

Summary: The Probabilistic Model of Age Assignment

February 28th 2017

Taehee Lee

0. Notations:

- $d^i = \{d_1^i, d_2^i, \dots, d_{L_i}^i\}$: the depth index of the data points of record i .
- $v^i = \{v_1^i, v_2^i, \dots, v_{L_i}^i\}$: $\delta^{18}\text{O}$ data of record i .
- $w^i = \{w_1^i, w_2^i, \dots, w_{L_i}^i\}$: radiocarbon data of record i .
- $A^i = \{A_1^i, A_2^i, \dots, A_{L_i}^i\}$: the ages in the stack that the data points of record i are aligned to.
- $R^i = \{R_1^i, R_2^i, \dots, R_{L_i}^i\}$: the inverse sedimentation ratio defined by following:

$$R_n^i \triangleq \frac{A_n^i - A_{n-1}^i}{d_n^i - d_{n-1}^i}$$

- τ^i : the mean shift on the $\delta^{18}\text{O}$ values of record i .
- $\mu = \{\mu_t\}_{1 \leq t \leq T}$: $\delta^{18}\text{O}$ values for the ages $1 \leq t \leq T$ in the target stack.
- $\sigma = \{\sigma_t\}_{1 \leq t \leq T}$: standard deviations of residuals of observed $\delta^{18}\text{O}$ values for the ages $1 \leq t \leq T$ in the target stack.
- $\theta = \{\theta_t\}_{1 \leq t \leq T}$: radiocarbon parameter for the ages $1 \leq t \leq T$ in the target stack.
- $N(a|b, c)$: the pdf of normal distribution at a with mean b and standard deviation c .
- $f(d|e)$: the pdf of radiocarbon data distribution at d with parameter e .
- $\rho: \{1, 2, \dots, R\} \times \{1, 2, \dots, R\} \rightarrow [0, 1]$: the transition function.

1. What we have observed, assumed, or known: $\{d^i\}_{i=1}^M, \{v^i\}_{i=1}^M, \{w^i\}_{i=1}^M, \theta, N, f, \rho$.

2. What we need to estimate: $\{A^i\}_{i=1}^M, \{R^i\}_{i=1}^M, \{\tau^i\}_{i=1}^M, \mu, \sigma$.

3. Comment:

Because it is reasonable to assume the finite number of sedimentation rates (say, 17) with pseudo-data for the missing time obtained by interpolating the adjacent ones, to make the algorithm faster, in this method we do not assume the continuous domain of ρ . Also, by using radiocarbon data, we can give not only a better estimation but also shorter running time of the algorithm in practice.

Let us assume that each $\delta^{18}\text{O}$ and radiocarbon data are conditionally independent given the time when they were generated. If not, we need to either learn their joint distribution or give a new model on the relationships of those data with some related parameter sets.

The pdf of radiocarbon data f is defined by taking 0 outside of $[\text{mean} - 5\text{std}, \text{mean} + 5\text{std}]$.

4. Model: it consists of two steps. One is the forward algorithm and the other is the backward sampling.

4.1. Forward Algorithm:

$$\begin{aligned} F^i(n, r_n^i, t_n^i) &\triangleq p(v_{1:n}^i, w_{1:n}^i, R_n^i = r_n^i, A_n^i = t_n^i | d_{1:L_i}^i) \\ &= N(v_n^i | \tau^i + \mu_{t_n^i}, \sigma_{t_n^i}) f(w_n^i | \theta_{t_n^i}) \sum_r \rho(r, r_n^i) F^i(n-1, r, t_n^i - r_n^i(d_n^i - d_{n-1}^i)) \end{aligned}$$

4.2. Backward Sampling:

For each i , get $\{\tilde{t}_n^{i,l}\}_{l=1}^{1000} = \{\{\tilde{t}_n^{i,l}\}_{n=1}^{L_i}\}_{l=1}^{1000}$ and $\{\tilde{r}_n^{i,l}\}_{l=1}^{1000} = \{\{\tilde{r}_n^{i,l}\}_{n=1}^{L_i}\}_{l=1}^{1000}$ by following:

$$\begin{aligned} \tilde{t}_n^{i,l} &\equiv \tilde{t}_{n+1}^{i,l} - \tilde{r}_{n+1}^{i,l}(d_{n+1}^i - d_n^i) \\ \mathbb{P}(R_n^{i,l} = r | A_{n+1}^{i,l} = \tilde{t}_{n+1}^{i,l}, R_{n+1}^{i,l} = \tilde{r}_{n+1}^{i,l}, v_{1:n+1}^i, w_{1:n+1}^i, d_{1:L_i}^i) \\ &\propto N(v_{n+1}^i | \tau^i + \mu_{\tilde{t}_{n+1}^{i,l}}, \sigma_{\tilde{t}_{n+1}^{i,l}}) f(w_{n+1}^i | \theta_{\tilde{t}_{n+1}^{i,l}}) \rho(r, \tilde{r}_{n+1}^{i,l}) F^i(n, r, \tilde{t}_{n+1}^{i,l} - \tilde{r}_{n+1}^{i,l}(d_{n+1}^i - d_n^i)) \end{aligned}$$

I.e., the very essence of this sampling is to sample $(\tilde{t}_{L_i}^{i,l})$ and $\tilde{r}_{L_i}^{i,l}$ iteratively: after that, $\tilde{t}_{1:L_i-1}^{i,l}$ is automatically assigned just after the corresponding $\tilde{r}_{2:L_i}^{i,l}$ is sampled.

