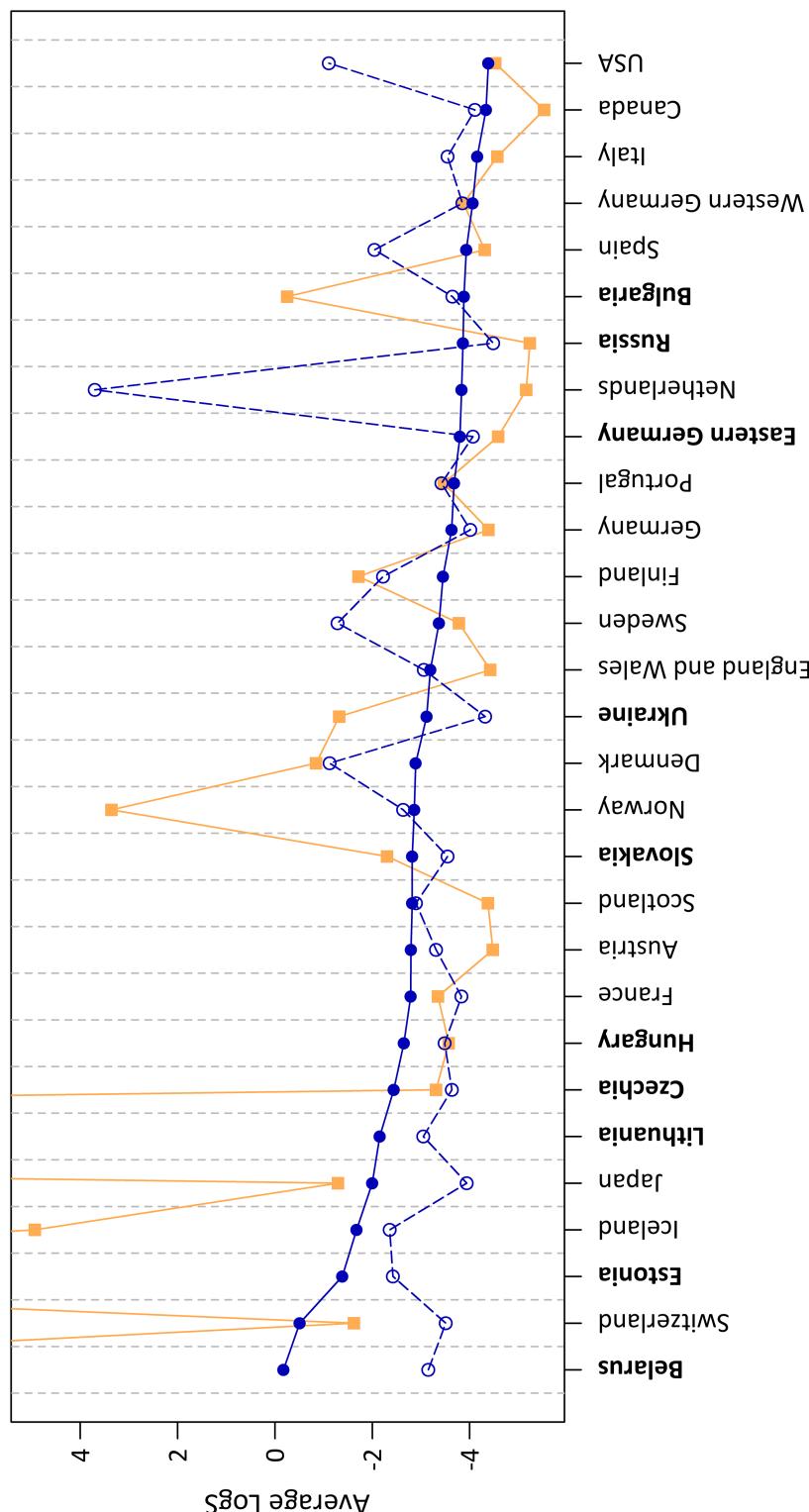
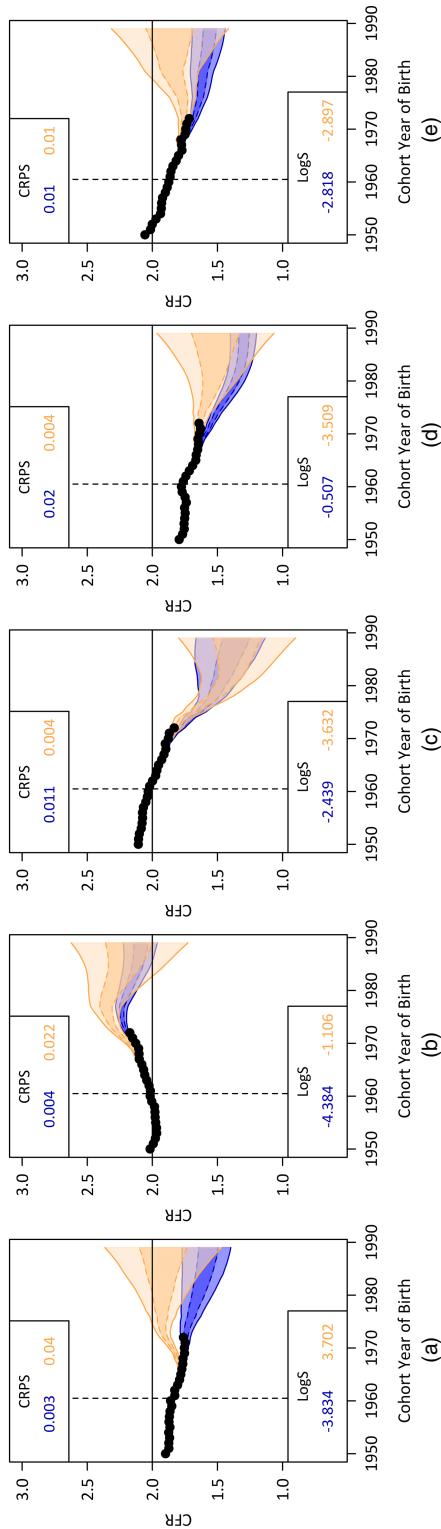


