

Particle filters in the dynamichazard package

BENJAMIN CHRISTOFFERSEN

JULY 12, 2018

This vignette covers the particle filter for the `dynamichazard` package in R. See <https://cran.r-project.org/web/packages/dynamichazard/vignettes/ddhazard.pdf> for more information on the package.

I assume that you are familiar with the setup in the `dynamichazard` package and the other estimation methods in the package. Further, some prior knowledge of particle filters is required. If you lack knowledge of particle filter then Doucet and Johansen (2009) provides a tutorial on particle filters and Kantas et al. (2015) covers parameter estimation with particle filters. This vignette relies heavily on Fearnhead et al. (2010).

1 Method

Model

The model is:

$$\begin{aligned} y_{it} &\sim P(y_{it} | \eta_{it}) \\ \eta_t &= \mathbf{X}_t \mathbf{R}^+ \boldsymbol{\alpha}_t + \boldsymbol{o}_t \\ \boldsymbol{\alpha}_t &= \mathbf{F} \boldsymbol{\alpha}_{t-1} + \mathbf{R} \boldsymbol{\epsilon}_t \quad \boldsymbol{\epsilon}_t \sim N(\mathbf{0}, \mathbf{Q}) \quad , \quad \begin{aligned} i &= 1, \dots, n_t \\ t &= 1, \dots, d \end{aligned} \\ \boldsymbol{\alpha}_0 &\sim N(\mathbf{a}_0, \mathbf{Q}_0) \end{aligned} \quad (1)$$

where I denote the conditional densities as $\mathbf{y}_t \sim g(\cdot | \mathbf{R}^+ \boldsymbol{\alpha}_t)$ and $\boldsymbol{\alpha}_t \sim f(\cdot | \boldsymbol{\alpha}_{t-1})$. We are in a survival analysis setting where the simplest model has an indicator of death of individual i in time t such that $y_{it} \in \{0, 1\}$ where η_{it} is the linear predictor where we use the logistic function as the link function. For each $t = 1, \dots, d$, we have a risk set given by $R_t \subseteq \{1, 2, \dots, n\}$. Further, we let $n_t = |R_t|$ and $n_{\max} = \max_{t=1, \dots, d} n_t$. The observed outcomes are denoted by $\mathbf{y}_t = \{y_{it}\}_{i \in R_t}$. \mathbf{X}_t is the design matrix of the covariates and $\boldsymbol{\alpha}_t$ is the vector of time-varying coefficients. The problems we are looking at have $n_{\max} \gg p \geq r$ (e.g. $n_{\max} = 100000$ and $r = 5$).

Superscript $+$ denotes the Moore-Penrose inverse, the i 'th row of \mathbf{X}_t is \mathbf{x}_{it} , $\mathbf{x}_{it}, \boldsymbol{\epsilon}_t \in \mathbb{R}^r$, $\boldsymbol{\alpha}_t \in \mathbb{R}^p$, $\mathbf{F} \in \mathbb{R}^{p \times p}$ and is invertible, $\mathbf{Q} \in \mathbb{R}^{r \times r}$ is a positive definite matrix, \boldsymbol{o}_t are known offsets, and $\mathbf{R} \in \{0, 1\}^{p \times r}$ with $p \geq r$ contains a subset of the columns of \mathbf{I}_p identity matrix with no duplicate columns in \mathbf{R} . The latter implies that $\mathbf{R}^+ = \mathbf{R}^\top$ and \mathbf{R} is left inverse (i.e., $\mathbf{R}^\top \mathbf{R} = \mathbf{I}_r$). $\mathbf{R} \mathbf{R}^\top$ is a $p \times p$ diagonal matrix with r diagonal entries with value 1 and $p - r$ with

value zero. We will let

$$\boldsymbol{\xi}_t = \mathbf{R}^+ \boldsymbol{\alpha}_t$$

I will use a particle filter and smoother to get smoothed estimates of $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_d$ given the outcomes $\mathbf{y}_{1:d} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_d\}$ and use an EM-algorithm to estimate \mathbf{Q} and \mathbf{a}_0 . One choice of smoother is the generalized two-filter smoothing in Fearnhead et al. (2010) and Briers et al. (2010). The rest of vignette is structured as follows: first I cover the particle filter and smoother. Then I cover the EM-algorithm and other miscellaneous topics. I will end with what is implemented at this point.

Considerations

Algorithm 1 shows one of the generalized two-filter smoother from Fearnhead et al. (2010). It requires that we specify the following importance densities and re-sampling weights (optimal values are given as the right hand side):

$$\begin{aligned} q\left(\boldsymbol{\alpha}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t\right) &= \mathbf{P}\left(\boldsymbol{\alpha}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t\right) \\ \beta_t^{(j)} &\propto \mathbf{P}\left(\mathbf{y}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}\right) w_{t-1}^{(j)} \\ \tilde{q}\left(\boldsymbol{\alpha}_t \middle| \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right) \tilde{\beta}_t^{(k)} &\approx \gamma_t(\boldsymbol{\alpha}_t) \mathbf{P}\left(\mathbf{y}_t \middle| \boldsymbol{\alpha}_t\right) \mathbf{P}\left(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \middle| \boldsymbol{\alpha}_t\right) \frac{\tilde{w}_{t+1}^{(k)}}{\gamma_{t+1}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})} \\ \tilde{q}\left(\hat{\boldsymbol{\alpha}}_t^{(i)} \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right) &= \mathbf{P}\left(\boldsymbol{\alpha}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right) \end{aligned} \quad (2)$$

Further, we need to define a backwards filter distribution approximation

$$\tilde{p}(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d}) \propto \gamma_t(\boldsymbol{\alpha}_t) \mathbf{P}(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_t) \quad (3)$$

with an artificial prior distribution $\gamma_t(\boldsymbol{\alpha}_t)$.

Given the models of interest we have that:

- Evaluating $\mathbf{P}(\mathbf{y}_t | \boldsymbol{\alpha}_t)$ is an expensive operation as $n_{\max} \gg r$ and it has a $\mathcal{O}(n_{\max} r)$ computational cost. Any $\mathcal{O}(n_{\max})$ operation is going to take considerable time.
- Evaluating $f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1})$ and $f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t+1})$ is cheap and can be done in closed form and sampling from these distribution can be done in closed form.
- The second example in Fearnhead et al. (2010) is close to the model here though with $n_{\max} = 1$ and $p = 2$.

The following sections will closely follow the example shown in Fearnhead et al. (2010) and the appendix of the paper.

