

Bayesian Forecasting of Cohort Fertility

Understanding of Methods

Yichen He¹

¹Department of Statistics and Data Science
National University of Singapore

All methods and charts come from Schmertmann et.al[1]

Table of Contents

- 1 Notation
- 2 General Idea
- 3 Shape Penalty
- 4 Time Penalty
- 5 Weighting
- 6 Final result
- 7 Proof
- 8 Reference

Table of Contents

- 1 Notation
- 2 General Idea
- 3 Shape Penalty
- 4 Time Penalty
- 5 Weighting
- 6 Final result
- 7 Proof
- 8 Reference

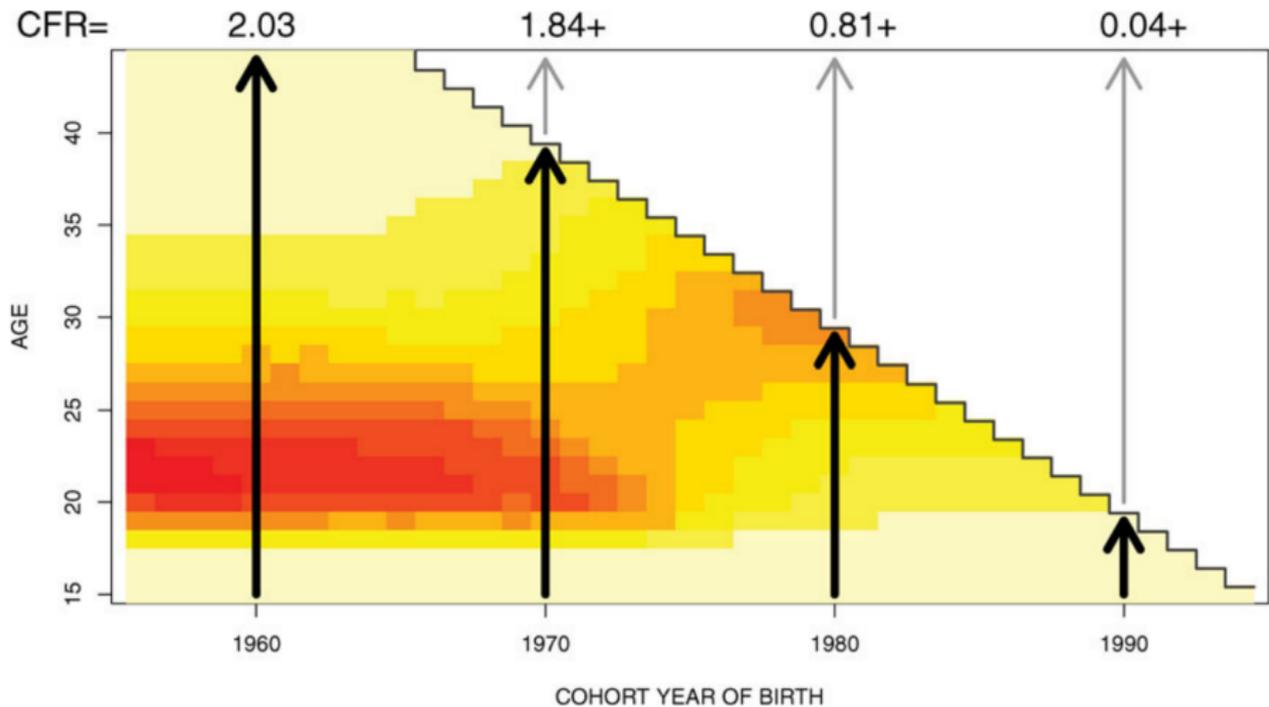
For a specific country, we have C birth cohorts of interest ($c = 1 \dots C$) over A reproductive ages ($a = 1 \dots A$).

