

Particle filters and smoothers in the dynamichazard package

BENJAMIN CHRISTOFFERSEN

AUGUST 20, 2019

This vignette covers the particle filters and smoothers implemented in the `dynamichazard` package in `R`. Some prior knowledge of particle filters is assumed. Doucet and Johansen (2009) provide a tutorial on particle filters and Kantas et al. (2015) cover parameter estimation with particle filters. See also Cappé et al. (2005) for a general introduction to Hidden Markov models. This vignette relies heavily on Fearnhead et al. (2010) and there is a big overlap between what is presented here and the paper.

1 Method

The models implemented in the package is survival analysis models for terminal events. These can be in discrete time where we have binary indicators $Y_{ik} = 1_{\{T_i \in (t_{k-1}, t_k]\}}$ which is one if the random event time of individual i denoted by $T_i \in (0, \infty)$ is in the interval $(t_{k-1}, t_k]$ and zero otherwise. It can also be in continuous time where we model the distribution of the event time of individual i , T_i , with a piecewise exponential distribution conditional on observable covariates and the path of a discrete latent variable. To be more concrete, the model is

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

where I denote the conditional densities as $g_t(\mathbf{y}_t|\boldsymbol{\alpha}_t) = g(\mathbf{y}_t|\mathbf{X}_t \mathbf{R}^+ \boldsymbol{\alpha}_t + \mathbf{o}_t + \mathbf{Z}_t \boldsymbol{\omega})$ and $f(\boldsymbol{\alpha}_t|\boldsymbol{\alpha}_{t-1})$. For each $t = 1, \dots, d$, we have risk set given by $R_t \subseteq \{1, 2, \dots, n\}$. Further, we let $n_t = |R_t|$ denote the number of observation at risk at time t and $n_{\max} = \max_{t \in \{1, \dots, d\}} n_t$. 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 state vector containing the time-varying coefficients. The \mathbf{Z}_t is the design matrix for the covariates with time-invariant coefficients and $\boldsymbol{\omega}$ are the corresponding coefficients.

The i 'th row of \mathbf{X}_t is \mathbf{x}_{it} , $\mathbf{x}_{it}, \boldsymbol{\epsilon}_t \in \mathbb{R}^r$, $\boldsymbol{\mu}_0, \boldsymbol{\alpha}_t \in \mathbb{R}^p$, $\mathbf{F} \in \mathbb{R}^{p \times p}$, $\mathbf{Q} \in \mathbb{R}^{r \times r}$ is a positive definite matrix, $\mathbf{Q}_0 \in \mathbb{R}^{p \times p}$ is a positive definite matrix, \mathbf{o}_t s are known offsets, and \mathbf{R} is a $p \times r$ matrix with $p \geq r$ which contains a subset of r columns of \mathbf{I}_p . We will order the entires of \mathbf{R} such that the first r columns are the first r columns of \mathbf{I}_p . I.e.,

$$\mathbf{R} = \begin{pmatrix} \mathbf{I}_r \\ \mathbf{0} \end{pmatrix}$$

Superscript $+$ denotes the Moore-Penrose inverse, $\mathbf{R}^+ = \mathbf{R}^\top$ 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 where the r first diagonal entries has value 1 and the rest of the diagonal entries are zero. The data sets we are working with have $n_{\max} \gg p \geq r$ (e.g. $n_{\max} = 1000$ and $r = 5$). We let

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

The above allows us to have an o th order vector autoregressions, VAR(o), by settings

$$\boldsymbol{\alpha}_t = (\boldsymbol{\xi}_t, \boldsymbol{\xi}_{t-1}, \dots, \boldsymbol{\xi}_{t-o+1})$$

$$\mathbf{F} = \begin{pmatrix} \mathbf{F}_1 & \cdots & \cdots & \mathbf{F}_{o-1} & \mathbf{F}_o \\ \mathbf{I}_r & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_r & \ddots & \vdots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{I}_r & \mathbf{0} \end{pmatrix}, \quad \mathbf{F}_i \in \mathbb{R}^{r \times r}$$

I will use a particle filter to get a discrete approximation of the conditional distribution 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} , $\boldsymbol{\omega}$, and $\boldsymbol{\mu}_0$. One choice of smoother is shown in Fearnhead et al. (2010) and another is the generalized two-filter smoother shown by Briers et al. (2009). The rest of vignette is structured as follows: first I give a brief introduction to the implemented particle filters and smoothers. Then I cover what the effect is of some the arguments to particle functions in \mathbf{R} in the packages. The implemented particle filter and smoother from Fearnhead et al. (2010) is presented next, followed by the used EM-algorithm and the smoother suggested by Briers et al. (2009). The last section covers the implemented approximations of the gradient and observed information matrix.

1.1 Overview

As a gentle introduction before the next sections, we will start by recalling an application of importance sampling, use this to motivate particle filtering, and give a brief idea of the implemented particle smoothers. Suppose we want to approximate a density $c(x) = \zeta \tilde{c}(x)$ where we only know $\tilde{c}(x)$ and not the normalization constant ζ . One way to approximate this density is to

- sample x_1, x_2, \dots, x_N from a distribution with density $b(x)$.
- Compute the unnormalized weights $\bar{w}_i = \tilde{c}(x_i)/b(x_i)$.
- Normalize the weights $w_i = \bar{w}_i / \sum_{i=1}^N \bar{w}_i$.