Forward filter (algorithm 2)

This section will cover some options for algorithm 2. Let $\mathcal{N}(\cdot|\cdot, \cdot)$ denote a multivariate normal distribution. We can select the proposal density as

$$q\left(\boldsymbol{\alpha}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t\right) = \mathcal{N}\left(\boldsymbol{\xi}_t \middle| \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{Q}\right) \quad (4)$$

which we can sample from in $\mathcal{O}(Np^2)$ time if we have a pre-computed Cholesky decomposition of \mathbf{Q} . This is often called the *bootstrap filter*. Another option is to use normal approximation of \mathbf{y}_t 's conditional density

$$\begin{aligned} g(\mathbf{y}_t | \boldsymbol{\xi}_t) &\simeq \tilde{g}_t(\mathbf{y}_t | \boldsymbol{\xi}_t) \\ &= \mathcal{N}\left(\mathbf{X}_t \boldsymbol{\xi}_t \middle| \mathbf{X}_t \mathbf{R}^+ \mathbf{F} \bar{\boldsymbol{\alpha}}_{t-1} - \mathbf{o}_t + \mathbf{G}_t (\mathbf{R}^+ \mathbf{F} \bar{\boldsymbol{\alpha}}_{t-1} + \mathbf{o}_t)^{-1} \mathbf{g}_t (\mathbf{R}^+ \mathbf{F} \bar{\boldsymbol{\alpha}}_{t-1} + \mathbf{o}_t), -\mathbf{G}_t (\mathbf{R}^+ \mathbf{F} \bar{\boldsymbol{\alpha}}_{t-1} + \mathbf{o}_t)^{-1}\right) \\ &= \mathcal{N}\left(\mathbf{X}_t \boldsymbol{\xi}_t \middle| \mathbf{X}_t \mathbf{R}^+ \bar{\boldsymbol{\alpha}}_t - \mathbf{o}_t + \mathbf{G}_t (\mathbf{R}^+ \bar{\boldsymbol{\alpha}}_t + \mathbf{o}_t)^{-1} \mathbf{g}_t (\mathbf{R}^+ \bar{\boldsymbol{\alpha}}_t + \mathbf{o}_t), -\mathbf{G}_t (\mathbf{R}^+ \bar{\boldsymbol{\alpha}}_t + \mathbf{o}_t)^{-1}\right) \\ &= \mathcal{N}\left(\mathbf{X}_t \boldsymbol{\xi}_t \middle| \mathbf{X}_t \bar{\boldsymbol{\xi}}_t - \mathbf{o}_t + \mathbf{G}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)^{-1} \mathbf{g}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t), -\mathbf{G}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)^{-1}\right) \end{aligned} \quad (5)$$

where

$$\begin{aligned} \mathbf{g}_t(\boldsymbol{\xi}) &= \left\{ \frac{\partial \log P(y_{it} | \eta_{it})}{\partial \eta_{it}} \middle|_{\eta_{it} = \mathbf{x}_{it}^\top \boldsymbol{\xi}} \right\}_{i \in R_t} \\ \mathbf{G}_t(\boldsymbol{\xi}) &= \text{diag} \left(\left\{ \frac{\partial^2 \log P(y_{it} | \eta_{it})}{\partial \eta_{it}^2} \middle|_{\eta_{it} = \mathbf{x}_{it}^\top \boldsymbol{\xi}} \right\}_{i \in R_t} \right) \end{aligned} \quad (6)$$

to get

$$q\left(\boldsymbol{\alpha}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t\right) \propto \mathcal{N}\left(\boldsymbol{\xi}_t \middle| \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{Q}\right) \tilde{g}\left(\mathbf{y}_t \middle| \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)}\right) \quad (7)$$

Thus, we need to sample from

$$\begin{aligned} q\left(\boldsymbol{\alpha}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t\right) &= \mathcal{N}\left(\boldsymbol{\xi}_t \middle| \boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t\right) \\ \bar{\boldsymbol{\xi}}_t &= \mathbf{R}^+ \mathbf{F} \bar{\boldsymbol{\alpha}}_{t-1} \\ \boldsymbol{\Sigma}_t^{-1} &= \mathbf{X}_t^\top (-\mathbf{G}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)) \mathbf{X}_t + \mathbf{Q}^{-1} \\ \boldsymbol{\mu}_t &= \boldsymbol{\Sigma}_t \left(\mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)} + \mathbf{X}_t^\top (-\mathbf{G}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)) (\mathbf{X}_t \bar{\boldsymbol{\xi}}_t - \mathbf{o}_t - \mathbf{G}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)^{-1} \mathbf{g}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)) \right) \\ &= \boldsymbol{\Sigma}_t \left(\mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)} + \mathbf{X}_t^\top (-\mathbf{G}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)) \mathbf{X}_t \bar{\boldsymbol{\xi}}_t - \mathbf{o}_t + \mathbf{g}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t) \right) \end{aligned} \quad (8)$$

We can select $\bar{\boldsymbol{\alpha}}_{t-1}$ as the weighted mean given the particle cloud at time $t-1$ (that is, $\{\boldsymbol{\alpha}_{t-1}^{(1)}, \boldsymbol{\alpha}_{t-1}^{(2)}, \dots, \boldsymbol{\alpha}_{t-1}^{(N)}\}$). This is similar to the second order random walk example in Fearnhead et al. (2010). They use the mode which is feasible as the outcome only depends on one element of the state vector. The downside is an $\mathcal{O}(n_{\max} p^2 + p^3)$ computational cost though independent of the number of particles, N . The total cost of sampling is $\mathcal{O}(n_{\max} p^2 + p^3 + N p^2)$. Another option is to set $\bar{\boldsymbol{\alpha}}_{t-1} = \boldsymbol{\alpha}_{t-1}^{(j)}$ for each particle $j = 1, 2, \dots, N$ in the particle cloud at time $t-1$. This will improve the Taylor expansion but yields an $\mathcal{O}(N(n_{\max} p^2 + p^3))$ computational cost. The extra $N n_{\max}$ factor makes this much slower.

Next, we have the re-sampling weights. A simple solution is not to use an auxiliary particle filter as in the examples of Fearnhead et al. (2010) and set:

$$\beta_t^{(j)} \propto w_{t-1}^{(j)} \quad (9)$$

which has an $\mathcal{O}(N)$ cost of sampling. Another options is to set

$$\begin{aligned} \beta_t^{(j)} &\propto \mathbf{P} \left(\mathbf{y}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)} \right) w_{t-1}^{(j)} \\ &\approx w_{t-1}^{(j)} \int_{\mathbb{R}^r} \mathcal{N} \left(\boldsymbol{\xi}_t \middle| \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{Q} \right) g \left(\mathbf{y}_t \middle| \boldsymbol{\xi}_t \right) \partial \boldsymbol{\xi}_t \\ &\approx w_{t-1}^{(j)} \int_{\mathbb{R}^r} \mathcal{N} \left(\boldsymbol{\xi}_t \middle| \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{Q} \right) \tilde{g}_t \left(\mathbf{y}_t \middle| \boldsymbol{\xi}_t \right) \partial \boldsymbol{\xi}_t \\ &\approx \frac{w_{t-1}^{(j)} \mathcal{N} \left(\boldsymbol{\mu}_t \middle| \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{Q} \right) \tilde{g}_t \left(\mathbf{y}_t \middle| \boldsymbol{\mu}_t \right)}{q \left(\boldsymbol{\mu}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t \right)} \end{aligned} \quad (10)$$

where $q \left(\boldsymbol{\mu}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t \right)$ is meant as the right hand side of the first line in equation (8) computed with $\boldsymbol{\mu}_t$ instead of $\boldsymbol{\xi}_t$. The re-sampling weights will differ as we condition of different particle $\boldsymbol{\alpha}_{t-1}^{(j)}$. This comes at an $\mathcal{O}(Np^2)$ computational cost assuming that we are using (8) already. Otherwise it has the same computational cost as mentioned when I covered the importance density since we have to make the Taylor expansion approximation.