Fig. 7. Average logarithmic score LogS against country for the 2004 CFR forecasts for the model proposed, hB (●), the models of Schmittmann *et al.* (2014a), SZGM (□), and Myrskylä *et al.* (2013a), MGC (○); countries are ordered by decreasing average LogS under model LogS; post-communist countries are in bold



**Fig. 8.** 2004 CFR posterior distributions for selected countries, with average CRPS and logarithmic score LogS for the model proposed, model hB, and the model of Schmidtmann *et al.* (2014a), model SZGM (●, start of the forecast period; ■, model hB 90% CI; □, model hB 50% CI; ▨, model SZGM 50% CI; ▨, model SZGM 50% CI; ●, Human Fertility Database (2019)): (a) the Netherlands; (b) the USA; (c) Czechia; (d) Switzerland; (e) Scotland

**Table 1.** Summary statistics calculated across all countries for the 2004 CFR-forecasts under the model proposed, hB, the models of Schmertmann *et al.* (2014a), SZGM, Myrskylä *et al.* (2013a), MGC, and de Beer (1985), dB, and the naive freeze rates approach†

Measure	Results for the following methods:				
	hB	SZGM	MGC	dB	Freeze rates
MAE (3 decimal places)	0.011	0.013	0.009	0.011	0.020
MAPE (%), 2 decimal places	0.63	0.72	0.49	0.62	1.15
RMSE (3 decimal places)	0.021	0.024	0.016	0.021	0.034
Coverage of 90% CI (%)	76	83	56	—	—
Coverage of 50% CI (%)	58	54	32	—	—

†MAE—mean absolute error; MAPE—mean absolute percentage error; RMSE—root-mean-square error.

### 3.2.2. Summary statistics

The scoring rules in Section 3.2.1 are single metrics that can quantify the overall performance of a forecast distribution in terms of accuracy and uncertainty; next we use various standard summary statistics to assess these two qualities separately, with the results presented in Table 1. For simplicity and ease of interpretation, these statistics are calculated across all the countries rather than being country specific as the average scores were in Section 3.2.1—in this way we can determine whether these results are consistent with our previous general findings.

We have given three measures of predictive accuracy, namely the mean absolute error MAE, mean absolute percentage error MAPE and the root-mean-square error RMSE; note that, for the models with probabilistic forecasts, we compute each error by using the median of the relevant forecast distribution. The MAE-values for the four models (excluding the freeze rates model because its MAE is substantially larger) indicate that the typical magnitude of the CFR-forecast error is 0.01, i.e. 10 children for every 1000 women in a given cohort over their reproductive lives. The freeze rates model actually performs worst by a long way under all three measures, which is not surprising given the analysis of Fig. 5 in Section 3.2.1 and the findings of Bohk-Ewald *et al.* (2018a). Focusing on the four models, we see that model MGC has the lowest value (and therefore best forecast accuracy) for all three statistics; this confirms our conclusions based on the computation of the CRPS in Section 3.2.1. In terms of the ordering of the remaining models, models hB and dB have nearly identical values, all slightly larger than those for model MGC; model SZGM is a little further behind again. This equivalence of models hB and dB does not necessarily follow from Section 3.2.1, where model dB appeared to perform significantly worse than model hB from Fig. 5—however, the fact that the scoring rules are not a measure of forecast accuracy alone (they also take into account forecast uncertainty) could explain this difference.

The reason that model SZGM performs relatively badly under these measures is likely to be because the frequent high scores in Fig. 5 pull up the averages. When we consider the results for just the non-PC countries (which are not presented in Table 1), they actually show a greater margin of improvement for model hB over model SZGM across the three statistics—in particular, MAPE decreases to 0.56% for model hB whereas it increases to 0.85% for model SZGM. Naturally this is countered by the opposite effect being observed for the PC countries (here

MAPE increases to 0.78% for model hB whereas it decreases to 0.45% for model SZGM, the lowest across all the models. So there appears to be evidence that model hB is better suited to forecasting CFR for countries with more stable contemporary fertility histories.