This gives us the following discrete approximation of the density

$$c(x) \approx \sum_{i=1}^N w_i \delta_{x_i}(x)$$

where δ_x is the Dirac delta function which has unit point mass at x . This is directly applicable to the model in Equation (1) as at time 1 we want to approximate

$$p(\boldsymbol{\alpha}_1 | \mathbf{y}_1) = \frac{g_1(\mathbf{y}_1 | \boldsymbol{\alpha}_1) \int f(\boldsymbol{\alpha}_1 | \mathbf{a}_0) \phi(\mathbf{a}_0 | \boldsymbol{\mu}_0, \mathbf{Q}_0) d\mathbf{a}_0}{p(\mathbf{y}_1)}$$

where $\phi(\cdot|\mathbf{m}, \mathbf{M})$ is the density function of a multivariate normal distribution with mean \mathbf{m} and covariance matrix \mathbf{M} . We can easily evaluate the numerator for each α_1 but not the normalization constant, $p(\mathbf{y}_1)$.

The extension to a particle filter (which I will call a forward particle filter) is that at time 2 we want to approximate

$$p(\alpha_{1:2}|\mathbf{y}_{1:2}) = p(\alpha_1|\mathbf{y}_1) \frac{g_2(\mathbf{y}_2|\alpha_2) f(\alpha_2|\alpha_1)}{p(\mathbf{y}_2|\mathbf{y}_1)}$$

Now, we can use the discrete approximation at time 1 of $p(\alpha_1|\mathbf{y}_1)$, sample α_2 given each sampled α_1 , and apply importance sampling again. We can repeat this with similar arguments at time 3, 4, ..., d giving us an approximation of $p(\alpha_{1:d}|\mathbf{y}_{1:d})$. We will call the last element of a sampled path at time t a particle. Further, we will denote the j th particle at time t and its associated weight by $\alpha_t^{(j)}$ and $w_t^{(j)}$ respectively.

One issue that may arise is that our samples (particles) may degenerate so essentially only one sampled path of $\alpha_{1:d}$ has any weight in the end. To avoid this, we may introduce a re-sampling step. One way to re-sample is using the weights and letting the re-sampling weights be $\beta_{t+1}^{(j)} = w_t^{(j)}$ where $\beta_{t+1}^{(j)}$ is the re-sampling weight of particle j at time t . We then sample with replacement using $\beta_{t+1}^{(j)}$. Another option when we re-sample the particles from time t is to use the information of the outcomes at time $t+1$, \mathbf{y}_{t+1} . This is called an auxiliary particle filter and is introduced by Pitt and Shephard (1999).

However, we may end up with few or only one unique value at the early time points (say α_1) when we re-sample. Thus, it will be useful to use a smoother to get a better approximation of the marginal density $p(\alpha_t|\mathbf{y}_{1:d})$. To do so, one idea is to use the two-filter formula from Kitagawa (1994). Though, this requires that we can evaluate $p(\mathbf{y}_{t:d}|\alpha_t)$. It turns out that we can approximate this up to a constant which is just what need. This is covered in further details in Section 6.

The approximation uses a particle filter which is run backwards in time and which approximates an artificial distribution. The arguments for the backward particle filter is very similar to the forward particle filter presented above. The k th particle in the backward particle filter at time t , its re-sampling weight, and the associated weight will be denoted by respectively $\tilde{\alpha}_t^{(k)}$, $\tilde{\beta}_{t-1}^{(k)}$ and $\tilde{w}_t^{(k)}$. The final i th smoothed particle and weight at time t will be denoted by $\hat{\alpha}_t^{(i)}$ and $\hat{w}_t^{(i)}$. The latter gives us the following approximation of the marginal density of $\alpha_t | \mathbf{y}_{1:d}$

$$p(\alpha_t|\mathbf{y}_{1:d}) \approx \sum_{i=1}^{N_S} \hat{w}_t^{(i)} \delta_{\hat{\alpha}_t^{(i)}}(\alpha_t)$$

if we sampled N_S smoothed particles at time t . The smoothing algorithm from Fearnhead et al. (2010) is shown in Algorithm 1, the forward particle filter is shown Algorithm 2, and the backward particle filter is shown in Algorithm 3.

1.2 Methods in the Package

The PF_EM method in the `dynamichazard` package contains an implementation of the described methods. 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 particles in the smoothing step, $N_s > N$, as pointed out in the discussion in Fearnhead

et al. (2010, p. 460-461). Further, selecting $N_f > N$ may be preferable to ensure coverage of the state space at time 0 and $d + 1$.

We do not need to sample the time 0 and $d + 1$ particles. Instead we can make a special proposal distribution for time 1 and time d . This is not implemented though...

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

- `"bootstrap_filter"` for a bootstrap filter. This is where we sample using Equation (5), (11) and (14). This is fast but the proposal distribution may be a poor approximation of the distribution we want to target.
- `"PF_normal_approx_w_cloud_mean"` and `"AUX_normal_approx_w_cloud_mean"` for the Taylor approximation of the conditional density of \mathbf{y}_t made using the mean of the parent particles and/or mean of the child particles. See Section 2. 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 approximation of the conditional density of \mathbf{y}_t made using the parent and/or child particle. See Section 2. 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 Re-sampling* (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.

It is not clear what will give the best performance for a given data set at a fixed computation cost. An advice is to use the `trace` argument and check the effective sample at each point in time during the estimation. `"bootstrap_filter"` may not be that much cheaper in terms of computation time as we still have to evaluate g_t in Equation (16), (17), and (19) which is $\mathcal{O}(n_{\max}Nr)$ or $\mathcal{O}(n_{\max}N_Sr)$ which is typically computationally expensive as n_{\max} is large. On the other hand, the `"..._w_particles"` methods have a computational complexity of $\mathcal{O}(n_{\max}Nr^2)$ or $\mathcal{O}(n_{\max}N_Sr^2)$ with a potentially much larger constant. Thus, the `"..._w_cloud_mean"` may be preferred.

The rest of the vignette covers the implemented methods. It is mainly included to show exactly what is computed and why. Further, I cover some currently not implemented extensions that may be implemented in the future.