Backward filter (algorithm 3)

We need to specify the artificial prior $\gamma_t(\boldsymbol{\alpha}_t)$. Briers et al. (2010, page 69 and 70) provides recommendation on the selection. This leads to

$$\begin{aligned} \gamma_t(\boldsymbol{\alpha}_t) &= \mathcal{N} \left(\boldsymbol{\alpha}_t \middle| \overleftarrow{\mathbf{m}}_t, \overleftarrow{\mathbf{P}}_t \right) \\ \overleftarrow{\mathbf{m}}_t &= \mathbf{F}^t \mathbf{a}_0 \\ \overleftarrow{\mathbf{P}}_t &= \begin{cases} \mathbf{Q}_0 & t = 0 \\ \mathbf{F} \overleftarrow{\mathbf{P}}_{t-1} \mathbf{F}^\top + \mathbf{R} \mathbf{Q} \mathbf{R}^\top & t > 0 \end{cases} \end{aligned} \quad (11)$$

Following Fearnhead et al. (2010) then we end with

$$\begin{aligned} \mathbf{P}(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t+1}) &= \mathcal{N} \left(\boldsymbol{\alpha}_t \middle| \overleftarrow{\mathbf{a}}_t, \overleftarrow{\mathbf{S}}_t \right) \\ \overleftarrow{\mathbf{S}}_t &= \overleftarrow{\mathbf{P}}_t \mathbf{F}^\top \overleftarrow{\mathbf{P}}_{t+1}^{-1} \mathbf{R} \mathbf{Q} \mathbf{R}^\top (\mathbf{F}^\top)^{-1} \\ \overleftarrow{\mathbf{a}}_t &= \overleftarrow{\mathbf{P}}_t \mathbf{F}^\top \overleftarrow{\mathbf{P}}_{t+1}^{-1} \boldsymbol{\alpha}_{t+1} + \overleftarrow{\mathbf{S}}_t \overleftarrow{\mathbf{P}}_t^{-1} \overleftarrow{\mathbf{m}}_t \end{aligned} \quad (12)$$

which simplifies for the first order random walk to

$$\begin{aligned} \overleftarrow{\mathbf{m}}_t &= \mathbf{a}_0 \\ \overleftarrow{\mathbf{P}}_t &= t \mathbf{Q} + \mathbf{Q}_0 \\ \overleftarrow{\mathbf{S}}_t &= (t \mathbf{Q} + \mathbf{Q}_0) ((t+1) \mathbf{Q} + \mathbf{Q}_0)^{-1} \mathbf{Q} \\ \overleftarrow{\mathbf{a}}_t &= (t \mathbf{Q} + \mathbf{Q}_0) ((t+1) \mathbf{Q} + \mathbf{Q}_0)^{-1} \boldsymbol{\alpha}_{t+1} + \overleftarrow{\mathbf{S}}_t \overleftarrow{\mathbf{P}}_t^{-1} \overleftarrow{\mathbf{m}}_t \end{aligned} \quad (13)$$

Further, setting $\mathbf{Q}_0 = \mathbf{Q}$ (only in the artificial prior where we may alter γ_0 – see Briers et al. (2010, page 70)) gives us

$$\begin{aligned}\overleftarrow{\mathbf{S}}_t &= \frac{t+1}{t+2} \mathbf{Q} \\ \overleftarrow{\mathbf{a}}_t &= \frac{t+1}{t+2} \boldsymbol{\alpha}_{t+1} + \frac{1}{t+2} \mathbf{a}_0\end{aligned}\tag{14}$$

We get the following backward version of equation (8)

$$\begin{aligned}\tilde{q}(\boldsymbol{\alpha}_t | \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) &= \mathcal{N}(\boldsymbol{\alpha}_t | \overleftarrow{\boldsymbol{\mu}}_t, \overleftarrow{\boldsymbol{\Sigma}}_t) \\ \tilde{\boldsymbol{\xi}}_t &= \mathbf{R}^+ \mathbf{F}^{-1} \tilde{\boldsymbol{\alpha}}_{t+1} \\ \overleftarrow{\boldsymbol{\Sigma}}_t^{-1} &= \overleftarrow{\mathbf{P}}_t^{-1} + \mathbf{R} \mathbf{X}_t^\top (-\mathbf{G}_t (\tilde{\boldsymbol{\xi}}_t + \mathbf{o}_t)) \mathbf{X}_t \mathbf{R}^\top + \mathbf{F}^\top \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F} \\ \overleftarrow{\boldsymbol{\mu}}_t &= \overleftarrow{\boldsymbol{\Sigma}}_t \left(\overleftarrow{\mathbf{P}}_t^{-1} \tilde{\mathbf{m}}_t + \mathbf{F}^\top \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} + \mathbf{R} \mathbf{X}_t^\top (-\mathbf{G}_t (\tilde{\boldsymbol{\xi}}_t + \mathbf{o}_t)) \mathbf{X}_t \tilde{\boldsymbol{\xi}}_t - \mathbf{o}_t + \mathbf{g}_t (\tilde{\boldsymbol{\xi}}_t + \mathbf{o}_t) \right)\end{aligned}\tag{15}$$

Notice that the dimension is now $p \geq r$. As an example, we can look at a second order univariate autoregressive model

$$\mathbf{F} = \begin{pmatrix} \theta_1 & \theta_2 \\ 1 & 0 \end{pmatrix}, \quad \mathbf{R} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \mathbf{Q} = (\sigma^2), \quad \mathbf{Q}_0 = \begin{pmatrix} s_1^2 & s_{12} \\ s_{12} & s_2^2 \end{pmatrix}$$

here

$$\mathbf{F}^\top \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} (\sigma^{-2} \quad 0) \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} = \begin{pmatrix} \sigma^{-2} \theta_1 \tilde{\xi}_{t+1}^{(k)} \\ \sigma^{-2} \theta_2 \tilde{\xi}_{t+1}^{(k)} \end{pmatrix}$$

We also find that

$$\mathbf{F}^\top \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F} = \mathbf{F}^\top \mathbf{R} \mathbf{Q}^{-1} \mathbf{R}^\top \mathbf{F} = \begin{pmatrix} \theta_1^2 \sigma^{-2} & \theta_1 \theta_2 \sigma^{-2} \\ \theta_1 \theta_2 \sigma^{-2} & \theta_2^2 \sigma^{-2} \end{pmatrix}$$

Further,

$$\overleftarrow{\mathbf{P}}_t = \mathbf{F}^t \mathbf{Q}_0 \mathbf{F}^{t\top} + \sum_{i=1}^t \mathbf{F}^{(i-1)} \mathbf{R} \mathbf{Q} \mathbf{R}^\top \mathbf{F}^{(i-1)\top} = \mathbf{F}^t \mathbf{Q}_0 \mathbf{F}^{t\top} + \sum_{i=1}^t (\mathbf{F}^{(i-1)})_{\cdot 1} \sigma^2 (\mathbf{F}^{(i-1)})_{\cdot 1}^\top.$$

where subscript $\mathbf{Z}_{k\cdot}$ is the k 'th row vector and $\mathbf{Z}_{\cdot k}$ is the k 'th column vector. Here, the first term has full rank if \mathbf{Q}_0 and \mathbf{F} has full rank, is not sparse, and tends toward zero if the largest eigenvalue of \mathbf{F} is less than one. Further the later terms has rank 1 and are not sparse either. The other terms and matrices only have one non-zero entry due to multiplication with \mathbf{R} or \mathbf{R}^\top which is the first entry for vectors and the (1, 1) entry for matrices. We can use Equation (15) to approximate

$$\begin{aligned}\tilde{q}(\boldsymbol{\alpha}_t | \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) &\propto g(\mathbf{y}_t | \boldsymbol{\xi}_t) f(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} | \boldsymbol{\alpha}_t) \frac{\gamma_t(\boldsymbol{\alpha}_t)}{\gamma_{t+1}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})} \\ &\approx \tilde{g}(\mathbf{y}_t | \boldsymbol{\xi}_t) \mathcal{N}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} | \mathbf{F} \boldsymbol{\alpha}_t, \mathbf{Q}) \frac{\mathcal{N}(\boldsymbol{\alpha}_t | \overleftarrow{\mathbf{m}}_t, \overleftarrow{\mathbf{P}}_t)}{\mathcal{N}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} | \overleftarrow{\mathbf{m}}_{t+1}, \overleftarrow{\mathbf{P}}_{t+1})}\end{aligned}\tag{16}$$

Can we sample from $\mathcal{N}(\cdot | \mathbf{R} \overleftarrow{\boldsymbol{\mu}}_t, \mathbf{R} \overleftarrow{\boldsymbol{\Sigma}}_t \mathbf{R}^\top)$ and set just the remaining $p - r$ entries corresponding to the non-included identity matrix columns in \mathbf{R} to $\overleftarrow{\boldsymbol{\mu}}_t$ for each particle $\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}$?

