

Bayesian estimation of total fertility from a population's age–sex structure

Carl P Schmertmann¹ and Mathew E Hauer¹

¹Center for Demography and Population Health, Florida State University, Tallahassee, FL, USA.

Abstract: We investigate a modern statistical approach to a classic deterministic demographic estimation technique. When vital event registration is missing or inadequate, it is possible to approximate a population's total fertility rate (TFR) from information about its distribution by age and sex. For example, if under-five child mortality is low then TFR is often close to seven times the child/woman ratio (CWR), the number of 0–4 year olds per 15–49-year-old woman. We analyse the formal relationship between CWR and TFR to identify sources of uncertainty in indirect estimates. We construct a Bayesian model for the statistical distribution of TFR conditional on the population's age–sex structure, in which unknown demographic quantities in the standard approximation are parameters with prior distributions. We apply the model in two case studies: to a small indigenous population in the Amazon region of Brazil that has extremely high fertility rates, and to the set of 159 counties in the US state of Georgia. A statistical approach yields important insights into the sources of error in indirect estimation, and their relative magnitudes.

Key words: Demography, fertility, indirect estimation, Bayesian modelling, age–sex distribution

1 A statistical approach to deterministic demographic estimators

Statistical analysis is central to demography. Censuses, surveys, vital records and other sources of demographic information have many imperfections. Understanding and dealing with the resulting uncertainty often requires probabilistic thinking. There is sometimes tension, however, between the statistical foundations of demography and the deterministic nature of many of its methods.

In this article, we demonstrate how re-thinking classic deterministic methods in explicitly statistical terms can yield useful results. We take apart a simple demographic approximation formula for the total fertility rate (TFR, the lifetime average number of children per woman at current age-specific rates). We rebuild the formula with statistical components, in a Bayesian model with TFR as one of several unknown parameters. As we demonstrate by example, an explicitly statistical approach yields

Address for correspondence: Carl P Schmertmann, Center for Demography and Population Health, Florida State University, Tallahassee, FL 32306-2240, USA.
E-mail: schmertmann@fsu.edu

important insights into the sources of error in indirect estimation and their relative magnitudes.

2 Indirect fertility estimation methods

Indirect methods are an integral part of applied demography. They are used to estimate demographic parameters for very small populations or in cases of partial, missing or intentionally masked data (United Nations, 1983; Nolin and Ziker, 2016). Fertility estimation, in particular, has received a significant amount of historical attention from demographers (United Nations, 1983; Arriaga et al., 1994; Brass, 1964; Coale and Trussell, 1974).

Indirect fertility estimation procedures that use age–sex counts as basic inputs are typically regression based. These include the Bogue–Palmore method (Bogue and Palmore, 1964; Palmore, 1978), the Rele method (Rele, 1967) and others (Gunasekaran and Palmore, 1984; Hanenberg, 1983). These techniques estimate TFR from the child–woman ratio (CWR; the number of young children per reproductive-age woman), from measures of marriage prevalence, from mortality indices such as life expectancy and from other proximate determinants of fertility (Bogue and Palmore, 1964).

Regression-based approaches have two main problems. The first relates to geographical scale: coefficients developed from analyses of national populations have historically underperformed when applied to sub-national geographies (Brunsdon et al., 1998; Tuchfeld et al., 1974). The second major problem involves changing relationships between predictors. Relationships between fertility, marriage, infant mortality, income, etc. change over time (Hauer et al., 2013), so that coefficients that initially predict fertility well may eventually have larger errors due to ‘relationship drift’.

Over the second half of the last century, improvements in data quality and availability gradually made indirect estimation less important. As a result, there has been little innovation in indirect fertility estimation techniques in the past 30 years. However, an increasing demand for specialized population estimates (and for small geographic areas in particular) has revived interest in indirect methods (Schmertmann et al., 2013). One newly derived method for indirectly estimating total fertility (Hauer et al., 2013) uses an algebraic rearrangement of the general fertility rate. Called the ‘implied Total Fertility Rate’ (iTFR), it estimates total fertility from age- and sex-specific population counts only. The iTFR is a particularly simple and effective estimator, and we use it as a baseline reference in much of the discussion that follows.

There are many contemporary demographic applications in which complete fertility data is unavailable or is not disaggregated appropriately for the problem at hand. In these cases, indirect estimates from the age–sex structure can be very valuable. Examples include populations without good birth registration systems, populations for which birth information is collected but masked for privacy reasons (such as for many small geographic areas in the United States) or cases in which

vital records do not include information about the mother's membership in a social category of interest (such as religion or income quantile).

In this article, we adopt a statistical approach to indirect estimation of total fertility for cases in which vital records are inadequate, but population counts by age and sex are reliable. The following sections describe the formal mathematical relationships between age–sex structure and TFR, embed those relationships in a Bayesian model with demographic priors, and apply the new model to two very different cases.

3 The formal relationship between total fertility, age–sex structure and the child/woman ratio

3.1 Derivation

Assume that fertility rates (births/woman-year) are positive over the age interval [15, 50) and zero at all other ages. Define

- F_a = average fertility rate over exact ages $[a, a + 5)$, which we call ‘age group a ’
- TFR = total fertility rate = $\sum_a 5 F_a$, where here and elsewhere summation is over reproductive ages $a \in 15, 20, \dots, 45$
- $\phi_a = \frac{5 F_a}{\text{TFR}}$ = fraction of total fertility occurring in age group a
- L_a = expected person-years lived in age group a , in a life table with a radix $l_0 = 1$
- $s = \frac{L_0}{5}$ = expected fraction still alive among children born in the past five years
- W_a = the observed number of women in age group a in a census or survey
- $W = \sum_a W_a$ = the total number of women at childbearing ages [15, 50).

Standard approximations used in cohort-component projection methods (e.g., Wachter, 2014) imply that the expected number of surviving 0–4-year-old children of both sexes per woman in age group $a = 15, 20, \dots, 50$ at the end of a five-year period, which we call K_a , is

$$\begin{aligned} K_a &= \left[\frac{L_{a-5}}{L_a} \cdot F_{a-5} + F_a \right] \frac{L_0}{2} \\ &= \text{TFR} \cdot \frac{L_0}{5} \cdot \frac{1}{2} \left(\frac{L_{a-5}}{L_a} \cdot \phi_{a-5} + \phi_a \right) \\ &= \text{TFR} \cdot s \cdot p_a \end{aligned} \tag{3.1}$$

The K_a term in (3.1) is unusual because it refers to surviving children per woman in age group a at the end, rather than the beginning, of a five-year period. The logic behind the demographic derivation of K_a is otherwise identical to that used to derive standard Leslie matrix terms. We assume that fertility does not begin until age 15, which implies that $\phi_{10} = 0$ in the calculation of K_{15} . The expression for K_{50} includes

$\phi_{45} > 0$, but in practice, K_{50} will always be ignorably small and we simply assume $K_{50} = 0$.

The right-hand side of (3.1) decomposes the expected number of children per woman in age group a into three multiplicative factors. The first two factors, total fertility and child survival, are identical for all age groups. The third factor p_a varies with age; it is an average of the proportions of total fertility experienced in age groups $a - 5$ and a , with a slightly higher weight on the earlier age group to account for the possible mortality of mothers over the previous five years.

In almost all populations, p_a and K_a will near zero for age groups $a = 15$ and $a = 45$, and will reach their highest values for age groups $a = 25$ and 30 . In other words, typical age patterns of fertility mean that women in their late 20s and early 30s are the most likely to have young children.

The expected total number of surviving 0–4-year-old children in a population with W_{15}, \dots, W_{45} women in childbearing age groups 15–19 through 45–49 is

$$C = \sum_a W_a K_a = \text{TFR} \cdot s \cdot \left(\sum_a W_a p_a \right) \quad (3.2)$$

and the CWR is the product of the three terms:

$$\frac{C}{W} = \text{TFR} \cdot s \cdot \left(\sum_a \frac{W_a}{W} p_a \right) = \text{TFR} \cdot s \cdot \bar{p}. \quad (3.3)$$

The third term on the right-hand side of (3.3) is a population-weighted average of p_a . \bar{p} depends on the age pattern of lifetime fertility ($\phi_{15}, \dots, \phi_{45}$), on the current ages of women in the population of interest (W_{15}, \dots, W_{45}) and (to a much lesser extent) on the potential mortality of women in childbearing ages (L_{10}, \dots, L_{45}). It is the average share of lifetime fertility that women in the population experienced over the past five years, after a small adjustment for the possible mortality of adult women.