1.3 Proposal Distributions and Re-sampling Weights

Algorithm 1 shows one of the particle smoothers shown by Fearnhead et al. (2010) in the first order state space model. In this situation $\mathbf{R} = \mathbf{I}_r$, $r = p$ and $\boldsymbol{\alpha}_t = \boldsymbol{\xi}_t$. We need to specify a series of proposal distributions and re-sampling weights. To show what is implemented and why, we first consider the model where

$$\mathbf{y}_t \mid \boldsymbol{\alpha}_t \sim \mathcal{N}(\mathbf{X}_t\boldsymbol{\alpha}_t + \mathbf{o}_t + \mathbf{Z}_t\boldsymbol{\omega}, \mathbf{H}_t)$$

for some known positive definite matrix \mathbf{H}_t . This is not implemented in this package but deriving optimal re-sampling weights and proposal distributions is possible in this case. In fact, it makes little sense to use a particle filter and particle smoother in this case since the

Kalman filter and an exact smoother can be applied. However, the results here will turn out to be useful to motivate the approximations we use later. The state space model is

$$\begin{aligned} \mathbf{y}_t &\sim \mathcal{N}(\boldsymbol{\eta}_t, \mathbf{H}_t) \\ \boldsymbol{\eta}_t &= \mathbf{X}_t \boldsymbol{\alpha}_t + \mathbf{o}_t + \mathbf{Z}_t \boldsymbol{\omega} & i = 1, \dots, n_t \\ \boldsymbol{\alpha}_t &= \mathbf{F} \boldsymbol{\alpha}_{t-1} + \mathbf{R} \boldsymbol{\epsilon}_t & \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}), \quad t = 1, \dots, d \\ & & \boldsymbol{\alpha}_0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \mathbf{Q}_0) \end{aligned}$$

We let $\mathbf{h}_t = \mathbf{o}_t + \mathbf{Z}_t \boldsymbol{\omega}$ such that $\boldsymbol{\eta}_t = \mathbf{X}_t \boldsymbol{\alpha}_t + \mathbf{h}_t$ to ease the notation. We first turn to the forward particle filter in Algorithm 2. Ideally, we want the re-sampling weights to be

$$\begin{aligned} \beta_t^{(j)} &\propto p(\mathbf{y}_t | \boldsymbol{\alpha}_{t-1}^{(j)}) w_{t-1}^{(j)} & (2) \\ &= \int g_t(\mathbf{y}_t | \boldsymbol{\alpha}_t) f(\mathbf{a}_t | \boldsymbol{\alpha}_{t-1}^{(j)}) d\mathbf{a}_t w_{t-1}^{(j)} \\ &= \phi(\mathbf{y}_t | \mathbf{X}_t \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)} + \mathbf{h}_t, \mathbf{X}_t \mathbf{Q} \mathbf{X}_t^\top + \mathbf{H}_t) w_{t-1}^{(j)} \end{aligned}$$

We can notice that setting $\beta_t^{(j)} = w_{t-1}^{(j)}$ yields the so-called sequential importance re-sampling algorithm. For the proposal distribution, the optimal proposal density is

$$q(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t) = p(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t)$$

where we find that

$$\begin{aligned} \log p(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t) &= \log p(\boldsymbol{\alpha}_t, \mathbf{y}_t | \boldsymbol{\alpha}_{t-1}^{(j)}) + \dots \\ &= \log g_t(\mathbf{y}_t | \boldsymbol{\alpha}_t) + \log f(\mathbf{a}_t | \boldsymbol{\alpha}_{t-1}^{(j)}) + \dots \\ &= -\frac{1}{2} (\mathbf{y}_t - \mathbf{X}_t \boldsymbol{\alpha}_t - \mathbf{h}_t)^\top \mathbf{H}_t^{-1} (\mathbf{y}_t - \mathbf{X}_t \boldsymbol{\alpha}_t - \mathbf{h}_t) \\ &\quad - \frac{1}{2} (\boldsymbol{\alpha}_t - \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)})^\top \mathbf{Q}^{-1} (\boldsymbol{\alpha}_t - \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)}) + \dots \\ &= -\frac{1}{2} \boldsymbol{\alpha}_t^\top \boldsymbol{\Sigma}_t^{-1} \boldsymbol{\alpha}_t + \boldsymbol{\alpha}_t^\top \boldsymbol{\Sigma}_t^{-1} \boldsymbol{\mu}(\boldsymbol{\alpha}_{t-1}^{(j)}) + \dots \\ \boldsymbol{\Sigma}_t &= (\mathbf{Q}^{-1} + \mathbf{X}_t^\top \mathbf{H}_t^{-1} \mathbf{X}_t)^{-1} & (3) \\ \boldsymbol{\mu}(\mathbf{x}) &= \boldsymbol{\Sigma}_t (\mathbf{Q}^{-1} \mathbf{F} \mathbf{x} + \mathbf{X}_t^\top \mathbf{H}_t^{-1} (\mathbf{y}_t - \mathbf{h}_t)) & (4) \end{aligned}$$

The ... are terms of the normalization constant. We recognize the multivariate normal distribution density and thus the optimal proposal density is

$$q(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t) = \phi(\boldsymbol{\alpha}_t | \boldsymbol{\mu}(\boldsymbol{\alpha}_{t-1}^{(j)}), \boldsymbol{\Sigma}_t)$$

Alternatively, we can use the so-called bootstrap filter and let

$$q(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t) = \phi(\boldsymbol{\alpha}_t | \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{Q}) \quad (5)$$

which we can sample from in $\mathcal{O}(Np^2)$ time if we have a pre-computed Cholesky decomposition of \mathbf{Q} . This is computationally cheap compared to optimal solution which is $\mathcal{O}(Np^2 + p^3 + n_{\max} p^2)$ but it is not optimal.