The re-sampling weights can be computed similarly to the forward filter using the backward transition density in equation (12) in the numerator. We can also use a bootstrap like filter with the importance density by sampling from (12).

Combining / smoothing (algorithm 1)

We need to specify the importance density $\tilde{q}(\cdot | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})$ (see Fearnhead et al. (2010, page 453)). Again, we can make a second order Taylor expansion as in Fearnhead et al. (2010) and choose

$$\begin{aligned} \tilde{q}(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) &= \mathcal{N}(\boldsymbol{\alpha}_t | \overleftarrow{\boldsymbol{\mu}}_t, \overleftarrow{\boldsymbol{\Sigma}}_t) \\ \overleftarrow{\boldsymbol{\Sigma}}_t^{-1} &= \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ + \mathbf{R} \mathbf{X}_t^\top (-\mathbf{G}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)) \mathbf{X}_t \mathbf{R}^\top + \mathbf{F}^\top \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F} \\ \overleftarrow{\boldsymbol{\mu}}_t &= \overleftarrow{\boldsymbol{\Sigma}}_t (\mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)} + \mathbf{F}^\top \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} + \mathbf{R} \mathbf{X}_t^\top (-\mathbf{G}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)) \mathbf{X}_t \bar{\boldsymbol{\xi}}_t - \mathbf{o}_t + \mathbf{g}_t (\bar{\boldsymbol{\xi}}_t + \mathbf{o}_t)) \end{aligned} \quad (17)$$

where $\bar{\boldsymbol{\xi}}_t$ can be a combined mean given the cloud means at time $t-1$ and $t+1$ or a mean for each of the two drawn particles in the (j, k_i) pairs. This is to approximate

$$\tilde{q}(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) \propto \tilde{g}(\mathbf{y}_t | \boldsymbol{\xi}_t) f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}) \frac{f(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} | \boldsymbol{\alpha}_t)}{\gamma_{t+1}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})} \quad (18)$$

Can we sample from $\mathcal{N}(\cdot | \mathbf{R} \overleftarrow{\boldsymbol{\mu}}_t, \mathbf{R} \overleftarrow{\boldsymbol{\Sigma}}_t^{-1} \mathbf{R}^\top)$ and set just the remaining $p - r$ entries corresponding to the non-included identity matrix columns in \mathbf{R} to $\overleftarrow{\boldsymbol{\mu}}_t$ for each particle pair?

We can also use a bootstrap like filter by sampling from

$$\begin{aligned} \tilde{q}(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) &= \mathcal{N}(\boldsymbol{\alpha}_t | \widetilde{\mathbf{m}}, \widetilde{\mathbf{S}}) \\ \widetilde{\mathbf{S}}^{-1} &= (\mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ + \mathbf{F}^\top \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F})^{-1} \\ \widetilde{\mathbf{m}} &= \widetilde{\mathbf{S}} (\mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)} + \mathbf{F}^\top \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) \end{aligned}$$

$\overleftarrow{\Sigma}_t$ may not have full rank. Take E.g, a third order univariate model. In this case

$$\begin{aligned}\mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ &= \begin{pmatrix} \sigma^{-2} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \\ \mathbf{R} \mathbf{X}_t^\top (-\mathbf{G}_t(\bar{\xi}_t)) \mathbf{X}_t \mathbf{R}^\top &= \begin{pmatrix} z & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \\ \mathbf{F} \mathbf{R}^{+\top} \mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F} &= \sigma^{-2} \begin{pmatrix} \theta_1^2 & \theta_1 \theta_2 & \theta_1 \theta_3 \\ \theta_1 \theta_2 & \theta_2^2 & \theta_2 \theta_3 \\ \theta_1 \theta_3 & \theta_2 \theta_3 & \theta_3^2 \end{pmatrix}\end{aligned}$$

The latter only has rank 1. I gather the current solution requires a full rank matrix...

Algorithm 1 $\mathcal{O}(N)$ generalized two filter smoother using the method in Fearnhead et al. (2010).

Input:

$\mathbf{Q}, \mathbf{Q}_0, \mathbf{a}_0, \mathbf{X}_1, \dots, \mathbf{X}_d, \mathbf{y}_1, \dots, \mathbf{y}_d, R_1, \dots, R_d, \mathbf{S}$

Importance density which optimally is (see Fearnhead et al. (2010, page 453)):

$$\tilde{q}\left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right) = \mathrm{P}\left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right) \quad (19)$$

Let $\boldsymbol{\alpha}_t^{(i)}$ denote particle i at time t , $w_t^{(i)}$ denote the weight of the particle and $\beta_t^{(i)}$ denote the re-sampling weight.

1: **procedure** FILTER FORWARD

2: Run a forward particle filter to get a particle clouds

$\left\{\boldsymbol{\alpha}_t^{(j)}, w_t^{(j)}, \beta_{t+1}^{(j)}\right\}_{j=1, \dots, N}$ approximating $\mathrm{P}\left(\boldsymbol{\alpha}_t \mid \mathbf{y}_{1:t}\right)$ for $t = 0, 1, \dots, d$. See algorithm 2.

3: **procedure** FILTER BACKWARDS

4: Run a similar backward filter to get $\left\{\tilde{\boldsymbol{\alpha}}_t^{(k)}, \tilde{w}_t^{(k)}, \tilde{\beta}_{t-1}^{(k)}\right\}_{k=1, \dots, d}$ approximating $\mathrm{P}\left(\boldsymbol{\alpha}_t \mid \mathbf{y}_{t:d}\right)$ for $t = d+1, d, d-1, \dots, 1$. See algorithm 3.

5: **procedure** SMOOTH (COMBINE)

6: **for** $t = 1, \dots, N$ **do**

Re-sample:

7: $i = 1, 2, \dots, N_s$ pairs of (j_i, k_i) where each component is independently sampled using re-sampling weights $\beta_t^{(j)}$ and $\tilde{\beta}_t^{(k)}$.

Propagate:

8: Sample particles $\hat{\boldsymbol{\alpha}}_t^{(i)}$ from importance density $\tilde{q}\left(\cdot \mid \boldsymbol{\alpha}_{t-1}^{(j_i)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}\right)$.