- $\theta_{ca} \in \mathbb{R}$, the true fertility rate for cohort c between exact ages a and $a + 1$;
- $\theta_c = (\theta_{c1} \dots \theta_{cA})^\top \in \mathbb{R}^A$, the fertility schedule for cohort c ;
- $\theta_a = (\theta_{1a} \dots \theta_{Ca})^\top \in \mathbb{R}^C$, the time series of rates at age a ;
- $\theta = (\theta_1^\top \dots \theta_C^\top)^\top \in \mathbb{R}^{CA}$, **the vector of all rates, sorted by age within cohort**;

Table of Contents

- 1 Notation
- 2 General Idea**
- 3 Shape Penalty
- 4 Time Penalty
- 5 Weighting
- 6 Final result
- 7 Proof
- 8 Reference

General Idea



Notation

1. $y \in \mathbb{R}^n$, a vector of published data for some subset of θ ;
2. $V \in \mathbb{R}^{n \times CA}$, a matrix of ones and zeroes such that $V\theta \in \mathbb{R}^n$ is the subset of parameters corresponding to y .

Assume our historical data y generated from the normal distribution:

$$y \mid \theta \sim N_n(V\theta, \Psi)$$

where $\Psi = \text{diag}_{i=1 \dots n} [y_i (1 - y_i) / W_i]$ and W_i is the number of a -year-old women in the (c, a) cell corresponding to the i -th rate.

Hence, we can based on the distribution, give out the log-likelihood function:

$$\ln L(y \mid \theta) = \text{const} - \frac{1}{2}(y - V\theta)^\top \Psi^{-1}(y - V\theta)$$

Notation

1. $y \in \mathbb{R}^n$, a vector of published data for some subset of θ ;
2. $V \in \mathbb{R}^{n \times CA}$, a matrix of ones and zeroes such that $V\theta \in \mathbb{R}^n$ is the subset of parameters corresponding to y .

According to the Lexis surface, we know that θ is not complete, so we still need a priori information for our θ :

$$\theta \sim N_{CA}(\underline{0}, K^{-1})$$

Therefore if we combine the prior information with log-likelihood together, we can get the expression for posterior for our θ in a Bayesian framework:

$$\ln P(\theta | y) = \text{const} + \ln L(y | \theta) + \ln f(\theta) \quad (1)$$

$$= \text{const} - \frac{1}{2}(y - V\theta)^\top \Psi^{-1}(y - V\theta) - \frac{1}{2}\theta^\top K\theta \quad (2)$$

Recall

$$\text{Posterior: } \ln P(\theta | y) = \text{const} - \frac{1}{2}(y - V\theta)^\top \Psi^{-1}(y - V\theta) - \frac{1}{2}\theta^\top K\theta$$

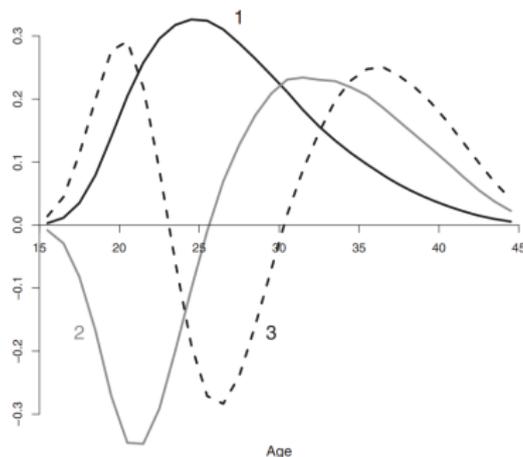
As we can tell from the above formula, the most important thing is the covariance matrix K . This part is aggregated by two different penalty terms, **Shape penalty** and **Time penalty**. We will explain these two specially-designed terms in the next two sections in detail.

Table of Contents

- 1 Notation
- 2 General Idea
- 3 Shape Penalty**
- 4 Time Penalty
- 5 Weighting
- 6 Final result
- 7 Proof
- 8 Reference

In the real case, we can set the number of ages and cohorts of interest as $A = 30$ ages (15, \dots , 44) and $C = 40$ cohorts (1956, \dots , 1995). In this way, for a specific country, we can define the historical array Φ .