Rearranged as an expression for TFR, (Equation 3.3) becomes

$$\text{TFR} = \frac{1}{s} \cdot \frac{1}{\bar{p}} \cdot \frac{C}{W} \quad (3.4)$$

so that one can calculate TFR as the product of the CWR and two factors: a *child mortality multiplier* ($1/s$) and an *age-structure multiplier* $1/\bar{p}$. If women of reproductive age are uniformly distributed over seven five-year age groups $a = 15, \dots, 45$, then $\bar{p} \approx \frac{1}{7}$ and the age-structure multiplier is ≈ 7 . However, from Equation (3.3) it is clear that \bar{p} will be higher (and the age-structure multiplier should be lower) if W_a and p_a are positively correlated across ages. In practice, this implies that we should use a higher age-structure multiplier in the $\frac{C}{W} \rightarrow \text{TFR}$ conversion if there are relatively few women in age groups $a = 25$ and 30 , and a lower multiplier if there are relatively many. We investigated the likely variation in this multiplier by calculating $1/\bar{p}$ for 2 054 fertility schedules and populations in the

Human Fertility Database (HFD; Max Planck Institute for Demographic Research and Vienna Institute of Demography, 2016). Multipliers ranged from 5.88 (Taiwan 1985) to 7.92 (Germany 2005); they were within 5% of 7 (6.65–7.35) in 70% of populations, and within 10% of 7 (6.30–7.70) in 96% of populations.

The simplest approximation to Equation (3.4) supposes that child mortality is near zero ($s \approx 1$) and that the age distribution of women is uniform across the reproductive-age groups ($\bar{p} \approx 1/7$). Under those assumptions, $\text{TFR} \approx 7 \cdot \frac{C}{W}$. This approximation is the iTFR measure used by Hauer et al. (2013), who demonstrated that it has smaller average errors across contemporary national populations than the regression-based Bogue–Palmore approach (1964). Fifty years ago, Rele (1967), as cited in Hanenberg (1983), developed a very similar approximation from an empirical study of simulated stable populations: for low-mortality populations, he proposed the estimator $\text{TFR} \approx 7.14 \cdot \frac{C}{W} - 0.06$.

3.2 Interpretation

Proper interpretation of any TFR estimator derived from CWRs requires understanding two fundamental points. First, the measure reflects average fertility over the period during which the children were born (typically the last five years). Second, a woman's demographic category, such as place of residence, marital status or education level, can change during that period. A TFR estimate from the CWR tells us about the *recent* fertility levels of those *currently* in a given category, which may differ conceptually from standard measures.

Estimation of TFR from a population's age–sex structure also depends on accurate census or survey information. Estimates are vulnerable to undercounts and omissions (specifically, to differential undercounts of women and children), to age misreporting and to errors in the reported place of residence or social status. In addition, surviving children must 'stay with' their mothers, in the sense that both must be counted as parts of the same population. The necessary assumptions for estimating TFR from age–sex distributions are commonly satisfied, but researchers must be aware of them.

4 A Bayesian model based on the formal relationship

The relationships in Section 3.1 are deterministic. But even with known vital rates and known numbers of potential mothers, the number of young children (C) is more appropriately modelled as random, especially in very small populations. In addition, demographic quantities such as s and \bar{p} are not truly constants, because fertility age patterns and mortality schedules are not known with certainty. A Bayesian approach to modelling the relationship between TFR, C and $W = (W_{15}, \dots, W_{45})'$ addresses both of these issues.

Figure 1 provides an overview of a Bayesian model for TFR, conditional on C and W . The model has five fundamental parameters (described broadly here and in detail in the following subsections). Scalar TFR and $\beta \in \mathbb{R}^2$ determine the level

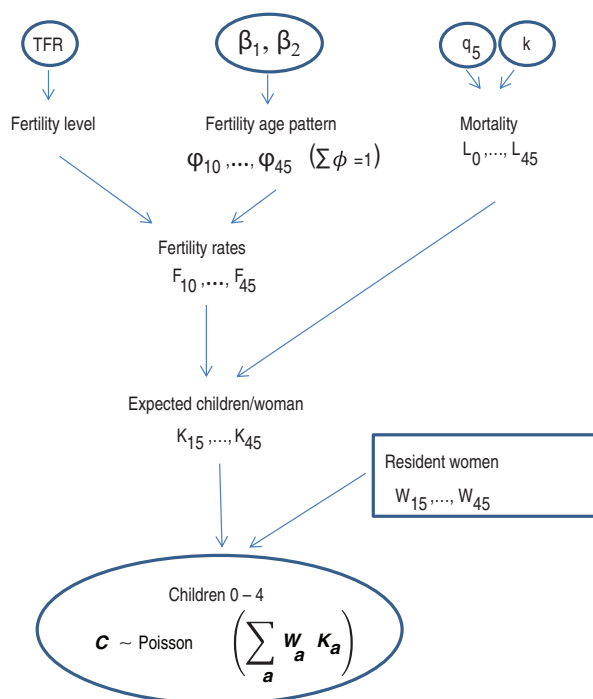


Figure 1 Graphical summary of the Bayesian model relating children 0–4 and women 15–49. Appendix A provides a more detailed tabular version

and age pattern of fertility rates, respectively. Scalar parameters q_5 and k determine age-specific mortality. Vital rates, in turn, determine the expected number of surviving 0–4-year-old children per women in each age group in Equation (3.1), and the expected total number of children in Equation (3.2). The observed number of 0–4 year olds is randomly distributed as a Poisson variable with a mean equal to the expected number. We also provide a tabular version of the complete model in Appendix A.

We specify prior distributions that represent existing demographic knowledge (and uncertainty) about q_5 , β and k . We then examine $P(\text{TFR} \mid C, \mathbf{W})$, the posterior distribution of TFR that arises from the combination of data and priors. This posterior distribution provides quantitative information about which TFR values are more and which are less likely, given the observed numbers of children and women in the population. The following subsections present the model in detail.

4.1 Fertility

4.1.1 Parameters

We separate the schedule of fertility rates for five-year age groups into the level and shape components, as

$$(F_{10}, F_{15}, \dots, F_{45}) = \frac{\text{TFR}}{5} \cdot (0, \phi_{15}, \dots, \phi_{45}) \quad (4.1)$$

where TFR is the total fertility rate and ϕ_a is the proportion of total lifetime fertility that occurs in age group a . We assume that fertility is negligible before age 15 ($F_{10} = 0$). The seven proportions $\phi_{15}, \dots, \phi_{45}$ must sum to 1, so we rewrite them in terms of indices $\gamma_a = \ln(\frac{\phi_a}{\phi_{15}})$ for $a = 15, \dots, 45$, such that $\phi_a(\boldsymbol{\gamma}) = \frac{\exp(\gamma_a)}{\sum_z \exp(\gamma_z)}$.

Finally, we model the $\boldsymbol{\gamma}$ indices as $\boldsymbol{\gamma} = \mathbf{m} + \mathbf{X}\boldsymbol{\beta}$, where $\mathbf{m} \in \mathbb{R}^7$ and $\mathbf{X} \in \mathbb{R}^{7 \times 2}$ are constants derived from empirical data (described in the following section), and $\boldsymbol{\beta} \in \mathbb{R}^2$ are the shape parameters. Thus, three fertility parameters (TFR, β_1, β_2) yield eight five-year fertility rates (F_{10}, \dots, F_{45}): $\boldsymbol{\beta} \rightarrow \boldsymbol{\gamma} \rightarrow \boldsymbol{\phi}$, and $\frac{\text{TFR}}{5} \cdot \boldsymbol{\phi} = \mathbf{F}$ as in Equation (4.1).

4.1.2 Priors for fertility parameters

We use a proper uniform prior for TFR that includes almost no information: $T \sim \text{Uniform}(0, 20)$. For the shape of the fertility schedule by age, we assign higher probability to more typical patterns by building the prior for $\boldsymbol{\beta}$ coefficients from information in the HFD and in the US Census Bureau's International Database (IDB; United States Census Bureau, 2016). In brief, we calculated $\boldsymbol{\gamma}$ indices for a large number of empirical $\{F_a\}$ schedules (226 from the IDB, 411 from the HFD), and then performed a singular-value decomposition on the (de-means) 6×637 $\boldsymbol{\gamma}$ array. This produced a model in which each of the 637 columns of $\boldsymbol{\gamma}$ could be well approximated; the mean vector plus a weighted sum of two principal components $\boldsymbol{\gamma}_i \approx \mathbf{m} + \mathbf{X}\boldsymbol{\beta}_i$. We scaled the two columns of \mathbf{X} so that $\boldsymbol{\beta}_i$ coefficients had zero means, unit variances and zero covariances over the empirical data $i = 1, \dots, 637$. These calculations produced constants $\mathbf{m} = (0 \ 1.39 \ 1.59 \ 1.23 \ 0.45 \ -0.89 \ -3.44)'$ and $\mathbf{X} = \begin{pmatrix} 0 & 0.27 & 0.54 & 0.73 & 0.88 & 1.04 & 1.52 \\ 0 & 0.32 & 0.51 & 0.51 & 0.35 & 0.05 & -0.72 \end{pmatrix}'$, with which we use the prior

$$\boldsymbol{\beta} \sim N(0, \mathbf{I}_2) \quad (4.2)$$

with support restricted to the range $[-2, +2]$ for each $\boldsymbol{\beta}$ coefficient, in order to better mimic the HFD distributions. The examination of the \mathbf{X} matrix shows that, roughly speaking, β_1 affects the mean age of childbearing and β_2 affects the variance. (Higher β_1 means that fertility is postponed, because the constants in the first column of \mathbf{X} are strictly increasing. Higher β_2 means that fertility is more concentrated at ages 20–29, because the largest values of constants in the second column of \mathbf{X} correspond to $a = 20$ and 25.) Figure 2 illustrates a sample of 25 $\boldsymbol{\phi}$ vectors corresponding to random draws from this prior distribution for $\boldsymbol{\beta}$.