Re-weight:

9: Assign each particle weights:

$$\hat{w}_t^{(i)} \propto \frac{\mathrm{P}\left(\hat{\boldsymbol{\alpha}}_t^{(i)} \mid \boldsymbol{\alpha}_{t-1}^{(j_i)}\right) \mathrm{P}\left(\mathbf{y}_t \mid \hat{\boldsymbol{\alpha}}_t^{(i)}\right) \mathrm{P}\left(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)} \mid \hat{\boldsymbol{\alpha}}_t^{(i)}\right) w_{t-1}^{(j_i)} \tilde{w}_{t+1}^{(k_i)}}{\tilde{q}\left(\hat{\boldsymbol{\alpha}}_t^{(i)} \mid \boldsymbol{\alpha}_{t-1}^{(j_i)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}\right) \beta_t^{(j_i)} \tilde{\beta}_t^{(k_i)} \gamma_{t+1}\left(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}\right)} \quad (20)$$

Algorithm 2 Forward filter due to Pitt and Shephard (1999). You can compare with Doucet and Johansen (2009, page 20 and 25). The version and notation below is from Fearnhead et al. (2010, page 449).

Input:

Importance density and specification of weights are optimally given by:

$$\begin{aligned} q\left(\boldsymbol{\alpha}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t\right) &= P\left(\boldsymbol{\alpha}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t\right) \\ \beta_t^{(j)} &\propto P\left(\mathbf{y}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j)}\right) w_{t-1}^{(j)} \end{aligned} \quad (21)$$

- 1: Sample $\boldsymbol{\alpha}_0^{(1)}, \dots, \boldsymbol{\alpha}_0^{(N_f)}$ particles from $\mathcal{N}(\cdot | \mathbf{a}_0, \mathbf{Q}_0)$ and set the weights $w_0^{(1)}, \dots, w_0^{(N_f)}$ to $1/N_f$.
- 2: **for** $t = 1, \dots, d$ **do**
- 3: **procedure** RE-SAMPLE
- 4: Compute re-sampling weights $\beta_t^{(j)}$ and re-sample according to $\beta_t^{(j)}$ to get indices j_1, \dots, j_N . If we do not re-sample then set $\beta_t^{(j)} = 1$.
- 5: **procedure** PROPAGATE
- 6: Sample new particles $\boldsymbol{\alpha}_t^{(i)}$ using importance density $q\left(\boldsymbol{\alpha}_t \middle| \boldsymbol{\alpha}_{t-1}^{(j_i)}, \mathbf{y}_t\right)$.
- 7: **procedure** RE-WEIGHT
- 8: Re-weight particles using $(w_{t-1}^{(j_i)} / \beta_t^{(j_i)})$ is added due to the auxiliary particle filter):

$$w_t^{(i)} \propto \frac{P\left(\mathbf{y}_t \middle| \boldsymbol{\alpha}_t^{(i)}\right) P\left(\boldsymbol{\alpha}_t^{(i)} \middle| \boldsymbol{\alpha}_{t-1}^{(j_i)}\right) w_{t-1}^{(j_i)}}{q\left(\boldsymbol{\alpha}_t^{(i)} \middle| \boldsymbol{\alpha}_{t-1}^{(j_i)}, \mathbf{y}_t\right) \beta_t^{(j_i)}} \quad (22)$$

Algorithm 3 Backwards filter. See Briers et al. (2010) and Fearnhead et al. (2010).

Input:

A backwards filter distribution approximation:

$$\tilde{p}(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d}) \propto \gamma_t(\boldsymbol{\alpha}_t) P(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_t) \quad (23)$$

with an artificial prior distribution $\gamma_t(\boldsymbol{\alpha}_t)$. This is introduced as $P(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_t)$ is not a density function in $\boldsymbol{\alpha}_t$.

Importance density and specification of weights (Fearnhead et al. (2010, page 451 – look in the example in the appendix)):

$$\tilde{q}(\boldsymbol{\alpha}_t | \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) \tilde{\beta}_t^{(k)} \approx \gamma_t(\boldsymbol{\alpha}_t) P(\mathbf{y}_t | \boldsymbol{\alpha}_t) P(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} | \boldsymbol{\alpha}_t) \frac{\tilde{w}_{t+1}^{(k)}}{\gamma_{t+1}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})} \quad (24)$$

where we want (see Briers et al. (2010, page 74)):

$$\begin{aligned} \tilde{q}(\boldsymbol{\alpha}_t | \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) &\propto P(\mathbf{y}_t | \boldsymbol{\alpha}_t) P(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} | \boldsymbol{\alpha}_t) \frac{\gamma_t(\boldsymbol{\alpha}_t)}{\gamma_{t+1}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})} \\ \tilde{\beta}_t^{(k)} &\propto \tilde{p}(\mathbf{y}_t | \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) \tilde{w}_{t+1}^{(k)} \end{aligned} \quad (25)$$

- 1: Sample $\tilde{\boldsymbol{\alpha}}_{d+1}^{(1)}, \dots, \tilde{\boldsymbol{\alpha}}_{d+1}^{(N_f)}$ particles from $\gamma_{d+1}(\cdot)$ and set the weights $\tilde{w}_{d+1}^{(1)}, \dots, \tilde{w}_{d+1}^{(N_f)}$ to $1/N_f$.
- 2: **for** $t = d, \dots, 1$ **do**
- 3: **procedure** RE-SAMPLE
- 4: Compute re-sampling weights $\tilde{\beta}_t^{(k)}$ and re-sample according to $\tilde{\beta}_t^{(k)}$ to get indices k_1, \dots, k_N . If we do not re-sample then set $\tilde{\beta}_t^{(k)} = 1$.
- 5: **procedure** PROPAGATE
- 6: Sample new particles $\tilde{\boldsymbol{\alpha}}_t^{(i)}$ using importance density $\tilde{q}(\boldsymbol{\alpha}_t | \tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}, \mathbf{y}_t)$.
- 7: **procedure** RE-WEIGHT
- 8: Re-weight particles using (see Briers et al. (2010, page 72) and $\tilde{w}_{t+1}^{(k_i)} / \beta_t^{(k_i)}$ is added due to the auxiliary particle filter):

$$\tilde{w}_t^{(i)} \propto \frac{P(\mathbf{y}_t | \tilde{\boldsymbol{\alpha}}_t^{(i)}) P(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)} | \tilde{\boldsymbol{\alpha}}_t^{(i)}) \gamma_t(\tilde{\boldsymbol{\alpha}}_t^{(i)}) \tilde{w}_{t+1}^{(k_i)}}{q(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}, \mathbf{y}_t) \gamma_{t+1}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}) \beta_t^{(k_i)}} \quad (26)$$

2 Log likelihood evaluation

We can evaluate the log likelihood for a particular value of $\theta = \{\mathbf{Q}, \mathbf{Q}_0, \mathbf{a}_0, \mathbf{F}\}$ as described in Doucet and Johansen (2009, page 5) and Malik and Pitt (2011, page 193) using the forward particle filter shown in algorithm 2.