To extract its rough pattern of previous data, we apply **Singular Value Decomposition (SVD)** on Φ , and choose the first three principle components, denote them as $X \in \mathbb{R}^{30 \times 3}$.



Notation

1. $\theta_c = (\theta_{c1} \dots \theta_{cA})^\top \in \mathbb{R}^A$, the fertility schedule for cohort c

To see to what extent our real rates θ_c could be explained by the first three principle components, we can do the projection on it:

$$\theta_c = X \left(X^\top X \right)^{-1} X^\top \theta_c + \varepsilon_c$$

After projection, we only need ε_c to measure to what degree the cohort matches with our historical pattern. Note that residual vector ε_c is:

$$\varepsilon_c = \left[I_A - X \left(X^\top X \right)^{-1} X^\top \right] \theta_c = M \theta_c$$

Construct Shape penalty

Recall

$$1. \varepsilon_c = \left[I_A - X (X^T X)^{-1} X^T \right] \theta_c = M \theta_c$$

To further quantify 'small' of residual vector, we consider calculating their average outer product :

$$\bar{\Omega} = \frac{1}{s} \sum_s \varepsilon_s \varepsilon_s^T$$

Then, these historical data allow us to establish a scalar penalty for the "badness" of each cohort schedule's shape:

$$\begin{aligned} \pi_c &= \varepsilon_c^T \bar{\Omega}^\dagger \varepsilon_c \\ &= \theta_c^T \left[M \bar{\Omega}^\dagger M \right] \theta_c \\ &= \theta^T \left[G_c^T M \bar{\Omega}^\dagger M G_c \right] \theta \\ &= \theta^T K_c \theta, \end{aligned}$$

Results

- An important feature of this π_c is that it's **improper**;
- By construction, the empirical average of π_c across the historical cohort schedules in Φ equals 27;
- The shape penalty term will give penalty on those ones which don't share similar trend of historic data, but will not be too heavy;

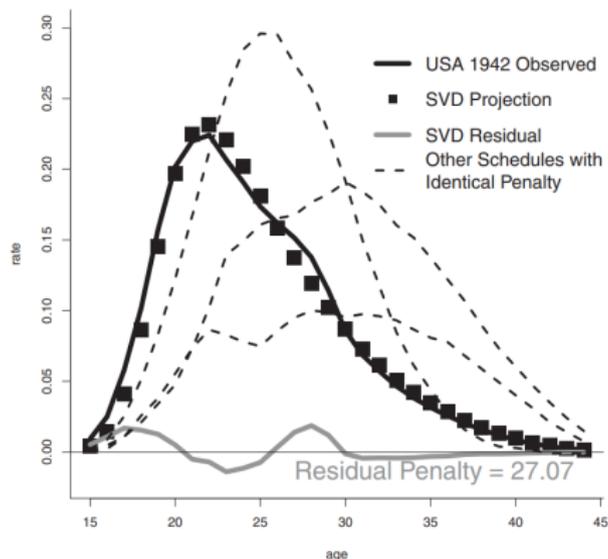


Table of Contents

- 1 Notation
- 2 General Idea
- 3 Shape Penalty
- 4 Time Penalty**
- 5 Weighting
- 6 Final result
- 7 Proof
- 8 Reference

Freeze rate

The freeze-rate method assumes that the most likely future value for the fertility rate at age a is simply the last observed rate at that age. It suggests that

$$\theta_{a,c+1} \approx \theta_{a,c}$$

which can be expressed in the numerical way, note that at each age on the Lexis surface, we define a vector of 30 freeze-rate residuals for cohorts 1966–1995:

$$\begin{aligned} u_a &= \begin{bmatrix} \theta_{a,1966} - \theta_{a,1965} \\ \vdots \\ \theta_{a,1995} - \theta_{a,1994} \end{bmatrix} \\ &= \begin{bmatrix} 0 & \cdots & -1 & 1 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & -1 & 1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & -1 & 1 \end{bmatrix} \theta_a \\ &= W_R \theta_a = W_R H_a \theta \end{aligned}$$