4.2 Mortality

4.2.1 Parameters

We model uncertainty in child and adult mortality with the two-parameter relational mortality model developed by Wilmoth et al. (2012). In this model, a mortality

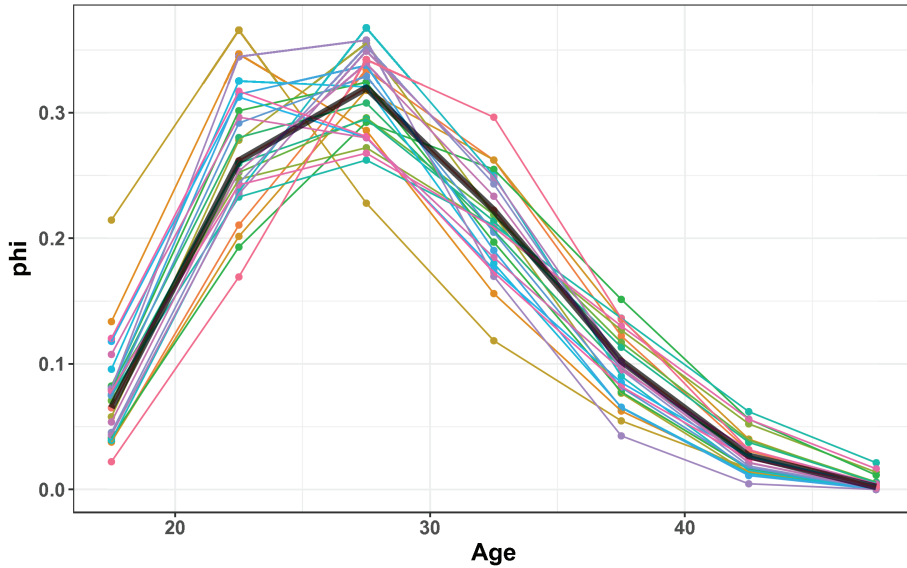


Figure 2 Illustrative ϕ -vectors for proportion of total fertility by five-year age group, based on random draws from prior distribution $\beta \sim N(0, I_2)$. Dark line corresponds to $\beta = 0$

schedule is indexed by the probability of death before age 5 (q_5) and a shape parameter k with typical values between -2 and $+2$. The model uses fixed constants $\{a_x, b_x, c_x, v_x\}$ estimated from schedules in the HMD (University of California, Berkeley (USA) and Max Planck Institute for Demographic Research (Germany), 2016):

$$\ln \mu_x(q_5, k) = a_x + b_x [\ln q_5] + c_x [\ln q_5]^2 + v_x k, \quad x = 0, 1, 5, 10, \dots, 45. \quad (4.3)$$

Mortality rates μ_0 and μ_1 refer to age intervals $[0, 1)$ and $[1, 5)$, respectively; all other rates μ_x refer to five-year age intervals $[x, x + 5)$. (Because $q_5 = 1 - l_5$ is a model parameter, there are no $\{a_1, b_1, c_1, v_1\}$ constants for calculating $\ln \mu_1$. Instead, $\mu_1 = -\frac{1}{4} [\mu_0 + \ln(1 - q_5)]$.)

We convert the log of mortality rates into life table person-years L_a for five-year intervals $[a, a + 5)$ using standard demographic approximations. Survival probabilities to exact ages are $l_0 = 1$, $l_1 = e^{-\mu_0}$, $l_5 = l_1 \cdot e^{-4\mu_1}$ and $l_x = l_{x-5} \cdot e^{-5\mu_{x-5}}$ for $x = 10, \dots, 45$. Life table person-years are $L_0 = \frac{1}{2} (l_0 + l_1) + \frac{4}{2} (l_1 + l_5)$ and $L_a = \frac{5}{2} (l_a + l_{a+5})$ for $a = 5, \dots, 45$. Thus, two mortality parameters (q_5, k) yield ten L_a values $(L_0, L_5, \dots, L_{45})$: $(q_5, k) \rightarrow \ln \mu \rightarrow l \rightarrow L$.

4.2.2 Priors for mortality parameters

We assume that there are one or more external estimates of q_5 , denoted \hat{q}_5 . We use the prior

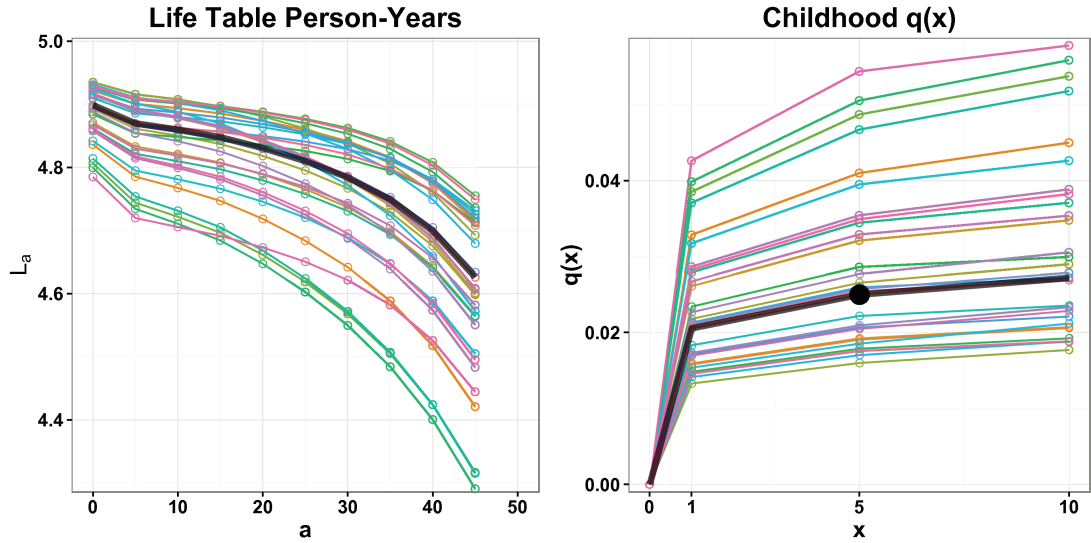


Figure 3 Illustrative draws from mortality priors when $\hat{q}_5 = 0.025$. Left panel shows calculated L_a values for 30 random draws from the joint (q_5, k) distribution. Right panel shows the corresponding child non-survival probabilities. Dark lines in both panels correspond to $(q_5, k) = (0.025, 0)$. Solid point in the right panel represents the estimate $\hat{q}_5 = 0.025$ on which the prior for q_5 is based

$$q_5 \sim \text{Beta} [a(\hat{q}_5), b(\hat{q}_5)] \quad (4.4)$$

where $a(\hat{q}_5)$ and $b(\hat{q}_5)$ are chosen so that $P[q_5 < \frac{1}{2} \min(\hat{q}_5)] = P[q_5 > 2 \max(\hat{q}_5)] = 0.05$. This conservatively assigns a 90% prior probability that under-five mortality q_5 is between one-half the minimum estimate and twice the maximum estimate. (If \hat{q}_5 is a scalar, then $\min(\hat{q}_5)$ and $\max(\hat{q}_5)$ are identical).

For the (much less influential) shape parameter k , we use the prior

$$k \sim N(0, 1) \quad (4.5)$$

which centres the distribution on zero and has a low probability of falling out of the typical $[-2, +2]$ range. We assume that the mortality parameters q_5 and k are independent. Figure 3 illustrates random draws from the joint prior for (q_5, k) .