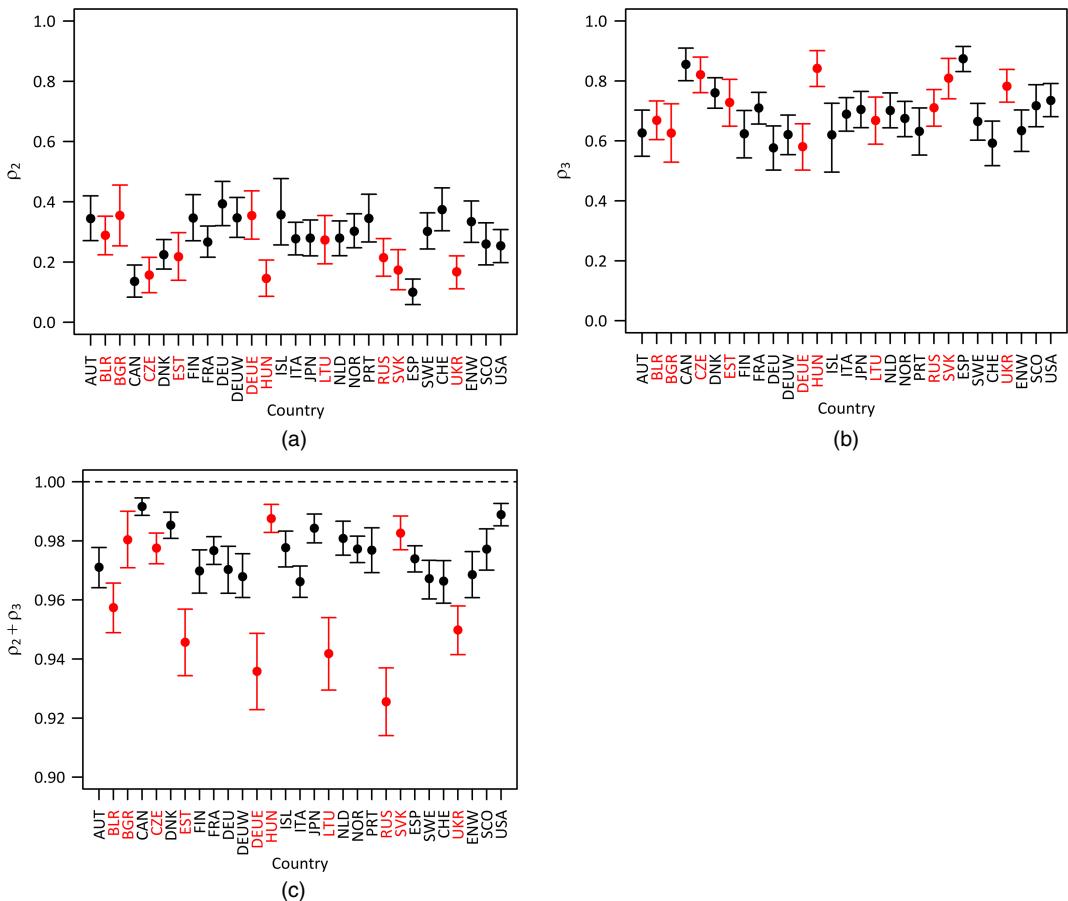
To quantify uncertainty we compute the coverage of the 90% and 50% CIs for the models with probabilistic forecasts (i.e. models hB, SZGM and MGC). The results are consistent with our findings from computing LogS in Section 3.2.1, with model SZGM closest to the nominal values, followed closely by model hB and then model MGC, which has very poor coverage. In terms of the coverages for the non-PC and PC countries separately, we see a similar trend to that observed for the predictive accuracy measures—the model hB coverages exceed those for model SZGM for the non-PC countries (79% versus 77% and 61% versus 44%) but are substantially lower for the PC countries (71% versus 96% and 54% versus 75%). This provides further evidence for our previous conclusion that the forecast performance of model hB is more favourable for the countries without a recent structural shift.

To summarize, the analysis in this section has confirmed our findings from Section 3.2.1, that model hB has forecast accuracy comparable with that of the current best cohort fertility forecasting models—indeed, only model MGC performs significantly better. Conclusions are more difficult to state with forecast uncertainty, as we have only the strong and weak coverage of models SZGM and MGC (established in Bohk-Ewald *et al.* (2018a)) to compare against. Model hB lies in between the two models in this regard but is undoubtedly closer to model SZGM than to model MGC; it is most competitive with model SZGM when considering the non-PC countries alone, where its coverage is slightly higher than that of model SZGM. Overall these results are positive for model hB; however, the poorer performance for PC countries is concerning and should be investigated—we do this in Section 3.2.3.

### 3.2.3. Stationarity and reversion

Lastly, we return to the stationarity discussion in Section 2.2 by presenting the posterior distributions of  $\rho_2$ ,  $\rho_3$  and their sum by country in Fig. 9. First, we note that, for each country, the  $\rho_2$  error bars lie beneath the  $\rho_3$  error bars (comparing Figs 9(a) and 9(b)). This means that the observed time series of age-specific rates in the contemporary Lexis surfaces are telling us that more weight should be put on the freeze slope approach, i.e. following the recent age-specific trends, compared with the freeze rate approach, i.e. remaining at the current age-specific level. Another interesting point is that all the countries choose their age-specific processes to be stationary, which we can see from the  $\rho_2 + \rho_3$  error bars in Fig. 9(c) all lying below 1. This means that our age-specific forecasts all revert to  $[X\beta]_a$  in the long term and so are unlikely to be explosive. Despite this, there does appear to be a difference between the average degree of stationarity, i.e. how close the sum is to 1, for the PC and non-PC countries. The former tend to have their distributions of  $\rho_2 + \rho_3$  at a lower level and hence exhibit a faster reversion.

We illustrate this observation by comparing the age-specific forecasts of Canada and Russia, which are countries with relatively large and small values of this sum respectively, in Fig. 10. The Canada forecasts have reasonably wide CIs with only slight evidence of reversion to the full line for the forecast period that is shown here. Conversely, the Russia forecasts show a very fast reversion and tend to have narrower CIs as a result. This provides some explanation for why our forecasts tend to underperform for the PC countries. It is especially damaging for the Russia forecasts at older ages, where the model hB intervals cannot cope sufficiently well with the continued trends that we observe in the forecast period as a result of this reversion. This seems counterintuitive when we consider that the contemporary data chose to put more weight