Backward filter (Algorithm 3)

We need to specify the artificial prior $\gamma_t(\boldsymbol{\alpha}_t)$ for our artificial backward distribution. Briers et al. (2009, p. 69-70) provides recommendation on the selection. One suggestion is the artificial density function

$$\begin{aligned}\gamma_t(\boldsymbol{\alpha}_t) &= \phi\left(\boldsymbol{\alpha}_t \mid \overleftarrow{\mathbf{m}}_t, \overleftarrow{\mathbf{P}}_t\right) \\ \overleftarrow{\mathbf{m}}_t &= \mathbf{F}^t \boldsymbol{\mu}_0 \\ \overleftarrow{\mathbf{P}}_t &= \begin{cases} \mathbf{Q}_0 & t = 0 \\ \mathbf{F} \overleftarrow{\mathbf{P}}_{t-1} \mathbf{F}^\top + \mathbf{Q} & t > 0 \end{cases}\end{aligned}\quad (6)$$

The backward arrows are added to stress that these are means and covariance matrices which we use in the artificial distribution we target in the backward particle filter. The artificial distribution we target in backward particle filters has the following conditional density functions

$$\begin{aligned}\tilde{p}(\boldsymbol{\alpha}_{t:d} \mid \mathbf{y}_{t:d}) &\propto \gamma_t(\boldsymbol{\alpha}_t) \prod_{i=t}^d g_i(\mathbf{y}_i \mid \boldsymbol{\alpha}_i) \prod_{i=t}^{d-1} f(\boldsymbol{\alpha}_{i+1} \mid \boldsymbol{\alpha}_i) \\ \tilde{p}(\boldsymbol{\alpha}_t \mid \mathbf{y}_{(t+1):d}) &\propto \gamma_t(\boldsymbol{\alpha}_t) \int \tilde{p}(\mathbf{a}_{t+1} \mid \mathbf{y}_{(t+1):d}) \frac{f(\mathbf{a}_{t+1} \mid \boldsymbol{\alpha}_t)}{\gamma_{t+1}(\mathbf{a}_{t+1})} d\mathbf{a}_{t+1} \\ \tilde{p}(\boldsymbol{\alpha}_t \mid \mathbf{y}_{t:d}) &\propto g_t(\mathbf{y}_t \mid \boldsymbol{\alpha}_t) \tilde{p}(\boldsymbol{\alpha}_t \mid \mathbf{y}_{(t+1):d}) \\ \tilde{p}(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t+1}) &= \frac{f(\boldsymbol{\alpha}_{t+1} \mid \boldsymbol{\alpha}_t) \gamma_t(\boldsymbol{\alpha}_t)}{\gamma_{t+1}(\boldsymbol{\alpha}_{t+1})}\end{aligned}\quad (7)$$

where we have left out some of the normalization constants. Sampling from this artificial distribution turns out to be useful as it gives us an approximation of a conditional density we need up to a constant (see Section 6). To derive the re-sampling weight, we first find an expression for the density $\tilde{p}(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t+1})$. We can observe that

$$\begin{aligned}\log \tilde{p}(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t+1}) &= \log f(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t+1}) + \log \gamma_t(\boldsymbol{\alpha}_t) + \dots \\ &= -\frac{1}{2} \boldsymbol{\alpha}_t^\top \overleftarrow{\mathbf{S}}_t^{-1} \boldsymbol{\alpha}_t - \boldsymbol{\alpha}_t^\top \overleftarrow{\mathbf{S}}_t^{-1} \overleftarrow{\mathbf{a}}_t(\boldsymbol{\alpha}_{t+1}) + \dots \\ \overleftarrow{\mathbf{S}}_t &= (\mathbf{P}_t^{-1} + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{F})^{-1} \\ \overleftarrow{\mathbf{a}}_t(\mathbf{x}) &= \overleftarrow{\mathbf{S}}_t (\mathbf{P}_t^{-1} \mathbf{m}_t + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{x})\end{aligned}$$

so

$$\tilde{p}(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t+1}) = \phi\left(\boldsymbol{\alpha}_t \mid \overleftarrow{\mathbf{a}}_t(\boldsymbol{\alpha}_{t+1}), \overleftarrow{\mathbf{S}}_t\right)\quad (8)$$

As shown by Fearnhead et al. (2010), we can show that

$$\begin{aligned}\overleftarrow{\mathbf{S}}_t &= \overleftarrow{\mathbf{P}}_t \mathbf{F}^\top \overleftarrow{\mathbf{P}}_{t+1}^{-1} \mathbf{Q} \mathbf{F}^{-\top} \\ \overleftarrow{\mathbf{a}}_t(\mathbf{x}) &= \overleftarrow{\mathbf{P}}_t \mathbf{F}^\top \overleftarrow{\mathbf{P}}_{t+1}^{-1} \mathbf{x} + \overleftarrow{\mathbf{S}}_t \overleftarrow{\mathbf{P}}_t^{-1} \overleftarrow{\mathbf{m}}_t\end{aligned}\quad (9)$$