4.3 Expected number of surviving children

Specific values of parameters (TFR, β, q_5, k) imply specific values K_a in (3.1). The expected number of surviving children for the W_a women observed in age group a is $W_a K_a$, and the observed number of their surviving children may be modelled as $C_a \sim \text{Poisson}(W_a K_a)$. It is reasonable to assume that C_a values are statistically independent, conditional on fertility and mortality rates, so that their sum

$C = \sum_a C_a$ is also a Poisson random variable. Thus,

$$C \mid \text{TFR}, \boldsymbol{\beta}, q_5, k \sim \text{Poisson} \left[\sum_a W_a K_a(\text{TFR}, \boldsymbol{\beta}, q_5, k) \right] \quad (4.6)$$

4.4 Posterior distribution of TFR

For TFR values in the permissible range $[0, 20]$, the posterior for parameters conditional on data is

$$P(\text{TFR}, \boldsymbol{\beta}, q_5, k \mid C) \propto L(C \mid \text{TFR}, \boldsymbol{\beta}, q_5, k) f_{\boldsymbol{\beta}}(\boldsymbol{\beta}) f_q(q_5) f_k(k) \quad (4.7)$$

where the likelihood on the right-hand side is the Poisson likelihood in Equation (4.6), and the f functions represent the prior densities implied by Equations (4.2), (4.4) and (4.5), respectively. The marginal posterior for TFR, which expresses the relative probabilities of alternative fertility levels, given the number of children C and the counts of women W_{15}, \dots, W_{45} , is

$$P(\text{TFR} \mid C) \propto \int L(C \mid \text{TFR}, \boldsymbol{\beta}, q_5, k) f_{\boldsymbol{\beta}}(\boldsymbol{\beta}) f_q(q_5) f_k(k) d\boldsymbol{\beta} dq_5 dk \quad (4.8)$$

In practice, we sample from the full posterior distribution in Equation (4.7) by applying Markov Chain Monte Carlo (MCMC) methods. Specifically, we programmed the model in the *Stan* MCMC language (Carpenter et al., 2017), as implemented in the *rstan* package in *R* (Stan Development Team, 2016; R Core Team, 2016). We use the empirical density of the sampled TFR values to estimate the marginal posterior of TFR in Equation (4.8). Experiments reported in Appendix B show that this posterior distribution is insensitive to the choice of priors about age patterns of fertility and mortality.

5 Example 1: A small indigenous population in the Brazilian Amazon

In 2010, the Kanamari do Rio Juruá Indigenous Territory in the Brazilian state of Amazonas (Terras Indígenas, 2017) had $C = 191$ resident children under age 5, and $\mathbf{W} = (40, 34, 29, 19, 14, 9, 8)'$ resident women in five-year age groups 15–19 through 45–49. Note that this population is very small (only 153 women of childbearing age), and that Kanamari women of childbearing age tend to be young (nearly half are under age 25). This data comes from geographically detailed population counts available online from the Brazilian census bureau (Instituto Brasileiro de Geografia e Estatística, 2016); the population counts mentioned earlier are totals for specific census sectors comprising the Kanamari do Rio Juruá territory.

The Kanamari CWR was $191/153 = 1.25$, so the iTFR estimate is $7 \cdot \frac{191}{153} = 8.74$.

The iTFR calculation does not account for child mortality, or for the concentration of 15–49 year-old Kanamari women in the young age groups.

The Kanamari territory covers approximately 6 000 km² and includes parts of four different municipalities in the Amazonas state. Estimated q_5 for the total populations (indigenous and non-indigenous) in those municipalities ranges from 23 to 32 per 1 000 (United Nations Development Program, 2013). On the basis of these external estimates, we use the prior $q_5 \sim \text{Beta}(3.99, 114.26)$ so that q_5 has a 90% prior probability of lying in [11.5, 64] per 1 000.

MCMC sampling of 2 000 values from the posterior distribution of TFR in (4.8), conditional on the Kanamari data for (C, W) , produces the data illustrated in Figure 4. The iTFR = 8.74 estimate is close to the posterior median of 8.49. The small size of the population, together with uncertainty about the level and pattern of mortality (especially childhood mortality) and about the age pattern of fertility, imply considerable uncertainty about the TFR. *A posteriori* there is a 50% probability that Kanamari TFR in the five-year period preceding Brazil's 2010 census was between 7.87 and 9.20 (the interquartile range of the posterior in Figure 4), and an 80% probability that the TFR was between 7.36 and 9.85 (the interval between the 10th and 90th percentiles).

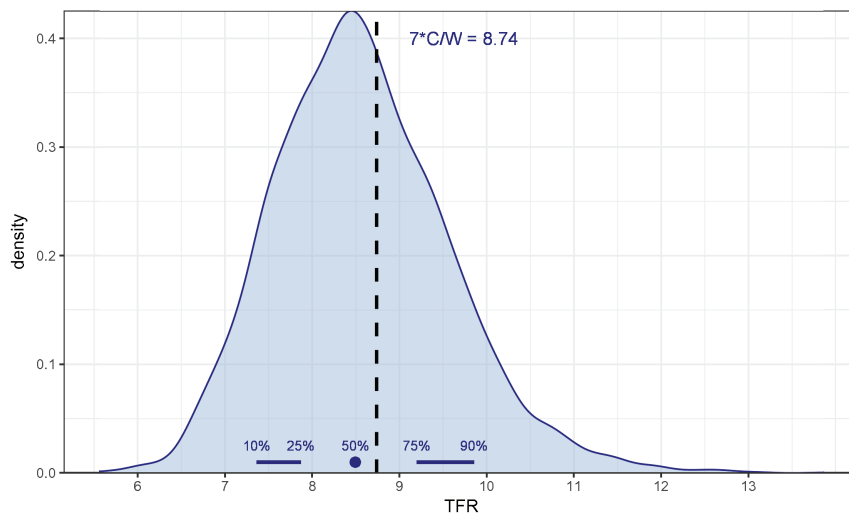


Figure 4 Posterior density of TFR for the Kanamari do Rio Juruá Indigenous Territory. Dashed vertical line corresponds to $\text{iTFR} = 7 \cdot \frac{C}{W}$ estimate from the child/woman ratio in the territory. Selected percentiles of the posterior distribution are marked near the horizontal axis

There are several important conclusions from this example. First, the unusual ratio of young children to reproductive-age women indicates that Kanamari TFR is extremely high. Second, the combination of a very small population with significant

uncertainty about important demographic parameters (especially child mortality levels) means that uncertainty about Kanamari TFR is also very high. Third and most importantly, a Bayesian approach to estimation automatically generates quantitative statements about what is and is not known about this population's fertility level. With a CWR of 191/153, Kanamari TFR is almost certainly higher than 7 (96% posterior probability), very likely higher than 8 (70%), possibly higher than 9 (31%), but almost certainly less than 11 (98%).

6 Example 2: 159 counties in Georgia

As a second example, we consider estimating total fertility from the numbers of children and women across a large number of related populations. Specifically, we use 2010 census population counts by sex and age group to estimate TFR for 159 counties in the US state of Georgia. These counties vary significantly in population size, but all are larger than the indigenous population in the previous example. The smallest population (Taliaferro County) had 335 resident women 15–49, and eight counties had fewer than 1 000 women in this age group. The largest county (Fulton, in the metropolitan Atlanta area) had more than 250 000 women 15–49, and 37 counties had more than 10 000.

The model structure is the same as that illustrated in Figure 1 and described earlier, with two exceptions. First, we add a hierarchical structure in which county fertility levels $\text{TFR}_1, \dots, \text{TFR}_{159}$ are independent draws from a common distribution with state-level parameters μ and σ :

$$\text{TFR}_i \sim N(\mu, \sigma^2) \quad (6.1)$$

This additional assumption means that we believe *a priori* that a set of TFR estimates ($\text{TFR}_1, \dots, \text{TFR}_{159}$) is more likely if the 159 values are similar. The practical effect of the hierarchical assumption is to shrink TFR estimates for counties with very low populations towards estimates in other counties. In other words, it will allow small counties like Taliaferro to ‘borrow strength’ from other locations that we believe—via (6.1)—to be similar. Second, we use distinct mortality priors for each location. Georgia public health data (<https://oasis.state.ga.us>) allows the calculation of county-level \hat{q}_5 estimates, which range from a low of 5 per 1 000 (Forsyth County) to a high of 24 per 1 000 (Quitman County). We used these estimates to construct a separate mortality prior for each county via Equation (4.4).

In the Georgia application, each parameter in Figure 1 becomes 159 separate parameters, one per county. The complete model therefore has 797 parameters: μ , σ , and $(\text{TFR}, \beta_1, \beta_2, q_5, k)_{i=1, \dots, 159}$. The data is similarly expanded: $(C, W_{15}, \dots, W_{45})_{i=1, \dots, 159}$. With the exception of TFR, we assume independence of all demographic parameters across counties. (We assume prior independence mainly to simplify exposition. Researchers could easily add spatial or social priors that assert higher probability for parameter sets in which geographically adjacent or socially similar areas had similar values).