Freeze slope

Similarly, the freeze-slope method assumes that trends, measured as fitted slopes over some recent period, will continue into the future. It suggests that

$$\theta_{a,c+1} \approx \theta_{a,c} + \hat{\Delta}_c$$

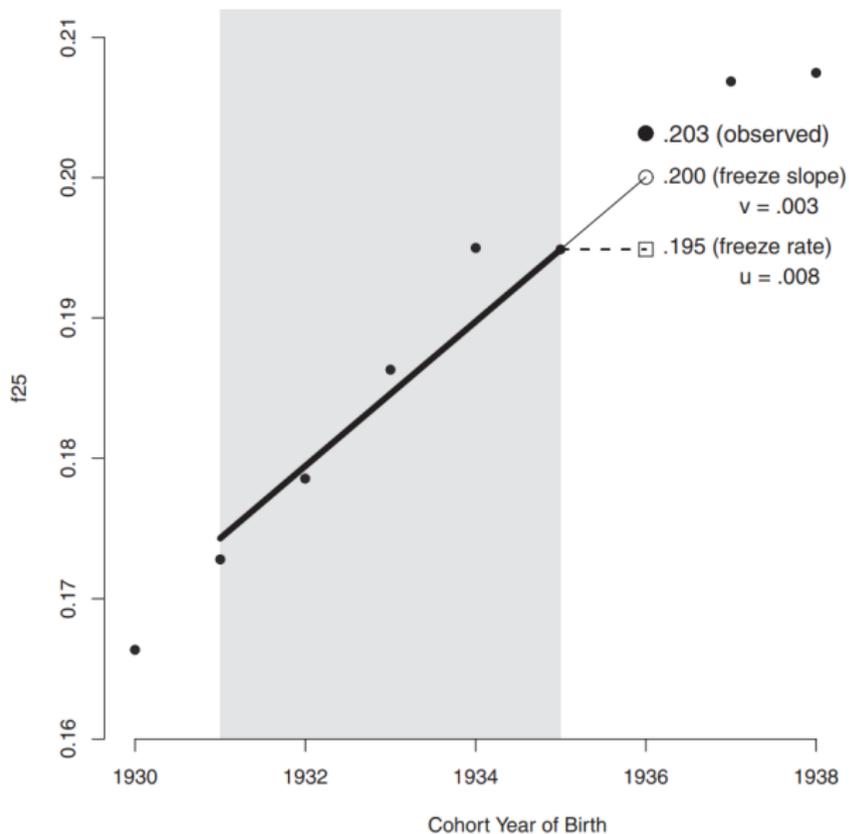
which can be expressed in the numerical way:

$$v_a = \begin{bmatrix} \theta_{a,1966} - (\theta_{a,1965} + \hat{\Delta}_{1965}) \\ \vdots \\ \theta_{a,1995} - (\theta_{a,1994} + \hat{\Delta}_{1994}) \end{bmatrix} = W_S \theta_a = W_S H_a \theta$$

where under this context, $\hat{\Delta}$ could be expressed by:

$$\hat{\Delta}_c = \frac{1}{30} (10\theta_{a,c} - \theta_{a,c-1} - 2\theta_{a,c-2} - 3\theta_{a,c-3} - 4\theta_{a,c-4})$$

Visualization



Construct Time penalty

Recall

1. Freeze rate: $u_a = W_R \theta_a = W_R H_a \theta$
2. Freeze slope: $v_a = W_S \theta_a = W_S H_a \theta$

As we did in Shape Penalty, to extrapolate the previous time trend into future in the horizontal direction, we also construct the time penalty by applying standardization and aggregate them into quadratic penalty term:

Freeze rate

$$\begin{aligned}\pi_{Ra} &= s_{Ra}^{-2} u_a^\top u_a \\ &= \theta^\top \left[s_{Ra}^{-2} H_a^\top W_R^\top W_R H_a \right] \theta \\ &= \theta^\top K_{Ra} \theta\end{aligned}$$

Freeze slope

$$\begin{aligned}\pi_{Sa} &= s_{Sa}^{-2} v_a^\top v_a \\ &= \theta^\top \left[s_{Sa}^{-2} H_a^\top W_S^\top W_S H_a \right] \theta \\ &= \theta^\top K_{Sa} \theta\end{aligned}$$

where s_{Ra}^{-2} , s_{Sa}^{-2} are the **average (mean) squared residuals** of freeze rate and freeze slope for each (age, method) combination respectively.

Table of Contents

- 1 Notation
- 2 General Idea
- 3 Shape Penalty
- 4 Time Penalty
- 5 Weighting**
- 6 Final result
- 7 Proof
- 8 Reference

Note that there are actually 90 penalty terms in total, including 30 **Shape Penalty** terms, 30 **Freeze Rate Penalty** terms and 30 **Freeze Slope Penalty** terms.

Since the residuals on which we base the penalties are not mutually independent, we need to assign non-unit weights to each term to better modify our prior distribution. We first assume the covariance matrix K to be constructed as:

$$K = \sum_{j=1}^{90} w_j K_j$$

where $K_j \in \mathbb{R}^{CA \times CA}$, are penalty matrices.

Theorem

Using notations, condition on θ is restricted to the column space of K , and given weights, we can calculate expectation of π_j based on the trace of multiplication of corresponding matrices and K , which is

$$E^* (\pi_j | w) = \text{trace} \left(K_j K^\dagger \right)$$

Table 1. Summary of penalties for rate surfaces over birth cohorts 1956–1995 ($C = 40$) and ages 15–44 ($A = 30$)

	Schedule shapes	Time-series (freeze rate)	Time-series (freeze slope)
Number of penalties	30	30	30
Penalty terms	$\pi_{1966} \dots \pi_{1995}$	$\pi_{R,15} \dots \pi_{R,44}$	$\pi_{S,15} \dots \pi_{S,44}$
Residuals	$\varepsilon_c = \mathbf{M} \theta_c$	$\mathbf{u}_a = \mathbf{W}_R \theta_a$	$\mathbf{v}_a = \mathbf{W}_S \theta_a$
Penalty matrices	$\mathbf{K}_{1966} \dots \mathbf{K}_{1995}$	$\mathbf{K}_{R,15} \dots \mathbf{K}_{R,44}$	$\mathbf{K}_{S,15} \dots \mathbf{K}_{S,44}$
A priori assumption	Incomplete schedules well approximated by SVD basis functions \mathbf{X}	Next cohort's rate at age a well predicted by current rate	Next cohort's rate at age a well predicted by recent trend
Calibration information from historical data	Projection errors from \mathbf{X}	One-ahead freeze-rate prediction errors	One-ahead freeze-slope prediction errors
Number of elements in each residual	30	30	30
Expected value of each penalty (= rank of \mathbf{M} or \mathbf{W})	27	30	30

Recall

1. $\text{target}_j = 27, j = 1, \dots, 30$
2. $\text{target}_j = 30, j = 31, \dots, 90$

- Initialize all weights at unity: $w_1 = w_2 = \dots = w_{90} = 1$;
- Calculate $K = \sum w_j K_j$, and its generalized inverse K^\dagger ;
- Calculate $E^*(\pi_j | w) = \text{trace}(K_j K^\dagger)$, for all $j = 1 \dots 90$;
- Update weights as $w_j^{\text{new}} = w_j \cdot \frac{E^*(\pi_j | w)}{\text{target}_j}, j = 1 \dots 90$;
- Stop if converged; otherwise return to Step 1.