You may have already noticed, however, that in fact there is a subtle but not negligible problem in this finite sedimentation rates setting: we assumed that ages $1 \leq t \leq T$ in the stack is discretized so the terms $t_n^i - r_n^i(d_n^i - d_{n-1}^i)$ in the forward algorithm or $\tilde{t}_n^{i,l} \equiv \tilde{t}_{n+1}^{i,l} - \tilde{r}_{n+1}^{i,l}(d_{n+1}^i - d_n^i)$ in the backward sampling may not be in the time set! If so, $F^i(n-1, r, t_n^i - r_n^i(d_n^i - d_{n-1}^i))$ or $(\mu_{\tilde{t}_n^{i,l}}, \sigma_{\tilde{t}_n^{i,l}}, \theta_{\tilde{t}_n^{i,l}})$ may even not be well-defined. To deal with it, if it is possible to assume that the forward terms $F^i(n, r, t)$ is ‘nice’ for small perturbations on t for each fixed r and the values $(\mu_t, \sigma_t, \theta_t)$ do not change much as t varies a little, then we can either choose the nearest t in the stack or define the interpolated terms from the nearest parameters as those values.

5. Other Possible Approaches:

5.1. Forward-Backward Algorithm:

$$\begin{aligned}
F^i(n, r_n^i, t_n^i) &\triangleq p(v_{1:n}^i, w_{1:n}^i, R_n^i = r_n^i, A_n^i = t_n^i | d_{1:L_i}^i) \\
&= N(v_n^i | \tau^i + \mu_{t_n^i}, \sigma_{t_n^i}) f(w_n^i | \theta_{t_n^i}) \sum_r \rho(r, r_n^i) F^i(n-1, r, t_n^i - r_n^i(d_n^i - d_{n-1}^i)) \\
B^i(n, r_n^i, t_n^i) &\triangleq p(v_{n+1:L_i}^i, w_{n+1:L_i}^i | R_n^i = r_n^i, A_n^i = t_n^i, d_{1:L_i}^i) \\
&= \sum_r \left[N(v_{n+1}^i | \tau^i + \mu_{r(d_{n+1}^i - d_n^i) + t_n^i}, \sigma_{r(d_{n+1}^i - d_n^i) + t_n^i}) \right. \\
&\quad \times f(w_n^i | \theta_{r(d_{n+1}^i - d_n^i) + t_n^i}) \rho(r_n^i, r) B^i(n+1, r, r(d_{n+1}^i - d_n^i) + t_n^i) \Big] \\
P^i(n, r_n^i, t_n^i) &\triangleq p(R_n^i = r_n^i, A_n^i = t_n^i | v_{1:L_i}^i, d_{1:L_i}^i) \propto F^i(n, t_{n-1}^i, t_n^i) B^i(n, t_{n-1}^i, t_n^i)
\end{aligned}$$

Then we can get the exact posterior probability on each element of data by following:

$$\mathbb{P}(A_n^i = t | v_{1:L_i}^i, d_{1:L_i}^i) = \sum_r P^i(n, r, t)$$

5.2. Blocked Gibbs Sampler: this is one of the other possible ways to get samples. One good aspect of this algorithm is that it does not require any burdensome precomputations, such as the forward algorithm in the above methods. However, it does depend on the initial conditions, so if we do not have any heuristic ways of choosing appropriate initial conditions, good sampling cannot be guaranteed.

$$\begin{aligned}
&\mathbb{P}(A_n^{i,l} = t | A_{n-2}^{i,l} = \tilde{t}_{n-2}^{i,l}, A_{n-1}^{i,l} = \tilde{t}_{n-1}^{i,l}, A_{n+1}^{i,l} = \tilde{t}_{n+1}^{i,l}, A_{n+2}^{i,l} = \tilde{t}_{n+2}^{i,l}, v_n^i, d_{1:L_i}^i) \\
&\propto N(v_n^i | \tau^i + \mu_t, \sigma_t) \rho\left(\frac{\tilde{t}_{n-1}^{i,l} - \tilde{t}_{n-2}^{i,l}}{d_{n-1}^i - d_{n-2}^i}, \frac{t - \tilde{t}_{n-1}^{i,l}}{d_n^i - d_{n-1}^i}\right) \rho\left(\frac{t - \tilde{t}_{n-1}^{i,l}}{d_n^i - d_{n-1}^i}, \frac{\tilde{t}_{n+1}^{i,l} - t}{d_{n+1}^i - d_n^i}\right) \rho\left(\frac{\tilde{t}_{n+1}^{i,l} - t}{d_{n+1}^i - d_n^i}, \frac{\tilde{t}_{n+2}^{i,l} - \tilde{t}_{n+1}^{i,l}}{d_{n+2}^i - d_{n+1}^i}\right)
\end{aligned}$$

Note that the above formula involves two blocks for the efficient sampling: $\{A_1^{i,l}, A_3^{i,l}, A_5^{i,l}, \dots\}$ and $\{A_2^{i,l}, A_4^{i,l}, A_6^{i,l}, \dots\}$.

6. How to estimate the parameters: by using the EM algorithm.

First, we assume that we already knew parameters, and infer the hidden states $\{A^i\}_{i=1}^M$. Second, we assume that we already knew the hidden states, and infer the parameters with ML estimators. Now go back to the first step if some predetermined convergence criteria have not been achieved yet.

7. Time complexity: $\mathcal{O}(MT^2)$.

If we assume that the domain of ρ is discretized by R elements, then it is $\mathcal{O}(MTR)$. According to Seonmin's paper, her algorithm took about one month with 200 computing nodes as it takes over 100 hours to complete one iteration: note that the codes in her github only used three sets of data as an example, but in the actual work we will need to deal with a lot of sets of those.