6.1 Results for Georgia counties

We generated posterior draws from the 797-vector of parameters with MCMC sampling. We then examined the marginal distributions of each of the 159 county TFR levels. Figure 5 illustrates these marginal distributions in condensed form, with TFR on the horizontal scale and counties stacked vertically. County TFRs are centred around a global mean of $\mu = 2.19$ (the posterior mean of μ). Using posterior medians (the dots in Figure 5) as point estimates, TFR estimates from CWRs exhibit a considerable range—from a low of 1.42 at the bottom of Figure 5 (Clarke County, site of a large state university) to a high of 2.72 at the top (Chattahoochee County, site of a military base). Most county TFR estimates, however, fall within a fairly narrow range, with 127 (80%) of the estimates between 1.9 and 2.4.

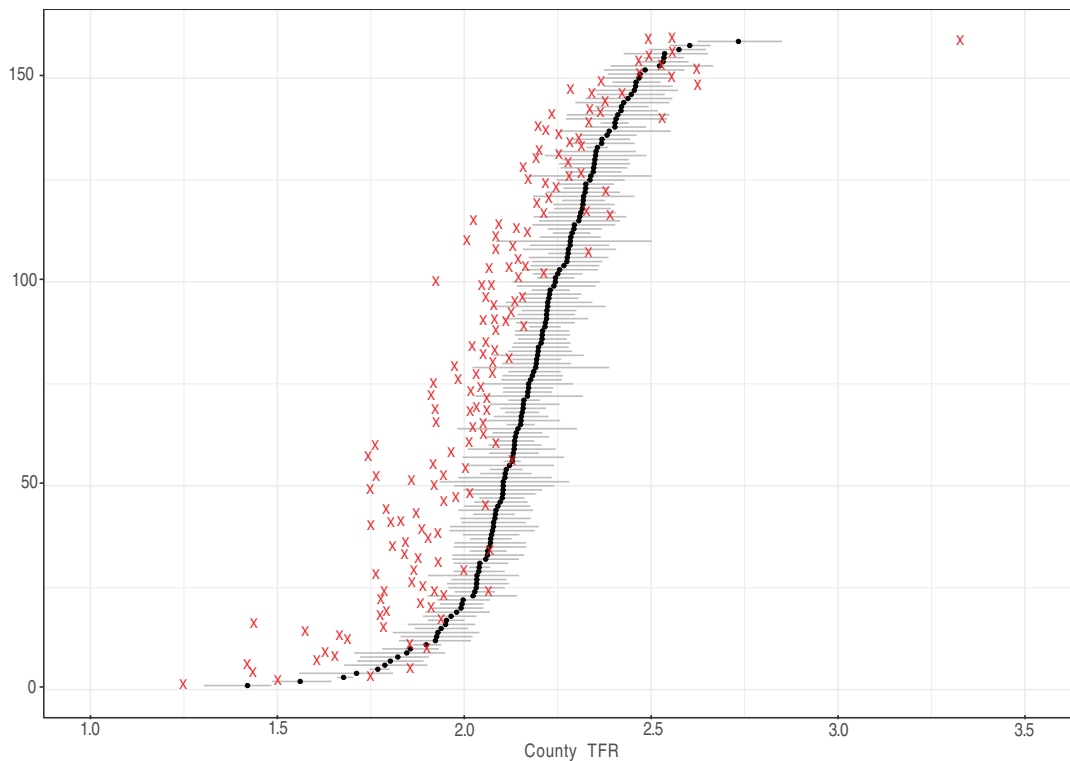


Figure 5 Posterior Distributions for TFR in 159 Georgia counties, 2010. Counties are sorted in ascending order of posterior median (dots). Horizontal bars represent 80% posterior probability intervals (10th to 90th percentile). Small xs represent iTFR estimates $7 \cdot \frac{C}{W}$

Figure 6 displays the posterior median TFRs in geographical form. As suggested by the relatively tight distribution of medians in Figure 5, there are no strong spatial gradients. However, there are a few clear spatial patterns. Fertility is generally a bit higher in the southern and southeastern part of the state, where there is a cluster of

counties with levels such as 2.4, 2.5 or 2.6. The Atlanta metropolitan area in the north-central region has lower fertility, indicated by a cluster of three counties with $\text{TFR} < 2$.

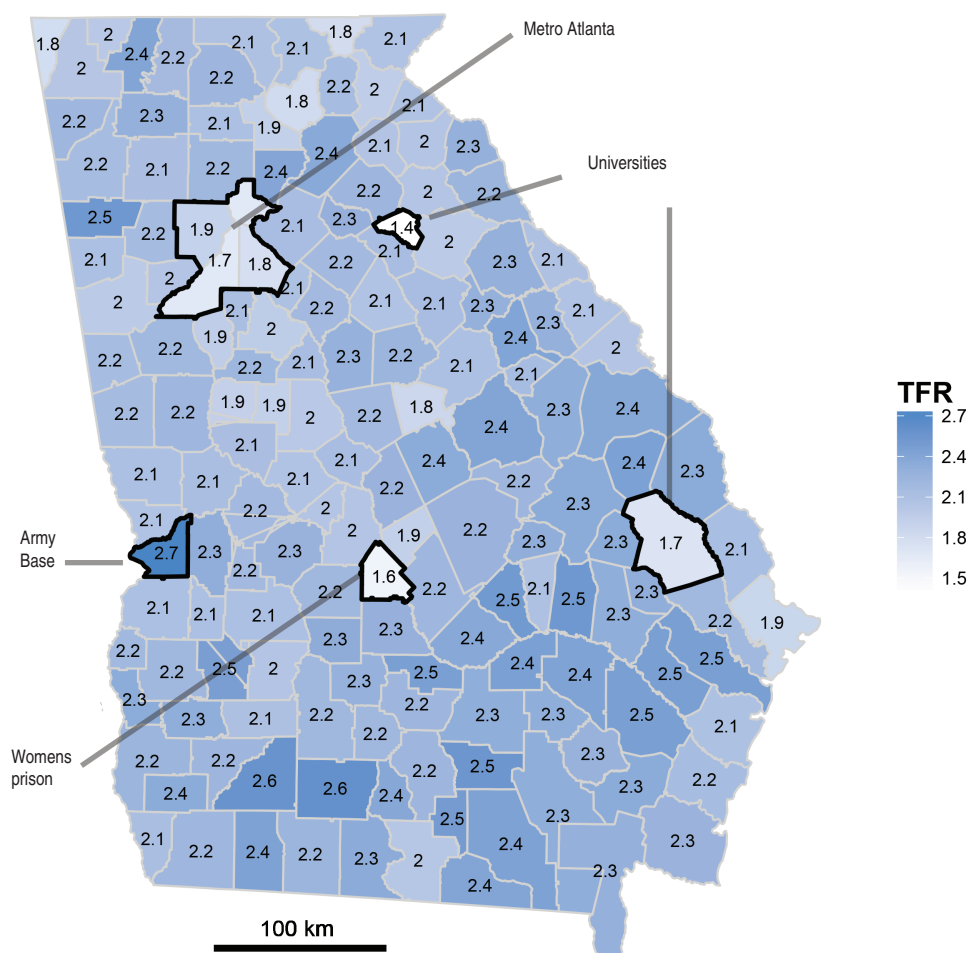


Figure 6 Map of posterior Median TFR, Georgia counties 2010

6.2 Analysis of Georgia results

6.2.1 Signal and noise in small-area estimates

Even in a low-mortality environment with moderately large populations, sampling variability and uncertainty about the age pattern of childbearing sometimes imply substantial uncertainty in TFR estimates based on CWRs. The widths of the 80%

credibility intervals in Figure 5 range from 0.04 for Fulton County [1.66, 1.70] to 0.42 for Taliaferro County [2.08, 2.50]. Importantly, differences between county-level estimates are much smaller than the uncertainty in those estimates. Thus, it would be very difficult to rank counties accurately by the fertility level.

6.2.2 The importance of age structure

Results from the Bayesian model also demonstrate the quantitative importance of accounting for age structure when estimating fertility levels from counts of women and children. From Figure 5, it is clear that $iTFR = 7 \cdot \frac{C}{W}$ is too low in most counties. Almost all iTFR estimates are below the Bayesian posterior median, and in many cases they fall below the 10 percentile of the posterior distribution. These underestimates occur because women in most counties are not uniformly distributed across age groups $a = 15, \dots, 45$. Figure 7 illustrates, showing the fraction of 15–49-year-old women in each county, who are in the high- K age groups $a = 25$ and 30 (horizontal axis) against the $\frac{C}{W}$ multipliers that would be needed to reproduce the Bayesian posterior medians.