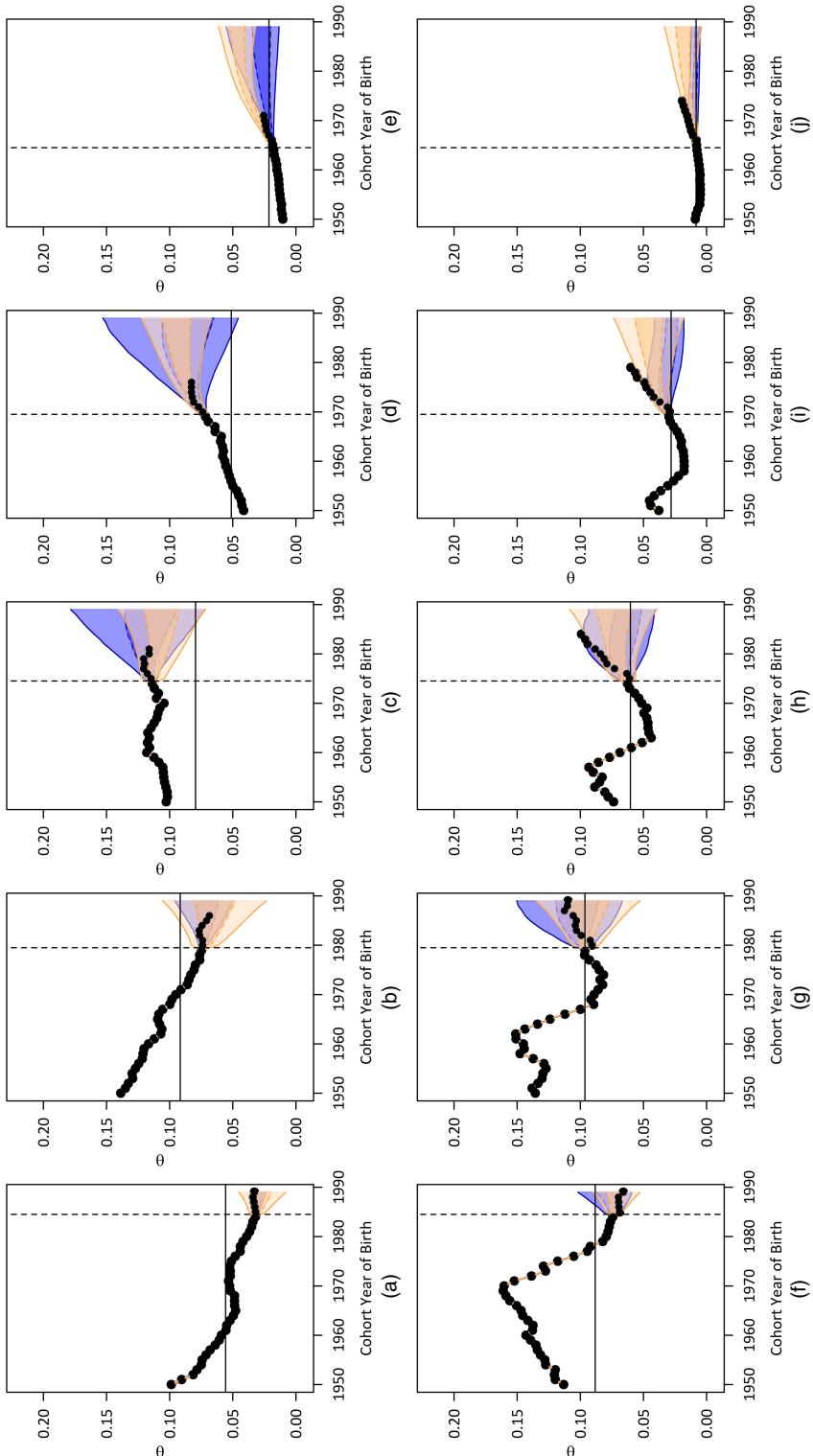


**Fig. 9.** 2004 posterior distribution summary of (a)  $\rho_2$ , (b)  $\rho_3$  and (c)  $\rho_2 + \rho_3$  by country: ●, at the sample median; └──, 90% CI; ━━, PC countries

on following the slope, but instead we have reverted quickly to the current level. Fig. 10 also enables us to assess how successfully we revert to the current level as imposed by the identifiability constraint (see Section 2.2). Canada and Russia present conflicting results, whereby the former has its median reversion values quite far from the current level compared with the latter, where they are much closer. This difference could be due to only constraining three linear combinations of the jump-off error terms to equal 0, causing the level of achievement of the desired reversion to vary across countries.

### 3.3. 2014 fertility forecasts

Following the 2004 forecasts, we now generate 2014 forecasts for ages 15–44 years, contemporary cohorts 1960–1999 and historical cohorts 1904–1953. Although some countries have data as recent as 2017 (e.g. Austria and Hungary), others have data only up to 2013 (e.g. Germany and Ukraine). We choose to use the data that were available in 2014 for these forecasts to ensure that we have 10 or 11 complete contemporary cohorts for each country. We do not use scoring rules to compare the forecasts because there are at most three additional



**Fig. 10.** (a)–(e) 2004 Canada and (f)–(i) Russia fertility forecasts at ages (a), (f) 20, (b), (g) 25, (c), (h) 30, (d), (i) 40 years;  $\vdash$ , start of the forecast period;  $\dashv$ , median revision value for the model proposed, model hB; ■, model hB 90% CI; □, model hB 50% CI; ●, model of Schmertmann et al. (2014a); model SZGM, 90% CI; ■, model SZGM 50% CI; ●, Human Fertility Database (2019)

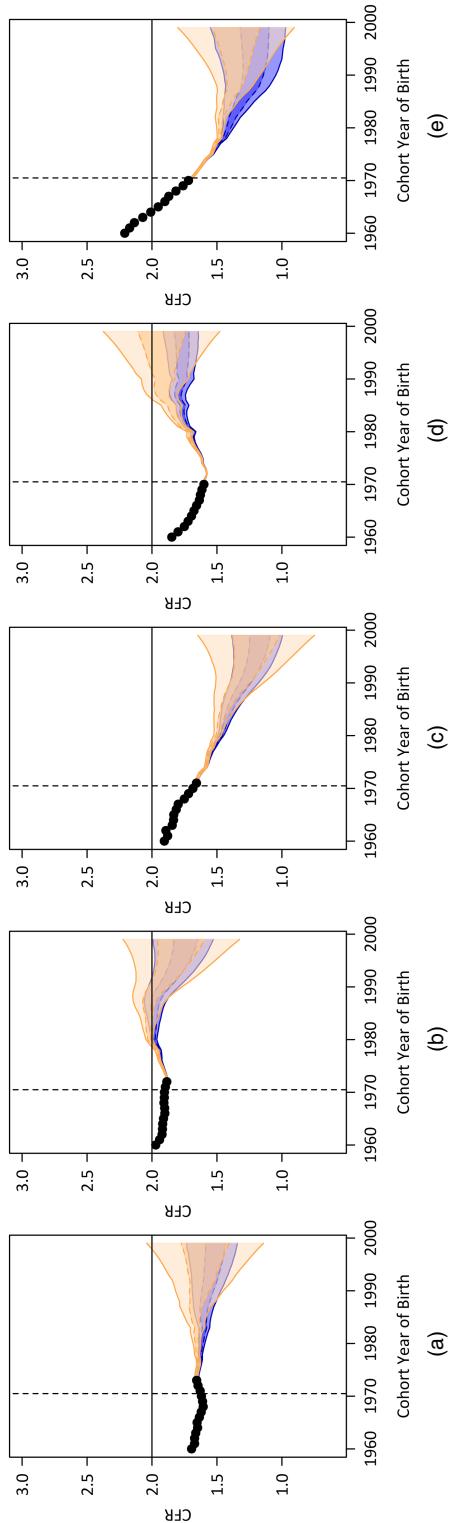
data points for any one country: too few to allow the average to be reliably interpreted. Also, differences in these averages are likely to be negligible because of the minimal uncertainty when completing the first few cohorts. Furthermore, countries with data up to 2013 or 2014 do not have any observed data to apply a scoring rule to, so we would not be able to compare all countries.