e.g., by

$$\begin{aligned}
& \left(\overleftarrow{\mathbf{P}}_t \mathbf{F}^\top \overleftarrow{\mathbf{P}}_{t+1}^{-1} \mathbf{Q} \mathbf{F}^{-\top} \right)^{-1} \left(\mathbf{P}_t^{-1} + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{F} \right)^{-1} \\
&= \mathbf{F}^\top \mathbf{Q}^{-1} \overleftarrow{\mathbf{P}}_{t+1} \mathbf{F}^{-\top} \overleftarrow{\mathbf{P}}_t^{-1} \left(\mathbf{P}_t^{-1} + \mathbf{F}^\top \mathbf{Q} \mathbf{F} \right)^{-1} \\
&= \mathbf{F}^\top \mathbf{Q}^{-1} \overleftarrow{\mathbf{P}}_{t+1} \mathbf{F}^{-\top} \overleftarrow{\mathbf{P}}_t^{-1} \left(\overleftarrow{\mathbf{P}}_t - \overleftarrow{\mathbf{P}}_t \mathbf{F}^\top \left(\mathbf{Q} + \mathbf{F} \overleftarrow{\mathbf{P}}_t \mathbf{F}^\top \right)^{-1} \mathbf{F} \overleftarrow{\mathbf{P}}_t \right) \\
&= \mathbf{F}^\top \mathbf{Q}^{-1} \overleftarrow{\mathbf{P}}_{t+1} \mathbf{F}^{-\top} \left(\mathbf{I} - \mathbf{F}^\top \overleftarrow{\mathbf{P}}_{t+1}^{-1} \mathbf{F} \overleftarrow{\mathbf{P}}_t \right) \\
&= \mathbf{F}^\top \mathbf{Q}^{-1} \overleftarrow{\mathbf{P}}_{t+1} \mathbf{F}^{-\top} - \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{F} \overleftarrow{\mathbf{P}}_t \\
&= \mathbf{F}^\top \mathbf{Q}^{-1} \left(\mathbf{F} \overleftarrow{\mathbf{P}}_t \mathbf{F}^\top + \mathbf{Q} \right) \mathbf{F}^{-\top} - \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{F} \overleftarrow{\mathbf{P}}_t \\
&= \mathbf{I}
\end{aligned}$$

where we assume that all matrices are non-singular and we use the Woodbury matrix identity. Similar arguments can be used for $\overleftarrow{\mathbf{a}}_t(\mathbf{x})$. Using the above, we can find that the optimal re-sampling weights are

$$\begin{aligned}
\tilde{\beta}_t^{(k)} &\propto \tilde{p} \left(\mathbf{y}_t \mid \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \right) \tilde{w}_{t+1}^{(k)} \\
&\propto \int g_t(\mathbf{y}_t \mid \mathbf{a}_t) \tilde{p} \left(\mathbf{a}_t \mid \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \right) d\mathbf{a}_t \tilde{w}_{t+1}^{(k)} \\
&= \phi \left(\mathbf{y}_t \mid \mathbf{X}_t \overleftarrow{\mathbf{a}}_t(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) + \mathbf{h}_t, \mathbf{X}_t \overleftarrow{\mathbf{S}}_t \mathbf{X}_t^\top + \mathbf{H}_t \right) \tilde{w}_{t+1}^{(k)}
\end{aligned} \tag{10}$$

or we can set the re-sampling weights proportional to $\tilde{\beta}_t^{(k)} \propto \tilde{w}_{t+1}^{(k)}$ and get a sequential importance re-sampling like algorithm. As for the proposal distribution, the optimal density is

$$\begin{aligned}
\log \tilde{q} \left(\boldsymbol{\alpha}_t \mid \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \right) &= \log \gamma_t(\boldsymbol{\alpha}_t) + \log g_t(\mathbf{y}_t \mid \boldsymbol{\alpha}_t) + \log f \left(\boldsymbol{\alpha}_{t+1}^{(k)} \mid \boldsymbol{\alpha}_t \right) + \dots \\
&= -\frac{1}{2} \boldsymbol{\alpha}_t^\top \overleftarrow{\boldsymbol{\Sigma}}_t^{-1} \boldsymbol{\alpha}_t + \boldsymbol{\alpha}_t^\top \overleftarrow{\boldsymbol{\Sigma}}_t^{-1} \overleftarrow{\boldsymbol{\mu}}_t(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) + \dots \\
\overleftarrow{\boldsymbol{\Sigma}}_t &= \left(\mathbf{P}_t^{-1} + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{F} + \mathbf{X}_t^\top \mathbf{H}_t^{-1} \mathbf{X}_t \right)^{-1} \\
\overleftarrow{\boldsymbol{\mu}}_t(\mathbf{x}) &= \boldsymbol{\Sigma}_t \left(\mathbf{P}_t^{-1} \mathbf{m}_t + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{x} + \mathbf{X}_t^\top \mathbf{H}_t^{-1} (\mathbf{y}_t - \mathbf{h}_t) \right)
\end{aligned}$$

Thus, we set

$$\tilde{q} \left(\boldsymbol{\alpha}_t \mid \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \right) = \phi \left(\boldsymbol{\alpha}_t \mid \overleftarrow{\boldsymbol{\mu}}_t(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}), \overleftarrow{\boldsymbol{\Sigma}}_t \right)$$

A computationally simpler but non-optimal option is to use a method like the bootstrap filter and set

$$\tilde{q} \left(\boldsymbol{\alpha}_t \mid \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \right) = \tilde{p} \left(\boldsymbol{\alpha}_t \mid \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \right) \tag{11}$$

Combining / smoothing (Algorithm 1)

We end this example with the conditional Gaussian observable outcome model with the proposal distribution needed for Algorithm 1. We want to select

$$\begin{aligned}
q \left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \right) &= p \left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \right) \\
&\propto g(\mathbf{y}_t \mid \boldsymbol{\alpha}_t) f \left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)} \right) f \left(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \mid \boldsymbol{\alpha}_t \right)
\end{aligned}$$