Concern about TFR estimates from CWRs has centred on bias caused by child mortality. For example, International Union for the Scientific Study of Population (2017) prominently warns that the CWR is not an accurate indicator of fertility because it can have a substantial downward bias when child mortality is high. Equation (3.4) makes the same point—for example, if 10% of newborns die before age five ($q_5 = 0.10$) then accurate estimation of TFR requires multiplying C/W by approximately $7/0.9 = 7.78$, rather than 7.00. But $q_5 = 0.10$ is a very high contemporary mortality rate, seen only in sub-Saharan Africa, Afghanistan and a handful of other countries. The information in Figure 7 illustrates that the age structure within the female population is a much more important source of potential bias when translating $C/W \rightarrow TFR$. Furthermore, the age-structure bias can go in either direction. Using the Bayesian posterior median as a point estimate of TFR, appropriate C/W multipliers for Georgia counties vary from 5.75 (Chattahoochee) to 9.50 (Fayette). This range of variation is far wider than that caused by plausible differences in child mortality, and it is due to county differences in the age distribution of women within the childbearing ages.

6.2.3 Comparison to vital statistics

Georgia publishes vital statistics reports of birth rates by county and mother's age (<https://oasis.state.ga.us>). This allows the comparison of results from the Bayesian hierarchical model with published estimates. It is not a comparison with 'true' local rates, however, because measured rates in small populations are also subject to coincidental sampling errors. Nevertheless, it is useful to compare estimates from the two methods.

Figure 8 shows county-level TFR estimates from vital records over 2006–2010, and posterior distributions from the Bayesian model based on the 2010 county age pyramids. The left-hand side of the figure contains data from 110 counties that had fewer than 15 resident thousand woman of childbearing age; the right-hand panel

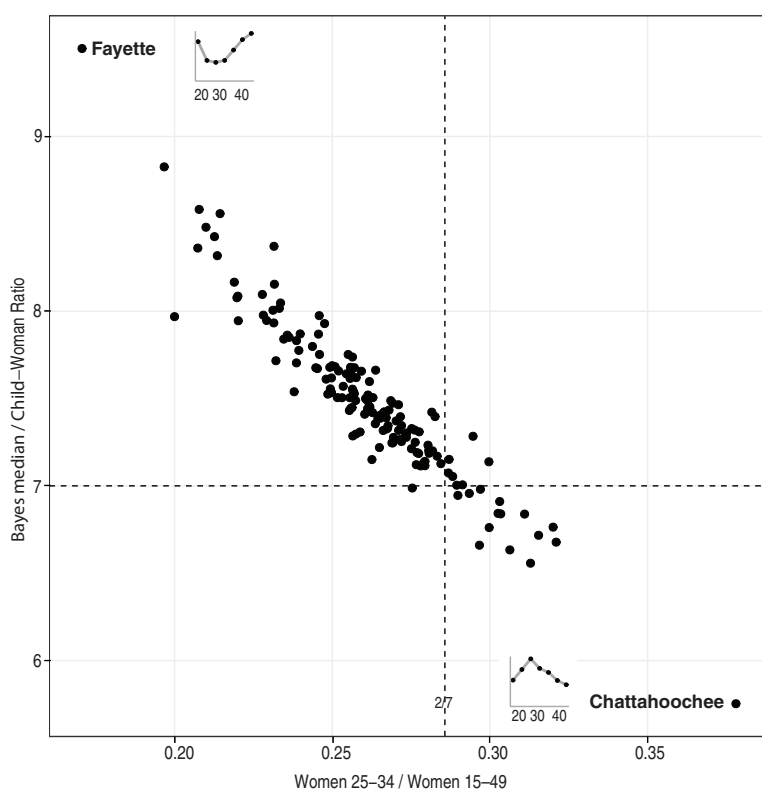


Figure 7 Estimated ratio of TFR to children per woman, as a function of local age structure, Georgia counties 2010. Horizontal axis represents fraction of women 25–34 among all women of childbearing age. Vertical axis represents ratio of posterior median TFR to the county’s child–woman ratio. A uniform distribution of women across age groups implies that $2/7$ of women are 25–34 and $TFR \approx 7 \frac{C}{W}$. Most counties have a lower proportion 25–34 and thus $TFR > 7 \frac{C}{W}$. Small insets illustrate the local age structure in the two most extreme counties, Fayette and Chattahoochee

contains data for 37 counties with more than 15 000 resident women in these age groups. (12 counties were so small that they had missing rates and TFRs even after accumulating five years of data).

For the small populations, the Bayesian model exhibits a classic shrinkage effect: counties with low TFR from vital statistics tend to have higher Bayesian estimates, and vice-versa. This regression-to-the-mean effect is desirable, because sampling errors affecting TFR are more likely to be negative in counties with low estimates and positive in counties with high estimates.

In the large populations in the right-hand panel of Figure 8, both vital statistics and Bayesian TFR estimates have higher precision. Consequently, there is almost no shrinkage effect from the hierarchical model for the large counties.

There is a clear tendency for Bayesian estimates to be lower than vital statistics estimates in the large counties. These differences are especially notable in the four

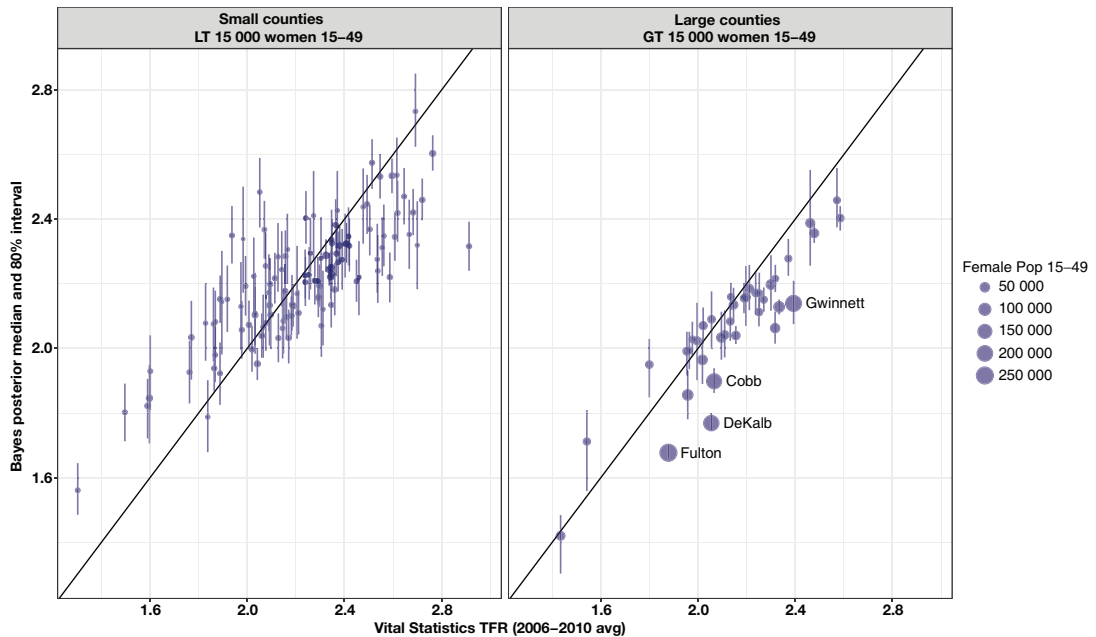


Figure 8 TFR estimates for Georgia counties from vital statistics (horizontal axis) versus 80% posterior intervals estimated from age–sex structure (vertical axis). Posterior medians indicated with dots proportional in size to female population 15–49. Left and right panels contain less and more populous counties, respectively

most populous counties—Fulton (vital statistics = 1.88, Bayes median = 1.68), DeKalb (2.05, 1.77), Cobb (2.07, 1.90) and Gwinnett (2.39, 2.14). Possible explanations for lower TFR estimates from age–sex distributions include (a) errors in census enumeration, with greater proportional undercount of young children than that of adult women; (b) fertility-related migration, with women more likely to move from populous urban counties to suburbs after having children; (c) residency errors in vital statistics, with non-resident mothers who give birth in large-county hospitals mistakenly tabulated as residents. These differences are important and require further investigation, beyond the scope of this article.

7 Discussion and conclusion

Recasting a traditional indirect estimator as a probabilistic model works well. Statistical modelling of demographic relationships can improve point estimates of TFR from age- and sex-specific population counts. It can also improve our analytical understanding of indirect estimators. Even more importantly, a statistical approach produces estimates of uncertainty about demographic parameters derived by indirect methods.