Final posterior expression

Finally, after calculating the K matrix for prior distribution, we can write the expression for posterior distribution based on the formula:

$$\ln P(\theta | y) = \text{const} - \frac{1}{2}(y - V\theta)^\top \Psi^{-1}(y - V\theta) - \frac{1}{2}\theta^\top K\theta$$

And through the procedure of maximizing the posterior distribution, we can get the optimal μ_{post} and Σ_{post} , it suggests that:

$$\theta | y \sim N \left\{ \begin{array}{l} \mu_{\text{post}} = [V^\top \Psi^{-1} V + K]^{-1} [V^\top \Psi^{-1} y] \\ \Sigma_{\text{post}} = [V^\top \Psi^{-1} V + K]^{-1} \end{array} \right\}$$

which is our final result.

Table of Contents

- 1 Notation
- 2 General Idea
- 3 Shape Penalty
- 4 Time Penalty
- 5 Weighting
- 6 Final result**
- 7 Proof
- 8 Reference

An example for Singapore cohort fertility

Singapore 2010 Forecast

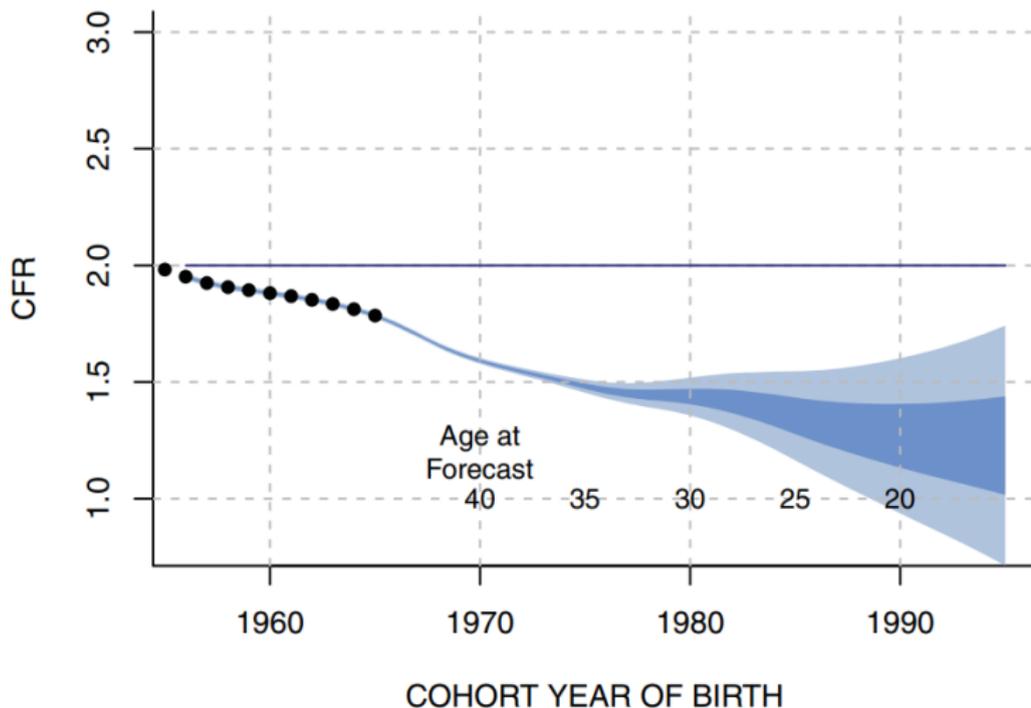


Table of Contents

- 1 Notation
- 2 General Idea
- 3 Shape Penalty
- 4 Time Penalty
- 5 Weighting
- 6 Final result
- 7 Proof**
- 8 Reference

Penalty Terms

Why expected value of shape penalty, freeze-rate penalty and freeze-slope penalty equal to 27, 30 and 30 respectively? (As in table 1)