Looking at the log density as we did before, we find that

$$\begin{aligned}
\log q\left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right) &= \log g\left(\mathbf{y}_t \mid \boldsymbol{\alpha}_t\right) + \log f\left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)}\right) + \log f\left(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} \mid \boldsymbol{\alpha}_t\right) + \dots \\
&= -\frac{1}{2} \boldsymbol{\alpha}_t^\top \overleftrightarrow{\boldsymbol{\Sigma}}_t^{-1} \boldsymbol{\alpha}_t + \boldsymbol{\alpha}_t^\top \overleftrightarrow{\boldsymbol{\Sigma}}_t^{-1} \overleftrightarrow{\boldsymbol{\mu}}_t\left(\boldsymbol{\alpha}_{t-1}^{(j)}, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right) + \dots \\
\overleftrightarrow{\boldsymbol{\Sigma}}_t &= \left(\mathbf{Q}^{-1} + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{F} + \mathbf{X}_t^\top \mathbf{H}^{-1} \mathbf{X}_t\right)^{-1} \tag{12}
\end{aligned}$$

$$\overleftrightarrow{\boldsymbol{\mu}}_t(\mathbf{x}, \overleftarrow{\mathbf{x}}) = \boldsymbol{\Sigma}_t \left(\mathbf{Q}^{-1} \mathbf{F} \mathbf{x} + \mathbf{F}^\top \mathbf{Q}^{-1} \overleftarrow{\mathbf{x}} + \mathbf{X}_t^\top \mathbf{H}^{-1} (\mathbf{y}_t - \mathbf{h}_t)\right) \tag{13}$$

so that

$$q\left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right) = \phi\left(\boldsymbol{\alpha}_t \mid \overleftrightarrow{\boldsymbol{\mu}}_t\left(\boldsymbol{\alpha}_{t-1}^{(j)}, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right), \boldsymbol{\Sigma}_t\right)$$

Alternatively, we can use a method like the bootstrap filter with a proposal distribution with

$$\begin{aligned}
\overleftrightarrow{\boldsymbol{\Sigma}}_t &= \left(\mathbf{Q}^{-1} + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{F}\right)^{-1} \\
\overleftrightarrow{\boldsymbol{\mu}}_t(\mathbf{x}, \overleftarrow{\mathbf{x}}) &= \boldsymbol{\Sigma}_t \left(\mathbf{Q}^{-1} \mathbf{F} \mathbf{x} + \mathbf{F}^\top \mathbf{Q}^{-1} \overleftarrow{\mathbf{x}}\right) \tag{14}
\end{aligned}$$

This is not optimal but faster.

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

Input:

$\mathbf{Q}, \mathbf{Q}_0, \mathbf{a}_0, \mathbf{X}_1, \dots, \mathbf{X}_d, \mathbf{Z}_1, \dots, \mathbf{Z}_d, \mathbf{o}_1, \dots, \mathbf{o}_d, \mathbf{y}_1, \dots, \mathbf{y}_d, R_1, \dots, R_d, \boldsymbol{\omega}$
 Proposal distribution with density

$$q\left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}\right) \quad (15)$$

1: **procedure** FILTER FORWARD

2: Run a forward particle filter to get particle clouds $\left\{\boldsymbol{\alpha}_t^{(j)}, w_t^{(j)}, \beta_{t+1}^{(j)}\right\}_{j=1, \dots, N}$ approximating the density $p(\boldsymbol{\alpha}_t \mid \mathbf{y}_{1:t})$ 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 the artificial density $\tilde{p}(\boldsymbol{\alpha}_t \mid \mathbf{y}_{t:d})$ for $t = d+1, d, d-1, \dots, 1$. See Algorithm 3.

5: **procedure** SMOOTH (COMBINE)

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

Re-sample

7: Sample $i = 1, 2, \dots, N_s$ pairs of $(j_i, k_i) \in N^2$ 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 the proposal distribution $\tilde{q}\left(\cdot \mid \boldsymbol{\alpha}_{t-1}^{(j_i)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}\right)$.

Re-weight

9: Assign each particle a weight

$$\hat{w}_t^{(i)} \propto \frac{f\left(\hat{\boldsymbol{\alpha}}_t^{(i)} \mid \boldsymbol{\alpha}_{t-1}^{(j_i)}\right) g_t\left(\mathbf{y}_t \mid \hat{\boldsymbol{\alpha}}_t^{(i)}\right) f\left(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)} \mid \hat{\boldsymbol{\alpha}}_t^{(i)}\right) w_{t-1}^{(j_i)} \tilde{w}_{t+1}^{(k_i)}}{\tilde{q}\left(\hat{\boldsymbol{\alpha}}_t^{(i)} \mid \boldsymbol{\alpha}_{t-1}^{(j_i)}, \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}\right) \beta_t^{(j_i)} \tilde{\beta}_t^{(k_i)} \gamma_{t+1}\left(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}\right)} \quad (16)$$

Algorithm 2 Forward filter as in Pitt and Shephard (1999). It is equivalent with Doucet and Johansen (2009, p. 20 and 25). The version and notation below is from Fearnhead et al. (2010, p. 449).

Input:

Proposal distribution with density

$$q\left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t\right)$$

Function h to compute re-sampling weights

$$\beta_t^{(j)} \propto h(\mathbf{y}_t, \boldsymbol{\alpha}_{t-1}^{(j)}) w_{t-1}^{(j)}$$

- 1: Sample $\boldsymbol{\alpha}_0^{(1)}, \dots, \boldsymbol{\alpha}_0^{(N_f)}$ particles from $\mathcal{N}(\boldsymbol{\mu}_0, \mathbf{Q}_0)$ 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 $\beta_t^{(j)}$ using h and re-sample according to $\beta_t^{(j)}$ to get indices j_1, \dots, j_N . If we do not re-sample then set $\beta_t^{(j)} = 1/N$ or $1/N_f$ at time $t = 1$.
- 5: **procedure** PROPAGATE
- 6: Sample new particles $\boldsymbol{\alpha}_t^{(i)}$ using the proposal distribution $q\left(\boldsymbol{\alpha}_t \mid \boldsymbol{\alpha}_{t-1}^{(j_i)}, \mathbf{y}_t\right)$.
- 7: **procedure** RE-WEIGHT
- 8: Re-weight particles using

$$w_t^{(i)} \propto \frac{g_t\left(\mathbf{y}_t \mid \boldsymbol{\alpha}_t^{(i)}\right) f\left(\boldsymbol{\alpha}_t^{(i)} \mid \boldsymbol{\alpha}_{t-1}^{(j_i)}\right) w_{t-1}^{(j_i)}}{q\left(\boldsymbol{\alpha}_t^{(i)} \mid \boldsymbol{\alpha}_{t-1}^{(j_i)}, \mathbf{y}_t\right) \beta_t^{(j_i)}} \quad (17)$$

Algorithm 3 Backwards filter. See Briers et al. (2009) and Fearnhead et al. (2010).