A Bayesian approach like the one we have introduced here is an extensible framework for modelling the data-generating process. It is a ‘plug-and-play’ system for statistical analysis. For example, in our application, one could add models and priors for census errors in the C and W data, for relationships between past and present fertility levels, for spatial patterns in TFR, for relationships between fertility and local socioeconomic covariates, and so forth.

Demography is ripe for marriages between analytical and statistical models for several reasons. First, the formal mathematical foundations of demographic estimators are well established. Second, recent developments in Bayesian modelling and software make it feasible to design and estimate fairly complex models that incorporate realistic demographic relationships. Finally, large, newly available demographic databases such as the HMD and HFD are invaluable sources of prior information. Patterns in these databases (and their variations) allow researchers to design and calibrate useful, informative priors for demographic parameters.

In addition to indirect estimators of total fertility, there are other old demographic dogs that could learn new statistical tricks. The classic UN Manual X (United Nations, 1983) includes indirect methods for estimating child mortality from survey information on the survival of children by mother’s age, for correcting fertility estimates using mothers’ reports about the number of children ever born, and so on. All of these methods could potentially be revived as Bayesian models and used in new applications, following the example in this article.

Appendix

A Tabularsummary of model and priors

The following table corresponds to Figure 1, providing a list of model parameters and relationships.

B Sensitivity to priors

Abbreviating the vector of all relevant parameters other than TFR as θ , Equation (4.8) becomes

$$P(\text{TFR}|C) \propto \int L(C|\text{TFR}, \theta) f(\theta) d\theta$$

For the prior $f(\theta)$ that we propose in this article, fertility proportions ϕ_a come from log ratios $\gamma = m + X\beta$, and the mortality schedule L_a comes from the the Wilmoth et al. (2012) model, as described in Section 4.2.1.

| Category | Type | Distribution or Demographic Relation |
|--------------|------------|--|
| Fertility | Parameter | TFR \sim Uniform(0, 20) |
| | Parameters | $\beta \sim N(0, I_2)$ |
| | Relation | $\gamma = m + X\beta \in \mathbb{R}^7$ |
| | Relation | $\phi_a = \frac{\exp(\gamma_a)}{\sum_{a=15}^{45} \exp(\gamma_a)}, a = 15, 20, \dots, 45$ |
| | Relation | $F_a = \frac{1}{5} \phi_a \text{TFR}, a = 15, 20, \dots, 45$ |
| Mortality | Parameter | $q_5 \sim \text{Beta}$ |
| | | $P(q_5 < 1/2 \min \hat{q}_5) = .05$ |
| | Parameter | $P(q_5 < 2 \max \hat{q}_5) = .95$ |
| | | $k \sim N(0, 1)$ |
| | Relation | $\ln \mu_x = a_x + b_x \ln q_5 + c_x (\ln q_5)^2 + v_x k, x = 0, 1, 5, \dots, 45$ |
| Children 0–4 | Relation | $L_a = L_a(\mu), a = 0, \dots, 45$ |
| | | |
| | Likelihood | $K_a = \text{TFR} \cdot \frac{L_a}{5} \cdot \frac{1}{2} \left(\frac{L_{a-5}}{L_a} \phi_{a-5} + \phi_a \right)$ $C \sim \text{Poisson} \left(\sum_a W_a K_a \right)$ |

For any prior $f(\theta)$, we can approximate the posterior of TFR empirically from a sample of parameter vectors $\theta_1^*, \dots, \theta_S^*$ drawn from $f(\theta)$, as

$$P(\text{TFR} | C) \approx \frac{1}{S} \sum_s L(C | \text{TFR}, \theta_s^*)$$

This emphasizes that the posterior distribution is a formal means of averaging over unknown demographic quantities θ to decide which TFR values are more and less likely.

In the experiments in this appendix, we define a fine grid of possible TFR values v_1, \dots, v_N , sample $\theta_1^*, \dots, \theta_S^*$ from alternative prior distributions $f(\theta)$, calculate $P_{is} = L(C | \text{TFR} = v_i, \theta = \theta_s^*)$ and then approximate the posterior of TFR up to a scalar multiple as

$$P(\text{TFR} = v_i | C) \propto \sum_s P_{is} \quad (\text{B.1})$$

The key question is whether this posterior distribution is sensitive to the choice of priors $f(\theta)$. In order to investigate, we considered two alternative pairs of fertility and mortality models.

Alternative Prior 1: Empirical For the first alternative prior, which we call *Empirical*, we use $\phi_{15}, \dots, \phi_{45}$ and L_0, \dots, L_{45} schedules drawn randomly from the HFD and HMD, respectively. The HFD has 2 054 complete age patterns, and in the *Empirical* prior, we assume that each of these patterns is equally likely. Thus for each simulation $s = 1, \dots, S$ we draw one ϕ schedule at random from the HFD, with each schedule

having a probability $1/2054$. The HMD has 921 complete L_a schedules. In the *Empirical* prior, we draw one of these schedules for each simulation, with higher probabilities of selection for schedules that better match the \hat{q}_5 estimates for the population. (Specifically, the selection probability for schedule i is proportional to the $q_5 \sim \text{Beta}$ density described in Equation (4.4)).

Alternative Prior 2: IDBHIV For the second alternative prior, which we call IDBHIV, we consider much wider possible variations in fertility and mortality patterns, and we derive these patterns from alternative sources. We draw fertility schedules randomly from a Dirichlet distribution $\phi \sim \text{Dirichlet}(\psi)$, where parameters $\psi = (4.7 \ 14.3 \ 17.7 \ 13.5 \ 7.1 \ 2.7 \ 0.7)'$ were fit by maximum likelihood to 226 schedules in the US Census Bureau's IDB (these schedules are available online at http://schmert.net/calibrated-spline/data/IDB_5fx.csv). IDB schedules span a much wider variety of shapes than those in the HFD, because they include many African, Asian, and Latin American countries. For the mortality patterns in the IDBHIV prior, we consider possible effects of HIV prevalence on the age pyramid and the CWR. For each simulation $s = 1, \dots, S$, we draw a q_5 value from the $q_5 \sim \text{Beta}[a(\hat{q}_5), b(\hat{q}_5)]$ distribution, and then use the *mortmod.5q0* function in R package *HIV.LifeTables* (Sharro, 2013) to calculate the corresponding L_0, \dots, L_{45} values in a population with a 10% adult HIV prevalence rate. In general, these mortality patterns imply higher ratios of surviving children to surviving mothers at any given level of TFR, thus raising the posterior probabilities of lower levels of TFR conditional on an observed CWR. The main question for the IDBHIV posterior is whether this effect would be large in a population with a high prevalence of HIV and a higher mortality of adult women.

Comparative Results: Posterior distributions are not sensitive to priors For each of the 159 counties in Georgia, and for the Kanamari indigenous territory, we drew $S = 1000$ fertility and mortality schedules for each of the three sets of priors: the priors proposed in the text, the *Empirical* priors, and the IDBHIV priors. We then calculated the implied posterior distribution of TFR for each population via equation (B.1).

Figure 9 illustrates the results of the sensitivity experiments for the 159 Georgia counties. (Results for the Kanamari indigenous population were extremely similar, but we omit them in order to provide a narrower scale and thus more detail in the figure.) The six panels in Figure 9 correspond to two alternative priors (*Empirical* priors for the top three panels, IDBHIV priors for the bottom three panels), and three quantiles of the posterior distribution (from left to right: 10, 50, and 90 percentiles). Within each panel the horizontal axis represents the posterior quantile when using the prior that we propose in the text, and the vertical axis represents the posterior quantile under one of the alternative priors. Each county in Georgia contributes one point in each panel.