With only qualitative comparisons possible, we present the 2014 CFR-forecasts for a range of countries in Fig. 11. Across these plots we see that, as in Fig. 8, the model hB forecasts tend to be more pessimistic and carry less uncertainty compared with the model SZGM forecasts. However, we also note that the forecast distributions consistently overlap at least partly across all 30 countries for which we obtained forecasts (see Fig. 14 in the on-line supporting information for the remaining CFR-plots that are not shown in Fig. 11). This is quite reassuring, as it suggests that the two approaches can make roughly similar inferences about the future based on the identical historical and contemporary data that have been fed into them for each country, albeit processed in different ways.

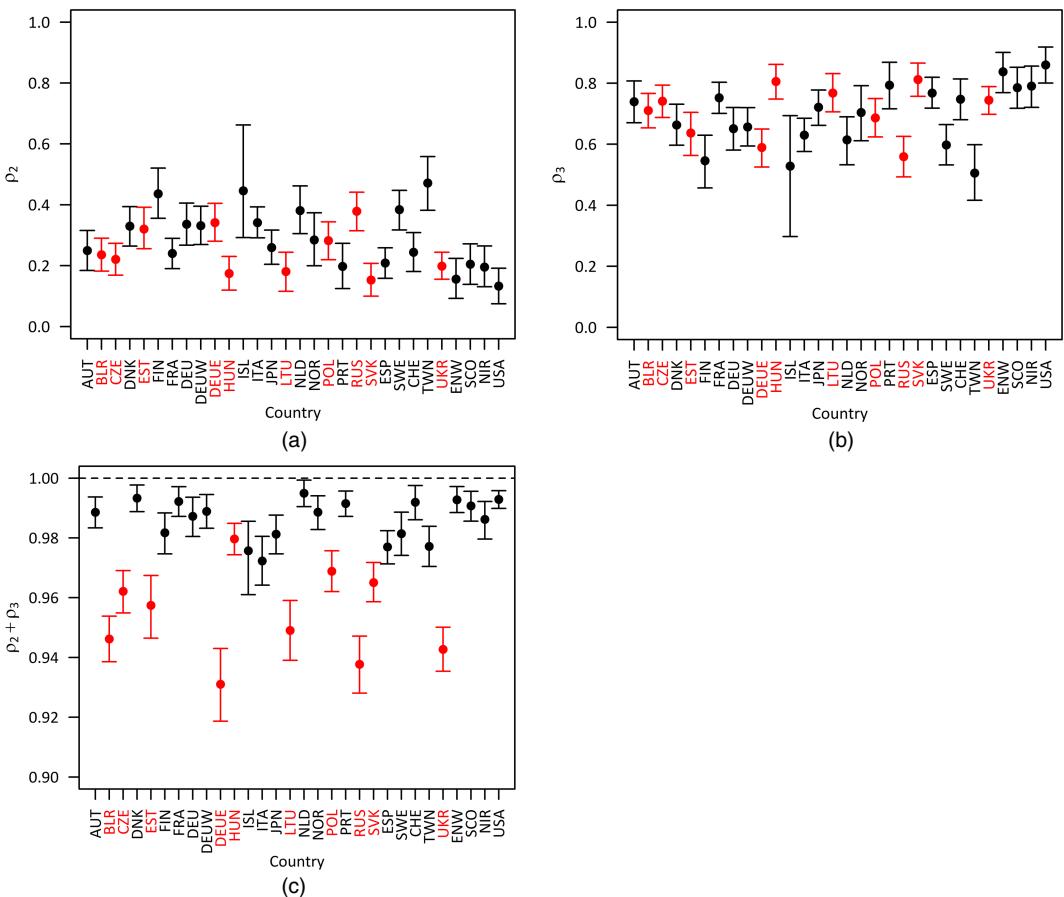
Fig. 11(a) for Austria and Fig. 11(b) for England and Wales have the most stationary observed CFR-values before the forecast period begins. However, whereas the Austria forecast remains without any particular direction, the England and Wales forecast initially shows an increase to replacement level which is then followed by a decline to subreplacement level for the younger cohorts. Figs 11(c), 11(d) and 11(e) for Portugal, Russia and Taiwan respectively are all forecasting from observed CFR-declines in the contemporary Lexis surfaces. Portugal is forecast to continue this decline almost linearly by both models but for Russia the opposite is true, with the forecasts reversing the downward trend; model SZGM in fact forecasts an optimistic return to replacement level. The subtle differences between the observed CFR-values are likely to explain the divergent forecasts for these two countries, namely that Russia shows evidence of the start of an upturn in CFR just as the forecast period begins whereas Portugal does not. Taiwan is an interesting case whereby like Portugal we see an initial continuation of the decline, but then the downturn stabilizes under both models; model SZGM forecasts this to happen slightly earlier than does model hB.

Regarding stationarity, the analogue of Fig. 9 for the 2014 forecasts presented in Fig. 12 again shows a preference for the freeze slope assumption compared with freeze rate. However, whereas for the 2004 forecasts we had no overlap between the error bars in Figs 9(a) and 9(b), here we see substantial overlap for Finland, Iceland and Taiwan. This suggests that, 10 years after the 2004 forecasts, there is some evidence of a move towards stability in the time series of age-specific fertility rates. Fig. 12(c) demonstrates again not only that all the countries are choosing to follow stationary processes as before, but also that the PC countries are still tending to revert faster once the forecast period begins through the comparatively smaller values of  $\rho_2 + \rho_3$ . However, even for a non-PC country such as Taiwan, which has a reasonably high level of  $\rho_2 + \rho_3$ , from the CFR-forecast in Fig. 11(e) we see clear evidence of the reversion kicking in as we move towards the younger cohorts. This is advantageous as it means that our model will not forecast a trend to continue indefinitely, which would be unrealistic.