3 Parameter estimation

This section shows an example of parameter estimation in a random walk. The formulas for parameter estimation with a given order random with the EM-algorithm (Dempster et al., 1977) are particularly simple. We need to estimate \mathbf{Q} and \mathbf{a}_0 elements of $\theta = \{\mathbf{Q}, \mathbf{Q}_0, \mathbf{a}_0, \mathbf{F}\}$. We do this by running algorithm 1 for the current θ . Then we compute the so-called *smoothed additive functional* or summary statistics

$$\begin{aligned} t_t^{(\theta)} &= \int_{\alpha_t} \alpha_t P_{\theta}(\alpha_t | \mathbf{y}_{1:d}) \partial \alpha_t \\ &\approx \sum_{i=1}^{N_s} \hat{\alpha}_t^{(i)} \hat{w}_t^{(i)} \approx \sum_{i=1}^{N_s} \alpha_t^{(j_i)} \hat{w}_{t+1}^{(i)} \approx \sum_{i=1}^{N_s} \tilde{\alpha}_t^{(k_i)} \hat{w}_{t-1}^{(i)} \\ T_t^{(\theta)} &= \int_{\mathbb{R}^{2p}} (\alpha_t - \mathbf{F}\alpha_{t-1}) (\alpha_t - \mathbf{F}\alpha_{t-1})^{\top} P_{\theta}(\alpha_{(t-1):t} | \mathbf{y}_{1:d}) \partial \alpha_{(t-1):t} \quad (27) \\ &\approx \sum_{i=1}^{N_s} \left(\hat{\alpha}_t^{(i)} - \mathbf{F}\hat{\alpha}_{t-1}^{(j_i)} \right) \left(\hat{\alpha}_t^{(i)} - \mathbf{F}\hat{\alpha}_{t-1}^{(j_i)} \right)^{\top} \hat{w}_t^{(i)} \\ &\approx \sum_{i=1}^{N_s} \left(\tilde{\alpha}_t^{(k_i)} - \mathbf{F}\tilde{\alpha}_{t-1}^{(i)} \right) \left(\tilde{\alpha}_t^{(k_i)} - \mathbf{F}\tilde{\alpha}_{t-1}^{(i)} \right)^{\top} \hat{w}_{t-1}^{(i)} \end{aligned}$$

where $\alpha_{s:t} = \{\alpha_s, \alpha_{s+1}, \dots, \alpha_t\}$ and the subscript in P denotes that it is the probability given the parameter θ . I use that the normalized weights $\hat{w}_t^{(i)}$ and the pairs $\{\alpha_{t-1}^{(j_i)}, \hat{\alpha}_t^{(i)}, \tilde{\alpha}_{t+1}^{(k_i)}\}$ form a discrete approximation of $P_{\theta}(\alpha_{(t-1):(t+1)} | \mathbf{y}_{1:d})$. That is,

$$\begin{aligned} P_{\theta}(\alpha_{(t-1):t+1} | \mathbf{y}_{1:d}) \\ \approx \sum_{i=1}^{N_s} \hat{w}_t^{(i)} \delta(\alpha_{t-1} - \alpha_{t-1}^{(j_i)}) \delta(\alpha_t - \hat{\alpha}_t^{(i)}) \delta(\alpha_{t+1} - \tilde{\alpha}_{t+1}^{(k_i)}) \end{aligned} \quad (28)$$

where $\delta(\cdot)$ is the Dirac delta function. The update of \mathbf{a}_0 and \mathbf{Q} given the summary statistics is

$$\mathbf{a}_0 = t_0^{(\theta)} \quad \mathbf{Q} = \frac{1}{d} \sum_{t=1}^d \mathbf{R}^+ T_t^{(\theta)} \mathbf{R}^{+\top} \quad (29)$$

We then repeat with the new \mathbf{a}_0 and \mathbf{Q} for a given number of iterations or till a convergence criteria is satisfied. See Kantas et al. (2015), Del Moral et al.

(2010) and Schön et al. (2011) for further details on parameter estimation with particle filters.

We can do parts of the probability estimates exact at time 1 by using

$$\begin{aligned}
P_{\theta}(\alpha_1, \alpha_0 | \mathbf{y}_{1:d}) &= P_{\theta}(\alpha_0 | \alpha_1, \mathbf{y}_{1:d}) P_{\theta}(\alpha_1 | \mathbf{y}_{1:d}) \\
&= \frac{P_{\theta}(\alpha_1, \mathbf{y}_{1:d} | \alpha_0) P_{\theta}(\alpha_0)}{P_{\theta}(\alpha_1, \mathbf{y}_{1:d})} P_{\theta}(\alpha_1 | \mathbf{y}_{1:d}) \\
&= \frac{P_{\theta}(\mathbf{y}_{1:d} | \alpha_1, \alpha_0) f_{\theta}(\alpha_1 | \alpha_0) P_{\theta}(\alpha_0)}{P_{\theta}(\alpha_1) P_{\theta}(\mathbf{y}_{1:d} | \alpha_1)} P_{\theta}(\alpha_1 | \mathbf{y}_{1:d}) \\
&= \frac{f_{\theta}(\alpha_1 | \alpha_0) P_{\theta}(\alpha_0)}{P_{\theta}(\alpha_1)} P_{\theta}(\alpha_1 | \mathbf{y}_{1:d}) \\
&= P_{\theta}(\alpha_0 | \alpha_1) P_{\theta}(\alpha_1 | \mathbf{y}_{1:d}) \\
&= \mathcal{N}(\alpha_0 | \mathbf{m}_{\theta}(\alpha_1), \mathbf{S}_{\theta}) P_{\theta}(\alpha_1 | \mathbf{y}_{1:d}) \\
&\approx \sum_{i=1}^{N_s} \mathcal{N}(\alpha_0 | \mathbf{m}_{\theta}(\hat{\alpha}_1^{(i)}), \mathbf{S}_{\theta}) \hat{w}_t^{(i)} \delta(\alpha_1 - \hat{\alpha}_1^{(i)})
\end{aligned} \tag{30}$$

where we use that $P_{\theta}(\mathbf{y}_{1:d} | \alpha_1, \alpha_0) = P_{\theta}(\mathbf{y}_{1:d} | \alpha_1)$ and let

$$\begin{aligned}
\mathbf{S}_{\theta}^{-1} &= \mathbf{Q}_0^{-1} + \mathbf{F}^{-1} \mathbf{R} \mathbf{Q}^{-1} \mathbf{R}^{\top} \mathbf{F}^{-\top} \\
\mathbf{m}_{\theta}(\alpha) &= \mathbf{S}_{\theta} (\mathbf{Q}_0^{-1} \mathbf{a}_0 + \mathbf{F}^{-1} \mathbf{R} \mathbf{Q}^{-1} \mathbf{R}^{\top} \mathbf{a})
\end{aligned} \tag{31}$$

Thus, we end with

$$\begin{aligned}
\mathbf{T}_1^{(\theta)} &\approx \sum_{i=1}^{N_s} \hat{w}_1^{(i)} \left(\hat{\alpha}_1^{(i)} \left(\hat{\alpha}_1^{(i)} \right)^{\top} \right. \\
&\quad \left. - \mathbf{F} \mathbf{m}_{\theta}(\hat{\alpha}_1^{(i)}) \left(\hat{\alpha}_1^{(i)} \right)^{\top} - \hat{\alpha}_1^{(i)} \mathbf{m}_{\theta}(\hat{\alpha}_1^{(i)})^{\top} \mathbf{F}^{\top} \right. \\
&\quad \left. + \mathbf{F} \mathbf{m}_{\theta}(\hat{\alpha}_1^{(i)}) \mathbf{m}_{\theta}(\hat{\alpha}_1^{(i)})^{\top} \mathbf{F}^{\top} \right) + \mathbf{S}_{\theta}
\end{aligned} \tag{32}$$

4 Other options

The $\mathcal{O}(N^2)$ two-filter smoother in Fearnhead et al. (2010) is going to be computationally expensive as an approximation is going to be needed for equation (8) in the article. For instance, only the Taylor expansion approximation around a single point to approximate g would be feasible. The non-auxiliary version in Briers et al. (2010) is more feasible as it only requires evaluation of f in the smoothing part of two-filter smoother (see equation (46) in the paper). Similar conclusions applies to the forward smoother in Del Moral et al. (2010) and the backward smoother as presented in Kantas et al. (2015). Both have a $\mathcal{O}(N^2)$ computational cost.