Input:

An artificial distribution

$$\tilde{p}(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d}) \propto \gamma_t(\boldsymbol{\alpha}_t) p(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_t) \quad (18)$$

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

Proposal distribution

$$\tilde{q}(\boldsymbol{\alpha}_t | \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})$$

Function h to compute re-sampling weights

$$\tilde{\beta}_t^{(k)} \propto h(\mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) \tilde{w}_{t+1}^{(k)}$$

- 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 $\tilde{\beta}_t^{(k)}$ using h 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 $\tilde{\beta}_t^{(k)} = 1/N$ or $1/N_f$ at time $t = d$.
- 5: **procedure** PROPAGATE
- 6: Sample new particles $\tilde{\boldsymbol{\alpha}}_t^{(i)}$ using the proposal distribution $\tilde{q}(\boldsymbol{\alpha}_t | \tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}, \mathbf{y}_t)$.
- 7: **procedure** RE-WEIGHT
- 8: Re-weight particles using

$$\tilde{w}_t^{(i)} \propto \frac{g_t(\mathbf{y}_t | \tilde{\boldsymbol{\alpha}}_t^{(i)}) f(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)} | \tilde{\boldsymbol{\alpha}}_t^{(i)}) \gamma_t(\tilde{\boldsymbol{\alpha}}_t^{(i)}) \tilde{w}_{t+1}^{(k_i)}}{q(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}, \mathbf{y}_t) \gamma_{t+1}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)}) \beta_t^{(k_i)}} \quad (19)$$

2 Non-linear Conditional Observation Model

If we go back to the model in Equation (1) then $\mathbf{y}_t \mid \boldsymbol{\alpha}_t$ is not a multivariate normal distribution for the implemented models. In this case, we have no closed form solutions for the optimal re-sampling weights, and we do not know the following conditional distributions: $\boldsymbol{\alpha}_t \mid \mathbf{y}_t, \boldsymbol{\alpha}_{t-1}$, $\boldsymbol{\alpha}_t \mid \mathbf{y}_t, \boldsymbol{\alpha}_{t+1}$ (in the artificial distribution $\tilde{\mathbb{P}}$), and $\boldsymbol{\alpha}_t \mid \mathbf{y}_t, \boldsymbol{\alpha}_{t-1}, \boldsymbol{\alpha}_{t+1}$. However, assume that $g_t(\mathbf{y}_t \mid \boldsymbol{\alpha}_t)$ is log-concave in $\boldsymbol{\alpha}_t$. If this is true then it is easy to show that all of the previous three conditional distributions are unimodal. Hence, we can make a multivariate normal approximation as in Pitt and Shephard (1999). To do so, we make a second order Taylor expansion around some value \mathbf{z} to get

$$\begin{aligned} k_t(\boldsymbol{\alpha}_t) &= \log g_t(\mathbf{y}_t \mid \boldsymbol{\eta}(\boldsymbol{\alpha}_t)), & \boldsymbol{\eta}(\boldsymbol{\alpha}_t) &= \mathbf{X}_t \boldsymbol{\alpha}_t + \mathbf{h}_t \\ \log g_t(\mathbf{y}_t \mid \boldsymbol{\alpha}_t) &\approx Dk_t(\mathbf{z})(\boldsymbol{\alpha}_t - \mathbf{z}) + \frac{1}{2}(\boldsymbol{\alpha}_t - \mathbf{z})^\top Hk_t(\mathbf{z})(\boldsymbol{\alpha}_t - \mathbf{z}) + \dots \\ &= \boldsymbol{\alpha}_t^\top Dk_t(\mathbf{z})^\top - \frac{1}{2}(\boldsymbol{\alpha}_t - \mathbf{z})^\top (-Hk_t(\mathbf{z}))(\boldsymbol{\alpha}_t - \mathbf{z}) + \dots \\ &= \boldsymbol{\alpha}_t^\top (-Hk_t(\mathbf{z}))(\mathbf{z} - Hk_t(\mathbf{z})^{-1}Dk_t(\mathbf{z})^\top) - \frac{1}{2}\boldsymbol{\alpha}_t^\top (-Hk_t(\mathbf{z}))\boldsymbol{\alpha}_t + \dots \end{aligned}$$

where \dots includes the zero order term, Dk_t is the Jacobian, and Hk_t denotes the Hessian. I add a subscript to D and H to which variable the Jacobian or Hessian is with respect to. We still assume that we use a first order state space model such that $r = p$. We notice that

$$\begin{aligned} Hk_t(\mathbf{z}) &= D_{\boldsymbol{\alpha}_t} \boldsymbol{\eta}(\mathbf{z})^\top H_\eta \log g_t(\mathbf{y}_t \mid \boldsymbol{\eta}(\mathbf{z})) D_{\boldsymbol{\alpha}_t} \boldsymbol{\eta}(\mathbf{z}) \\ &= \mathbf{X}_t^\top (-\mathbf{G}_t(\mathbf{z})) \mathbf{X}_t, & \mathbf{G}_t(\mathbf{z}) &= -H_\eta \log g_t(\mathbf{y}_t \mid \boldsymbol{\eta}(\mathbf{z})) \end{aligned}$$

which follows from the chain rule and we use that $H_{\boldsymbol{\alpha}_t} \boldsymbol{\eta}(\mathbf{z}) = \mathbf{0}$. Thus,

$$\begin{aligned} \log g_t(\mathbf{y}_t \mid \boldsymbol{\alpha}_t) &\approx \boldsymbol{\alpha}_t^\top \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{u}_t(\mathbf{z}) - \frac{1}{2} \boldsymbol{\alpha}_t^\top \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{X}_t \boldsymbol{\alpha}_t \\ \mathbf{u}_t(\mathbf{z}) &= \mathbf{X}_t \mathbf{z} - \mathbf{X}_t Hk_t(\mathbf{z})^{-1} Dk_t(\mathbf{z})^\top \end{aligned}$$