For example, the top-right point in all six panels is for Chattahoochee County. The 10th percentile of the posterior distribution for TFR in this county is 2.72 with our proposed priors, 2.69 with *Empirical* priors, and 2.73 with IDBHIV priors. Thus

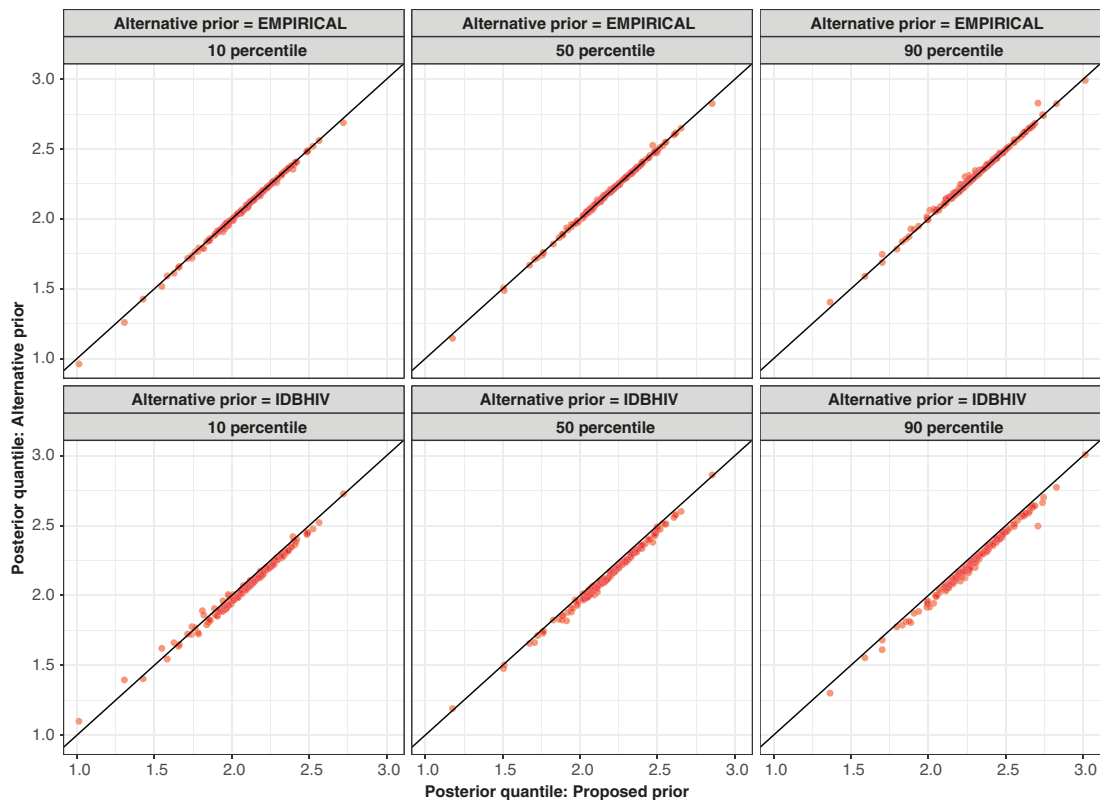


Figure 9 Marginal posterior quantiles of TFR under alternative priors for age patterns of fertility and mortality. Horizontal axis of each panel corresponds to the priors proposed in the main text of this article. Vertical axes correspond to quantiles for the alternative priors described in this appendix

the top-left panel contains the point (2.72, 2.69), and the bottom-left panel contains (2.72, 2.73). All other points in Figure 9 are analogous.

All points in Figure 9 lie close to the 45° line, which means that posterior medians and 80% intervals are very similar under all three priors. Elevated adult mortality, as in populations with HIV prevalence in the bottom panels, tends to shift posterior TFR distributions towards lower values. This effect is very slight, however.

Our strong conclusion from this exercise is that the posterior for TFR is not sensitive to the choice of priors. Bayesian TFR estimates are driven mainly by variation in the age–sex distribution, and are similar under many plausible distributions of fertility and mortality age patterns.

Declaration of conflicting interests

The authors declared no potential conflicts of interest with respect to the research, authorship and/or publication of this article.

Funding

The authors received no financial support for the research, authorship and/or publication of this article.

References

- Arriaga EE, Johnson PD and Jamison E (1994) *Population analysis with microcomputers*, Vol. 1. Suitland, MD: United States Census Bureau.
- Bogue D and Palmore JA (1964) Some empirical and analytic relations among demographic fertility measures, with regression models for fertility estimation. *Demography*, **1**, 316–38.
- Brass W (1964) *Uses of census or survey data for the estimation of vital rates* (United Nations publication E/CN.14/CAS/7). New York, NY: United Nations.
- Brunsdon C, Fotheringham S and Charlton M (1998) Geographically weighted regression. *Journal of the Royal Statistical Society: Series D (The Statistician)*, **47**, 431–43.
- Carpenter B, Gelman A, Hoffman M, Lee D, Goodrich B, Betancourt M, Brubaker M, Guo J, Li P and Riddell A (2017) Stan: A probabilistic programming language. *Journal of Statistical Software*, **76**, 1–32. doi: 10.18637/jss.v076.i01. URL <https://www.jstatsoft.org/index.php/jss/article/view/v076i01> (last accessed 20 September 2018).
- Coale AJ and Trussell TJ (1974) Model fertility schedules: Variations in the age structure of childbearing in human populations. *Population Index*, **40**, 185–258.
- Gunasekaran S and Palmore J (1984) Regression estimates of the gross reproduction rate using moments of the female age distribution. *Asian and Pacific Census Forum*, **10**, 5–10.
- Hanenberg R (1983) Estimates of the total fertility rate based on the child-woman ratio. *Asian and Pacific Census Forum*, **10**, 5–11.
- Hauer M, Baker J and Brown W (2013) Indirect estimates of total fertility rate using child woman/ratio: A comparison with the Bogue–Palmore method. *PLOS ONE*, **8**, e67226. doi: 10.1371/journal.pone.0067226. URL <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0067226> (last accessed 20 September 2018).
- Instituto Brasileiro de Geografia e Estatística (2016) *Resultados do universo agregados por setores censitários* [Short-form aggregate results for census sector]. URL ftp://ftp.ibge.gov.br/Censos/Censo_Demografico_2010/Resultados_do_Universo/Agregados_por_Setores_Censitarios/ (last accessed 20 September 2018).
- International Union for the Scientific Study of Population (2017) *PAPP–Population Analysis for Policies & Programmes*. URL <http://papp.iussp.org/index.html> (last accessed 20 September 2018).
- Max Planck Institute for Demographic Research and Vienna Institute of Demography (2016) *Human fertility database*. URL <http://www.humanfertility.org/cgi-bin/main.php> (last accessed 20 September 2018).
- Nolin DA and Ziker JP (2016) Reproductive responses to economic uncertainty. *Human Nature*, **27**, 351–71.
- Palmore JA (1978) Regression estimates of changes in fertility, 1955–60 to 1965–75, for most major nations and territories. East-West Center, Honolulu, Hawaii.
- R Core Team (2016) *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/> (last accessed 20 September 2018).

- Rele J (1967) *Fertility analysis through extension of stable population concepts*. Population monograph series. Institute of International Studies, University of California. URL <https://books.google.com/books?id=nC-1nQEACAAJ> (last accessed 20 September 2018).
- Schmertmann CP, Cavenaghi SM, Assunção RM and Potter JE (2013) Bayes plus Brass: Estimating total fertility for many small areas from sparse census data. *Population Studies*, 67, 255–73. doi: 10.1080/00324728.2013.795602. URL <http://dx.doi.org/10.1080/00324728.2013.795602> (last accessed 20 September 2018).
- Sharrow D (2013) *HIV.LifeTables: HIV calibrated model life tables for countries with generalized HIV epidemics*. URL <https://CRAN.R-project.org/package=HIV.LifeTables> (last accessed 20 September 2018). R package version 0.1.
- Stan Development Team (2016) *RStan: The R interface to Stan*. URL <http://mc-stan.org/> (last accessed 20 September 2018). R package version 2.14.1.
- Terras Indígenas (2017) Terras indígenas do Brasil [Brazilian Indigenous Territories]. URL <https://terrasindigenas.org.br/es/terras-indigenas/3718> (last accessed 20 September 2018).
- Tuchfeld BS, Guess LL and Hastings DW (1974) The Bogue–Palmore technique for estimating direct fertility measures from indirect indicators as applied to Tennessee counties, 1960 and 1970. *Demography*, 11, 195–205.
- United Nations (1983) *Manual X: Indirect techniques for demographic estimation* (United Nations publication, Sales No. E.83.XIII.2). New York, NY: Author.
- United Nations Development Program (2013) *Atlas do Desenvolvimento Humano no Brasil*. [Atlas of Human Development in Brazil] URL <http://www.atlasbrasil.org.br/2013/> (last accessed on 20 September 2018).
- United States Census Bureau (2016) *U.S. Census Bureau International Database*. URL <https://www.census.gov/population/international/data/idb> (last accessed 20 September 2018).
- University of California, Berkeley (USA) and Max Planck Institute for Demographic Research (Germany) (2016) *Human mortality database*. URL <http://www.mortality.org> (last accessed 20 September 2018).
- Wachter KW (2014) *Essential demographic methods*. Harvard University Press. URL <http://www.jstor.org/stable/j.ctt6wps5v> (last accessed 20 September 2018).
- Wilmoth J, Zureick S, Canudas-Romo V, Inoue M and Sawyer C (2012) A flexible two-dimensional mortality model for use in indirect estimation. *Population Studies*, 66, 1–28. doi: 10.1080/00324728.2011