To summarize, even though it is not easy to draw substantive conclusions from the 2014 CFR-forecasts because of the lack of available validation data, there is a consistent overlap of the model hB and SZGM CIs across countries. Also, the presence of stationarity in the age-specific forecasts for every country, as in the 2004 forecasts in Section 3.2.3, means that we do not need to be concerned about explosive behaviour in our age-specific or CFR-forecasts. Therefore with the little information that we have, the 2014 forecasts appear to be plausible and have a well-calibrated level of uncertainty.



**Fig. 11.** 2014 CFR posterior distributions for (a) Austria, (b) England and Wales, (c) Portugal, (d) Taiwan and (e) Russia and model hB, 90% CI; ■, model of Schmertmann *et al.* (2014a), model SZGM, 90% CI; □, model SZGM 50% CI; ●, Human Fertility Database (2019)



**Fig. 12.** 2014 posterior distribution summary of (a)  $\rho_2$ , (b)  $\rho_3$  and (c)  $\rho_2 + \rho_3$  by country: ●, at the sample median; └──, 90% CI; ━━━━, PC countries

#### 4. Discussion

The aim of this paper is to propose a transparent and intuitive hierarchical Bayesian model (model hB) for forecasting cohort fertility in the spirit of the highly successful model of Schmertmann *et al.* (2014a) that can compete with the current best cohort fertility forecasting models in terms of forecast accuracy and uncertainty. We incorporate our assumptions, which are similar to those made by Schmertmann *et al.* (2014a), explicitly into the model structure through a coherent auto-regressive time series prior for the error terms (see Section 2.2); the resulting hierarchical form of our model also allows the borrowing of strength across the contemporary cohort–age combinations. The precise specification of the prior is determined by the data, which enables us to learn about the relative weights on staying at the current level (freeze rate approach) *versus* following the recent trend (freeze slope approach) for each country, and subsequently dictating the degree of stationarity of the age-specific processes. The presence of stationarity for all countries in both sets of forecasts makes the reversion level very important—our decision for this level to be as close to jump-off as possible was in the spirit of the recent literature (for example, see Myrskylä *et al.* (2013a)). However, this reversion *does* appear to suppress the desire of the data to follow the recent slope in some cases (see Section 3.2.3).

To be considered an important contribution to the literature, it is necessary that our proposed model performs sufficiently well in its purpose, i.e. forecasting age-specific fertility rates and in particular CFR. To determine this, we carry out an extensive validation of model hB by comparing its 2004 CFR-forecasts for 29 countries against those generated from the models of Schmertmann *et al.* (2014a), Myrskylä *et al.* (2013a) and de Beer (1985, 1989): three of the top four models determined by Bohk-Ewald *et al.* (2018a) in terms of forecast accuracy; in addition to this we compare against the naive freeze rates method, which any justifiable fertility forecasting method should be able to outperform easily. We show that, when quantifying the probabilistic accuracy of these forecasts through scoring rules (see Section 3.2.1), there is strong evidence that model hB is highly competitive in terms of forecast accuracy and uncertainty, as well as unquestionably able to outperform the freeze rates method. The calculation of summary statistics regarding point accuracy and coverage in Section 3.2.2 support these conclusions. For the 2014 forecasts a quantitative comparison is not possible; however, the forecasts look reasonable in terms of level and uncertainty (see Section 3.3). So, on the whole, model hB can compete well with the current best models in the field. It is important to note, however, that as the model hB forecasts have been assessed over only a select set of countries at one time point, further validation will be necessary to obtain firmer conclusions regarding forecast performance.

A key advantage of the competing models is their low computational cost, as it takes seconds to produce a fit for one country. Model hB requires Markov chain Monte Carlo methodology and therefore large numbers of iterations are sometimes necessary to obtain results of a suitable quality. This can be computationally expensive; however, thanks to the state of the art Hamiltonian Monte Carlo methods and the RStan software package (Stan Development Team, 2018a), posterior sampling can be conducted with reasonable efficiency. To quantify this, we found that the average fitting time per country for the 2004 and 2014 forecasts were 15 and 21 min respectively, with no country taking longer than 1 h to fit; this is not an unreasonable length of time in practice, especially if it makes the underlying model assumptions more realistic and provides adequate levels of uncertainty. Regarding the simple extrapolation method of Myrskylä *et al.* (2013a), it may perform well in terms of forecast accuracy but does not have well-calibrated levels of uncertainty (see Section 3.2.2); the model of de Beer (1985, 1989) and the freeze rates method do not provide any uncertainty quantification and therefore can have only limited use as deterministic forecasting approaches. Hence, we believe that the use of complex statistical methods in cohort fertility forecasting models such as model hB and the conjugate Bayesian model of Schmertmann *et al.* (2014a) is worth the effort, in response to the question that was posed by Bohk-Ewald *et al.* (2018a). This is especially important if the aim is to obtain long-term fertility forecasts, as long time series of rates are necessary in order to have appropriate uncertainty.

The finding regarding the greater success in forecasting for non-PC countries deserves some discussion. Clearly the contemporary Lexis surfaces of PC countries provide a stronger forecasting challenge due to the sharp decline in the time series of age-specific fertility rates following the regime change. Model hB seems to respond to this by decreasing the combined sum of the weights on the freeze rate and freeze slope approaches, leading to an increased degree of stationarity and therefore a faster reversion in the forecast period (see Section 3.2.3). This suggests that the sensitivity of model hB to recent data could be a disadvantage when there has been a recent structural shift in the country of interest. To decrease this sensitivity, we experimented with allowing  $\rho_2$  and  $\rho_3$  to borrow strength across countries; however, the PC countries appeared only able to tolerate very slight increases in the value of  $\rho_2 + \rho_3$ , which were insufficient to influence the forecasts noticeably. Further investigation into the reasons behind this inconsistency in performance, as well as its reduction, will be required.