Despite the $\mathcal{O}(N^2)$ cost of the method in Briers et al. (2010) and Del Moral et al. (2010) they are still worthy candidates as the computational cost is independent of the number of observations, n . Further, the computational cost can be reduced to $\mathcal{O}(N \log(N))$ with the approximations in Klaas et al. (2006).

The method in Malik and Pitt (2011, see particularly section 6.2 on page 203) can be used to do continuous likelihood evaluation. I am not sure how well these method scale with higher state dimension, p .

Kantas et al. (2015) show empirically that it may be worth just using a forward filter. However, the example is with a univariate outcome ($n = 1$ – not to be confused with the number of periods d). The cost here of the forward filter is at least $\mathcal{O}(dNn_{\max}p)$. Every new particle yields an $\mathcal{O}(dn_{\max}p)$ cost which is expensive due to the large number of outcomes, n . Thus, the considerations are different and a $\mathcal{O}(dNn_{\max}p + N^2)$ method will not make a big difference unless N is large. Another alternative is to add noise to θ_t and use the methods in Andrieu and Doucet (2002).

5 Briers et al. (2010)

The $\mathcal{O}(N^2)$ smother from Briers et al. (2010) is also implemented as it is feasible for a moderate number of particles (though, we can use the approximations in Kantas et al. (2015) to reduce the computational complexity). It is shown in algorithm 4. The weights in equation (36) comes from the two-filter formula:

$$\begin{aligned}
P(\alpha_t | \mathbf{y}_{1:d}) &= \frac{P(\alpha_t | \mathbf{y}_{1:t-1}) P(\mathbf{y}_{t:d} | \alpha_t)}{P(\mathbf{y}_{t:d} | \mathbf{y}_{1:t-1})} \\
&\propto P(\alpha_t | \mathbf{y}_{1:t-1}) P(\mathbf{y}_{t:d} | \alpha_t) \\
&= P(\alpha_t | \mathbf{y}_{1:t-1}) \frac{P(\alpha_t | \mathbf{y}_{t:d}) P(\mathbf{y}_{t:d})}{P(\alpha_t)} \\
&\propto P(\alpha_t | \mathbf{y}_{1:t-1}) \frac{P(\alpha_t | \mathbf{y}_{t:d})}{P(\alpha_t)} \\
&= P(\alpha_t | \mathbf{y}_{t:d}) \frac{\left[\int_{\alpha_{t-1}} P(\alpha_{t-1} | \mathbf{y}_{1:t-1}) f(\alpha_t | \alpha_{t-1}) \partial \alpha_{t-1} \right]}{P(\alpha_t)} \\
&\approx \sum_{i=1}^N \tilde{w}_t^{(i)} \delta(\alpha_t - \tilde{\alpha}_t^{(i)}) \frac{\left[\sum_{j=1}^N w_{t-1}^{(j)} f(\tilde{\alpha}_t^{(i)} | \alpha_{t-1}^{(j)}) \right]}{\gamma_t(\tilde{\alpha}_t^{(i)})}
\end{aligned} \tag{33}$$

Similar arguments leads to:

$$\begin{aligned}
& P(\boldsymbol{\alpha}_t, \boldsymbol{\alpha}_{t-1} | \mathbf{y}_{1:d}) \\
& \propto P(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d}) \frac{P(\boldsymbol{\alpha}_{t-1} | \mathbf{y}_{1:t-1}) f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1})}{P(\boldsymbol{\alpha}_t)} \\
& \approx \sum_{i=1}^N \sum_{j=1}^N \tilde{w}_t^{(i)} \delta(\boldsymbol{\alpha}_t - \tilde{\boldsymbol{\alpha}}_t^{(i)}) \frac{\left[\sum_{j=1}^N w_{t-1}^{(j)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)}) \right]}{\gamma_t(\tilde{\boldsymbol{\alpha}}_t^{(i)})} \frac{w_{t-1}^{(j)} \delta(\boldsymbol{\alpha}_{t-1} - \boldsymbol{\alpha}_{t-1}^{(j)}) f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)})}{\left[\sum_{j=1}^N w_{t-1}^{(j)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)}) \right]} \\
& = \sum_{i=1}^N \sum_{j=1}^N \hat{w}_t^{(i,j)} \delta(\boldsymbol{\alpha}_t - \tilde{\boldsymbol{\alpha}}_t^{(i)}) \delta(\boldsymbol{\alpha}_{t-1} - \boldsymbol{\alpha}_{t-1}^{(j)})
\end{aligned} \tag{34}$$

where

$$\hat{w}_t^{(i,j)} = \tilde{w}_t^{(i)} \frac{w_{t-1}^{(j)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)})}{\left[\sum_{j=1}^N w_{t-1}^{(j)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)}) \right]} \tag{35}$$

The above is what we need for the EM-algorithm.

Algorithm 4 $\mathcal{O}(N^2)$ generalized two filter smoother using the method in Briers et al. (2010).

Input:

$\mathbf{Q}, \mathbf{Q}_0, \mathbf{a}_0, \mathbf{X}_1, \dots, \mathbf{X}_d, \mathbf{y}_1, \dots, \mathbf{y}_d, R_1, \dots, R_d, \mathbf{S}$

1: **procedure** FILTER FORWARD

2: Run a forward particle filter to get a particle clouds

$\left\{ \boldsymbol{\alpha}_t^{(j)}, w_t^{(j)}, \beta_{t+1}^{(j)} \right\}_{j=1, \dots, N}$ approximating $P(\boldsymbol{\alpha}_t | \mathbf{y}_{1:t})$ for $t = 0, 1, \dots, d$. See algorithm 2.

3: **procedure** FILTER BACKWARDS

4: Run a similar backward filter to get $\left\{ \tilde{\boldsymbol{\alpha}}_t^{(k)}, \tilde{w}_t^{(k)}, \tilde{\beta}_{t-1}^{(k)} \right\}_{k=1, \dots, N}$ approximating $P(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d})$ for $t = d+1, d, d-1, \dots, 1$. See algorithm 3.

5: **procedure** SMOOTH (COMBINE)

6: **for** $t = 1, \dots, d$ **do**

7: Assign each backward filter particle a smoothing weight given by:

$$\hat{w}_t^{(i)} \propto \tilde{w}_t^{(i)} \frac{\left[\sum_{j=1}^N w_{t-1}^{(j)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)}) \right]}{\gamma_t(\tilde{\boldsymbol{\alpha}}_t^{(i)})} \tag{36}$$

6 Implementation

The `PF_EM` method in the `dynamichazard` package contains an implementation of the above described method. You specify the number of particles by the `N_first`, `N_fw_n_bw` and `N_smooth` argument for respectively the N_f , N and N_s in the algorithm 1-3. We may want more particle in the smoothing step, $N_s > N$, as pointed out in the discussion in Fearnhead et al. (2010, page 460 and 461). Further, selecting $N_f > N$ may be preferable to ensure coverage of the state space at time 0 and $d + 1$.

The `method` argument specify how the filters are set up. The argument can take the following values:

- `"bootstrap_filter"` for a bootstrap filter.
- `"PF_normal_approx_w_cloud_mean"` and `"AUX_normal_approx_w_cloud_mean"` for the Taylor expansion normal approximation of the conditional density of \mathbf{y}_t made around the weighted mean of the previous cloud. The PF and AUX prefix specifies whether or not the auxiliary version should be used.
- `"PF_normal_approx_w_particles"` and `"AUX_normal_approx_w_particles"` for the Taylor expansion normal approximation of the conditional density of \mathbf{y}_t made around the parent (or/and child) particle. The PF_ and AUX_ prefix specifies whether or not the auxiliary version should be used.