Pf: 1) **Shape Penalty**

According to the paper, X is the first three principal components of historical array $\Phi \in \mathbb{R}^{A \times S}$, where S is the number of historical complete cohorts in 1900 to 1949, i.e.

$$\text{Denote } \Phi = UDV^T, \quad X = U(:, 1:3) \in \mathbb{R}^{A \times 3}$$

Then we consider the residual vector between θ_c and θ_c 's linear projection on X 's column space ($c = 1, \dots, S$), ie.

$$\varepsilon_c = \left[I_A - X \left(X^T X \right)^{-1} X^T \right] \theta_c \triangleq M \theta_c$$

Proof

Note that $M (= I_A - X (X^T X)^{-1} X^T)$ is an idempotent matrix, therefore,

$$\begin{aligned}\text{rank}(M) &= \text{trace}(M) \\ &= A - \text{trace} \left(X (X^T X)^{-1} X^T \right) \\ &= A - \text{trace} \left((X^T X)^{-1} X^T X \right) \\ &= A - 3\end{aligned}$$

Still, according to the paper, $\bar{\Omega} = \frac{1}{S} \sum_{c=1}^S \varepsilon_c \varepsilon_c^T$ is covariance matrix for complete cohorts between 1900-1949, then,

$$\begin{aligned}\frac{1}{S} \sum_{c=1}^S \varepsilon_c^T \bar{\Omega}^\dagger \varepsilon_c &= \frac{1}{S} \text{trace} \left(\bar{\Omega}^\dagger \sum_{c=1}^S \varepsilon_c \varepsilon_c^T \right) \\ &= \text{trace} \left(\bar{\Omega}^\dagger \bar{\Omega} \right) = \text{rank}(M) \\ &= A - 3\end{aligned}$$

2) Freeze-rate Penalty

$$\begin{aligned}\frac{1}{A} \sum_{a=1}^A s_{R_a}^{-2} u_a^\top u_a &= \frac{1}{A} \sum_{a=1}^A \left(\frac{\sum_{c=1}^S u_{a,c}^2}{S_{R_a}^2} \right) \\ &= \frac{1}{A} \sum_{a=1}^A S \\ &= S\end{aligned}$$

3) Freeze-slope Penalty

Similarly,

$$\frac{1}{A} \sum_{a=1}^A s_{s_a}^{-2} v_a^\top v_a = S$$

To conclude, by designing prior above, during forecasting procedure, we believe they generate from

for shape penalty,

$$\varepsilon_c \sim N(\underline{0}, \bar{\Omega})$$

for freeze rate penalty,

$$u_a \sim N\left(\underline{0}, \begin{pmatrix} s_{R_1}^2 & & \\ & \dots & \\ & & s_{R_a}^2 \end{pmatrix}\right)$$

for freeze rate penalty,

$$v_a \sim N\left(\underline{0}, \begin{pmatrix} s_{S_1}^2 & & \\ & \dots & \\ & & s_{S_A}^2 \end{pmatrix}\right)$$

Hence, when $a = 1, \dots, 30(15, \dots, 44)$, $C = 1, \dots, 30(1966, \dots, 1995)$, we have

$$\mathbb{E}[\pi_c \mid \text{history}] = A - 3 = 27$$

$$\mathbb{E}[\pi_{R_a} \mid \text{history}] = C = 30$$

$$\mathbb{E}[\pi_{S_a} \mid \text{history}] = C = 30$$

Table of Contents

- 1 Notation
- 2 General Idea
- 3 Shape Penalty
- 4 Time Penalty
- 5 Weighting
- 6 Final result
- 7 Proof
- 8 Reference**



Schmertmann, C., Zagheni, E., Goldstein, J. R., Myrskylä, M. (2014). Bayesian forecasting of cohort fertility. *Journal of the American Statistical Association*, 109(506), 500-513.