As mentioned earlier, further validation of model hB needs to be carried out. This should involve expanding the set of countries that are considered so that we can assess performance in a wider range of circumstances, e.g. high fertility settings. Additionally, multiple forecasts should be generated at various time points to enable quantitative assessment of the performance of model hB across time and in the long term. We also aim to make further attempts to improve the forecast performance of model hB for countries exhibiting a recent structural break. For this, potential avenues to explore are restricting the contemporary data that are used for such countries, incorporating expert opinion and imposing a constraint on the shape of the cohort schedules; the shape constraint suggestion is the one key assumption of the model of Schmertmann *et al.* (2014a) that we did not build into our proposed model. The superior performance of the model of Myrskylä *et al.* (2013a) in terms of forecast point accuracy (see Section 3.2.2) provides evidence that such constraints are not vital to achieve reasonable success; however, it could still be useful to investigate, especially the implications on long-term forecasts. Given the results from this paper, when performing fertility forecasting we recommend fitting a selection of models that have been shown to be competitive in the literature, and fully exploring the causes of any divergences occurring among the forecasts.

In conclusion, our hierarchical Bayesian approach to forecasting cohort fertility is successful through its transparent specification and competitive forecast performance when compared against three of the current best models in the field according to Bohk-Ewald *et al.* (2018a), in particular for countries without a recent structural shift. In addition, it demonstrates how advanced computational methods can be used to fit hierarchical Bayesian models with an atypical set-up. This not only cements the position of hierarchical Bayesian methods at the forefront of population forecasting methods but also makes a valuable contribution to the fertility modelling and forecasting literature.

## Acknowledgements

The doctoral programme of the first author is funded by the Engineering and Physical Sciences Research Council (award reference 1801045). The work of the second and third authors is partly supported by the Economic and Social Research Council Centre for Population Change—phase II (grant ES/K007394/1). The authors thank Jakub Bijak, Jason Hilton and Peter Smith for their feedback during the initial research and writing phases of the project. The authors also acknowledge Carl Schmertmann and the reviewers, who provided helpful comments on earlier versions of this paper.

## Appendix A: Quantifying uncertainty for a conjugate model

The posterior distribution for the conjugate model of Schmertmann *et al.* (2014a) is  $(\theta|y_{\text{country}}) \sim N(\mu_{\text{post}}, \Sigma_{\text{post}})$ , using the same notation as theirs. To incorporate the additional variation described in Section 2.3, we add a modified version of  $\Psi$ , the covariance matrix for the observed rates, to  $\Sigma_{\text{post}}$ . We denote this by  $\Psi^* := \text{diag}_{j=1,\dots,CA}([\mu_{\text{post}}]_j/W_j^*)$ . We modify  $\Psi$  by first extending it to all CA cohort–age combinations, with index  $j$  corresponding to the  $j$ th combination when ordered by age within cohort. We then evaluate the numerator of each entry at its corresponding value of  $\mu_{\text{post}}$  rather than  $y$ , but the denominator  $W_j^*$  remains as the  $j$ th exposure; in the same spirit as Section 2.3,  $W_j^*$  is taken to be its most recently observed value at age  $a$  if unobserved. We then compute the 90% and 50% CIs for the empirical birth rates as  $\mu_{\text{post}} \pm z\sqrt{\text{diag}(\Sigma_{\text{post}} + \Psi^*)}$ , where  $z \approx 1.64$  and  $z \approx 0.67$  respectively. This will have a noticeable effect only for small countries with comparatively low exposures for some cohort–age combinations.

## References

- de Beer, J. (1985) A time series model for cohort data. *J. Am. Statist. Ass.*, **80**, 525–530.