The smoother is selected with the `smoother` argument. `"Fearnhead_0_N"` gives the smoother in algorithm 1 and `"Brier_0_N_square"` gives the smoother in algorithm 4.

The *Systematic Resampling* (Kitagawa, 1996) is used in all re-sampling steps. See Douc and Cappé (2005) for a comparison of re-sampling methods. The rest of the arguments to `PF_EM` are similar to those of the `ddhazard` function.

6.1 Linear maps

The methods describe above involves many linear maps. These are implemented with C++ abstract classes with specialized `map` member functions for particular problem to decrease the computation. An alternative would have been to use a sparse matrix implementation. A.s. an example we have a mapping matrix \mathbf{A} which C++ abstract member on the main data object used in the package called `xyz`. E.g., this could \mathbf{F} with name `state_trans`. Then the following operations are implemented

- `map()`: Returns \mathbf{A} .
- `map(const arma::vec &x, bool tranpose)`: Returns $\mathbf{A}^\top \mathbf{x}$ if `tranpose == true` and $\mathbf{A}\mathbf{x}$ otherwise.
- `map(const arma::mat &X, side s, bool tranpose)`: Let $\mathbf{B} = \mathbf{A}^\top$ if `tranpose == true` and otherwise $\mathbf{B} = \mathbf{A}$. Then the result is $\mathbf{B}\mathbf{X}\mathbf{B}^\top$ if `s == both`, $\mathbf{B}\mathbf{X}$ if `s == left`, and $\mathbf{X}\mathbf{B}^\top$ if `s == right`.

These are implemented for

C++ member name	Matrix A
<code>err_state</code>	\mathbf{R}
<code>err_state_inv</code>	$\mathbf{R}^+ = \mathbf{R}^\top$
<code>state_trans</code>	\mathbf{F}
<code>state_trans_err</code>	$\mathbf{R}^+ \mathbf{F} = \mathbf{R}^\top \mathbf{F}$
<code>state_trans_inv</code>	\mathbf{F}^{-1}

Further, we will need function to compute terms

$$\overleftarrow{\mathbf{P}}_t \mathbf{F}^\top \overleftarrow{\mathbf{P}}_{t+1}^{-1} \alpha_{t+1} + \overleftarrow{\mathbf{S}}_t \overleftarrow{\mathbf{P}}_t^{-1} \overleftarrow{\mathbf{m}}_t$$

for an arbitrary α_{t+1} , $\overleftarrow{\mathbf{S}}_t = \overleftarrow{\mathbf{P}}_t \mathbf{F}^\top \overleftarrow{\mathbf{P}}_{t+1}^{-1} \mathbf{R} \mathbf{Q} \mathbf{R}^\top (\mathbf{F}^\top)^{-1}$, $\overleftarrow{\mathbf{P}}_t^{-1} \overleftarrow{\mathbf{m}}_t$, and $\overleftarrow{\mathbf{P}}_t^{-1}$ for equation (12) and (15). This is done with the methods `bw_mean(signed int, const arma::vec&)`, `bw_covar(signed int)`, `uncond_mean_term(signed int)`, and `uncond_covar_inv(signed int)`. The terms and factors that can will computed once using equation (12).

6.2 Approximations

The approximations in equation (8), (15) and (17) can be done as follow

1. Taylor develop around a given value ξ_t and compute $\mathbf{k} = \mathbf{X}_t^\top (-\mathbf{G}_t(\xi_t) \mathbf{X}_t \bar{\xi}_t + \mathbf{g}_t(\xi_t))$, $\mathbf{K} = \mathbf{X}_t^\top (-\mathbf{G}_t(\xi_t)) \mathbf{X}_t$.
2. Given one or two pairs of means and covariance matrices $(\mathbf{S}_1, \mathbf{a}_1) \in \mathbb{R}^{p \times p} \times \mathbb{R}^p$ and $(\mathbf{S}_2, \mathbf{a}_2) \in \mathbb{R}^{p \times p} \times \mathbb{R}^p$ then compute

$$\begin{aligned} \Sigma^{-1} &= \mathbf{S}_1^{-1} + \mathbf{S}_2^{-1} + \mathbf{R} \mathbf{K} \mathbf{R}^\top \\ \mu &= \Sigma (\mathbf{S}_1^{-1} \mathbf{a}_1 + \mathbf{S}_2^{-1} \mathbf{a}_2 + \mathbf{R} \mathbf{k}) \end{aligned}$$

So a common task in Equation (15) and (17) is to compute Σ^{-1} and \mathbf{k} given one or two $(\mathbf{S}_i, \mathbf{a}_i)$ pairs. Equation (8) only the needs the $r \times r$ dimensional output since there is $p - r$ elements that are constant. Thus, an `is_forward` flag is used to indicate whether the r dimensional output is returned. The function `taylor_normal_approx` returns Σ^{-1} and $\Sigma \mathbf{R} \mathbf{k}$ (the latter to save some computation later). If `is_forward = true` then the Σ is as in Equation (8) and $\Sigma \mathbf{k}$.

References

- Christophe Andrieu and Arnaud Doucet. Particle filtering for partially observed gaussian state space models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(4):827–836, 2002.
- Mark Briers, Arnaud Doucet, and Simon Maskell. Smoothing algorithms for state-space models. *Annals of the Institute of Statistical Mathematics*, 62(1): 61–89, 2010.
- Pierre Del Moral, Arnaud Doucet, and Sumeetpal Singh. Forward smoothing using sequential monte carlo. *arXiv preprint arXiv:1012.5390*, 2010.
- Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society. Series B (methodological)*, pages 1–38, 1977.
- Randal Douc and Olivier Cappé. Comparison of resampling schemes for particle filtering. In *Image and Signal Processing and Analysis, 2005. ISPA 2005. Proceedings of the 4th International Symposium on*, pages 64–69. IEEE, 2005.
- Arnaud Doucet and Adam M Johansen. A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of nonlinear filtering*, 12(656-704): 3, 2009.
- Paul Fearnhead, David Wyncoll, and Jonathan Tawn. A sequential smoothing algorithm with linear computational cost. *Biometrika*, 97(2):447–464, 2010.
- Nikolas Kantas, Arnaud Doucet, Sumeetpal S Singh, Jan Maciejowski, Nicolas Chopin, et al. On particle methods for parameter estimation in state-space models. *Statistical science*, 30(3):328–351, 2015.
- Genshiro Kitagawa. Monte carlo filter and smoother for non-gaussian nonlinear state space models. *Journal of computational and graphical statistics*, 5(1): 1–25, 1996.
- Mike Klaas, Mark Briers, Nando De Freitas, Arnaud Doucet, Simon Maskell, and Dustin Lang. Fast particle smoothing: If i had a million particles. In *Proceedings of the 23rd international conference on Machine learning*, pages 481–488. ACM, 2006.
- Sheheryar Malik and Michael K Pitt. Particle filters for continuous likelihood evaluation and maximisation. *Journal of Econometrics*, 165(2):190–209, 2011.
- Michael K Pitt and Neil Shephard. Filtering via simulation: Auxiliary particle filters. *Journal of the American statistical association*, 94(446):590–599, 1999.
- Thomas B Schön, Adrian Wills, and Brett Ninness. System identification of nonlinear state-space models. *Automatica*, 47(1):39–49, 2011.