- de Beer, J. (1989) Projecting age-specific fertility rates by using time-series methods. *Eur. J. Popln*, **5**, 315–346.
- Bijak, J. and Bryant, J. (2016) Bayesian demography 250 years after Bayes. *Popln Stud.*, **70**, 1–19.
- Billingsley, S. (2010) The post-communist fertility puzzle. *Popln Res. Poly Rev.*, **29**, 193–231.
- Bohk-Ewald, C., Li, P. and Myrskylä, M. (2018a) Forecast accuracy hardly improves with method complexity when completing cohort fertility. *Proc. Natn. Acad. Sci. USA*, **115**, 9187–9192.
- Bohk-Ewald, C., Li, P. and Myrskylä, M. (2018b) Forecast accuracy hardly improves with method complexity when completing cohort fertility. *R Code*. Max Planck Institute for Demographic Research, Rostock. (Available from <https://github.com/fertility-forecasting/validate-forecast-methods>.)
- Bongaarts, J. and Feeney, G. (1998) On the quantum and tempo of fertility. *Popln Devlpmnt Rev.*, 271–291.
- Booth, H. (2006) Demographic forecasting: 1980 to 2005 in review. *Int. J. Forecast.*, **22**, 547–581.
- Brass, W. (1960) The graduation of fertility distributions by polynomial functions. *Popln Stud.*, **14**, 148–162.
- Brass, W. (1974) Perspectives in population prediction: illustrated by the statistics of England and Wales (with discussion). *J. R. Statist. Soc. A*, **137**, 532–583.
- Bröcker, J. (2012) Evaluating raw ensembles with the continuous ranked probability score. *Q. J. R. Meteorol. Soc.*, **138**, 1611–1617.
- Chandola, T., Coleman, D. A. and Hiorns, R. W. (1999) Recent European fertility patterns: fitting curves to ‘distorted’ distributions. *Popln Stud.*, **53**, 317–329.
- Coale, A. J. and Trussell, T. J. (1974) Model fertility schedules: variations in the age structure of childbearing in populations. *Popln Indx*, **40**, 185–258.
- Czado, C., Delwarde, A. and Denuit, M. (2005) Bayesian Poisson log-bilinear mortality projections. *Insur. Math. Econ.*, **36**, 260–284.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A. and Rubin, D. B. (2014) *Bayesian Data Analysis*, 3rd edn. Boca Raton: Chapman and Hall–CRC.
- Grosi, F. and King, G. (2008) *Demographic Forecasting*. Princeton: Princeton University Press.
- Gneiting, T. and Raftery, A. E. (2007) Strictly proper scoring rules, prediction, and estimation. *J. Am. Statist. Ass.*, **102**, 359–378.
- Good, I. J. (1952) Rational decisions. *J. R. Statist. Soc. B*, **14**, 107–114.
- Hadwiger, H. (1940) Eine analytische Reproduktionsfunktion für biologische Gesamtheiten [An analytic reproduction function for biological groups]. *Scand. Act. J.*, **23**, 101–113.
- Hoem, J. M., Madsen, D., Nielsen, J. L., Ohlsen, E., Hansen, H. O. and Rennermalm, B. (1981) Experiments in modelling recent Danish fertility curves. *Demography*, **18**, 231–244.
- Human Fertility Database (2019) Human fertility database. Max Planck Institute for Demographic Research, Rostock, and Vienna Institute of Demography. (Available from <http://www.humanfertility.org>.)
- Hyndman, R. J. and Ullah, M. S. (2007) Robust forecasting of mortality and fertility rates: a functional data approach. *Computnl Statist. Data Anal.*, **51**, 4942–4956.
- Jasilioniene, A., Jdanov, D. A., Sobotka, T., Andreev, E. M., Zeman, K., Shkolnikov, V. M., Goldstein, J., Nash, E. J., Philipov, D. and Rodriguez, G. (2015) Methods protocol for the Human Fertility Database. (Available from <http://www.humanfertility.org/Docs/methods.pdf>.)
- Jordan, A., Krüger, F. and Lerch, S. (2017) Evaluating probabilistic forecasts with scoringRules. *Preprint arXiv:1709.04743*. University of Bern, Bern.
- Krüger, F., Lerch, S., Thorarinsdottir, T. L. and Gneiting, T. (2019) Predictive inference based on Markov chain Monte Carlo output. *Preprint arXiv:1608.06802*.
- Lee, R. D. (1992) Stochastic demographic forecasting. *Int. J. Forecast.*, **8**, 315–327.
- Li, N. and Wu, Z. (2003) Forecasting cohort incomplete fertility: a method and an application. *Popln Stud.*, **57**, 303–320.
- Matheson, J. E. and Winkler, R. L. (1976) Scoring rules for continuous probability distributions. *Mangmnt Sci.*, **22**, 1087–1096.
- Myrskylä, M., Goldstein, J. R. and Cheng, Y. A. (2013a) New cohort fertility forecasts for the developed world: rises, falls, and reversals. *Popln Devlpmnt Rev.*, **39**, 31–56.
- Myrskylä, M., Goldstein, J. R. and Cheng, Y. A. (2013b) New cohort fertility forecasts for the developed world: rises, falls, and reversals. *Stata Code*. Max Planck Institute for Demographic Research, Rostock. (Available from <https://www.demogr.mpg.de/go/cohort.fertility/>.)
- National Records of Scotland (2019) Uses and limitations of population projections. National Records of Scotland, Edinburgh. (Available from <https://tinyurl.com/ub6hlfz>.)
- Ní Bhrolcháin, M. (2011) Tempo and the TFR. *Demography*, **48**, 841–861.
- Office for National Statistics (2017) National population projections consultation—2016-based national population projections: fertility. Office for National Statistics, Newport. (Available from <https://tinyurl.com/y23dzv4a>.)
- Office for National Statistics (2019a) National population projections quality and methodology information (QMI). Office for National Statistics, Newport. (Available from <https://tinyurl.com/ybj58awe>.)
- Office for National Statistics (2019b) Births quality and methodology information (QMI). Office for National Statistics, Newport. (Available from <https://tinyurl.com/yawcgzya>.)

- Peristera, P. and Kostaki, A. (2007) Modeling fertility in modern populations. *Demog. Res.*, **16**, 141–194.
- Population Reference Bureau (2001) Understanding and using population projections. Population Reference Bureau, Washington DC. (Available from [https://www.prb.org/wp-content/uploads/2001/12/UnderStndPopProj\\_Eng.pdf](https://www.prb.org/wp-content/uploads/2001/12/UnderStndPopProj_Eng.pdf).)
- Schmertmann, C., Zagheni, E., Goldstein, J. R. and Myrskylä, M. (2014a) Bayesian forecasting of cohort fertility. *J. Am. Statist. Ass.*, **109**, 500–513.
- Schmertmann, C., Zagheni, E., Goldstein, J. R. and Myrskylä, M. (2014b) Bayesian forecasting of cohort fertility. *Project Website*. Florida State University, Tallahassee. (Available from <http://schmert.net/cohort-fertility/>.)
- Ševčíková, H., Li, N., Kantorová, V., Gerland, P. and Raftery, A. E. (2016) Age-specific mortality and fertility rates for probabilistic population projections. In *Dynamic Demographic Analysis* (ed. R. Schoen), pp. 285–310. Cham: Springer.
- Shang, H. L. (2012) Point and interval forecasts of age-specific fertility rates: a comparison of functional principal component methods. *J. Popln Res.*, **29**, 249–267.
- Stan Development Team (2018a) RStan: the R interface to Stan. *R Package Version 2.17.4*. Stan Development Team. (Available from <http://mc-stan.org/>.)
- Stan Development Team (2018b) *Stan Modeling Language: User's Guide and Reference Manual*, version 2.18.0. Stan Development Team.
- Wiśniowski, A., Smith, P. W. F., Bijak, J., Raymer, J. and Forster, J. J. (2015) Bayesian population forecasting: extending the Lee-Carter method. *Demography*, **52**, 1035–1059.

#### *Supporting information*

Additional ‘supporting information’ may be found in the on-line version of this article:

‘Supporting information for “Forecasting of cohort fertility under a hierarchical Bayesian approach”’.