This yields the following multivariate normal approximation

$$g_t(\mathbf{y}_t \mid \boldsymbol{\alpha}_t) \approx \phi(\mathbf{X}_t \boldsymbol{\alpha}_t \mid \mathbf{u}_t(\mathbf{z}), \mathbf{G}_t(\mathbf{z})^{-1})$$

The Taylor approximation is easily used in the proposal distributions. E.g., for given \mathbf{z} , we get the following mean and covariance matrix analogues to Equation (3) and (4) in the proposal distribution in the forward particle filter

$$\begin{aligned} \boldsymbol{\Sigma}_t(\mathbf{z}) &= (\mathbf{Q}^{-1} + \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{X}_t)^{-1} \\ \boldsymbol{\mu}_t(\mathbf{x}, \mathbf{z}) &= \boldsymbol{\Sigma}_t(\mathbf{z}) (\mathbf{Q}^{-1} \mathbf{F} \mathbf{x} + \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{u}_t(\mathbf{z})) \end{aligned}$$

As for the re-sampling weights, we can use

$$\begin{aligned}
\hat{\boldsymbol{\alpha}} &= \boldsymbol{\mu}_t(\boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{z}) \\
\beta_t^{(j)} &\propto p(\mathbf{y}_t | \boldsymbol{\alpha}_{t-1}^{(j)}) w_{t-1}^{(j)} \\
&= \frac{p(\mathbf{y}_t | \boldsymbol{\alpha}_{t-1}^{(j)})}{g_t(\mathbf{y}_t | \hat{\boldsymbol{\alpha}}) f(\hat{\boldsymbol{\alpha}} | \boldsymbol{\alpha}_{t-1}^{(j)})} g_t(\mathbf{y}_t | \hat{\boldsymbol{\alpha}}) f(\hat{\boldsymbol{\alpha}} | \boldsymbol{\alpha}_{t-1}^{(j)}) w_{t-1}^{(j)} \\
&= \frac{g_t(\mathbf{y}_t | \hat{\boldsymbol{\alpha}}) f(\hat{\boldsymbol{\alpha}} | \boldsymbol{\alpha}_{t-1}^{(j)}) w_{t-1}^{(j)}}{p(\hat{\boldsymbol{\alpha}} | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t)} \\
&\approx \frac{g_t(\mathbf{y}_t | \hat{\boldsymbol{\alpha}}) f(\hat{\boldsymbol{\alpha}} | \boldsymbol{\alpha}_{t-1}^{(j)}) w_{t-1}^{(j)}}{q(\hat{\boldsymbol{\alpha}} | \boldsymbol{\alpha}_{t-1}^{(j)}, \mathbf{y}_t)}
\end{aligned}$$

as in Fearnhead et al. (2010). We can approximate the backwards particle filter re-sampling weights in equation (10) in a similar way

$$\begin{aligned}
\tilde{\beta}_t^{(k)} &\propto \tilde{p}(\mathbf{y}_t | \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) \tilde{w}_{t+1}^{(k)} \\
&\approx \frac{g_t(\mathbf{y}_t | \hat{\boldsymbol{\alpha}}) \tilde{p}(\hat{\boldsymbol{\alpha}} | \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) \tilde{w}_{t+1}^{(k)}}{\tilde{q}(\hat{\boldsymbol{\alpha}} | \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})} \\
&= \frac{g_t(\mathbf{y}_t | \hat{\boldsymbol{\alpha}}) f(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} | \hat{\boldsymbol{\alpha}}) \gamma_t(\hat{\boldsymbol{\alpha}}) \tilde{w}_{t+1}^{(k)}}{\tilde{q}(\hat{\boldsymbol{\alpha}} | \mathbf{y}_t, \tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}) \gamma_{t+1}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})} \tag{20}
\end{aligned}$$

$$\begin{aligned}
\hat{\boldsymbol{\alpha}} &= \hat{\boldsymbol{\mu}}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}, \mathbf{z}) \\
\hat{\boldsymbol{\mu}}(\mathbf{x}, \mathbf{z}) &= \hat{\boldsymbol{\Sigma}}_t(\mathbf{z}) (\mathbf{P}_t^{-1} \mathbf{m}_t + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{x} + \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{u}_t(\mathbf{z})) \tag{21} \\
\hat{\boldsymbol{\Sigma}}_t(\mathbf{z}) &= (\mathbf{P}_t^{-1} + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{F} + \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{X}_t)^{-1} \tag{22}
\end{aligned}$$

We may also use a multivariate t -distribution for the proposal distribution to get heavier tails than we do with the multivariate normal distribution. This is important as too light tailed proposal distributions (relative to the target) can yield few large importance weights.

2.1 Where to Make the Expansion

An options is to make the Taylor expansion at a mode for each particle or particle pair in the smoothing step. This yields

$$\begin{aligned}
\mathbf{z}^{(j)} &= \arg \max_{\mathbf{z}} g_t(\mathbf{y}_t | \mathbf{z}) f(\mathbf{z} | \boldsymbol{\alpha}_{t-1}^{(j)}) \\
\mathbf{z}^{(k)} &= \arg \max_{\mathbf{z}} g_t(\mathbf{y}_t | \mathbf{z}) \gamma_t(\mathbf{z}) f(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} | \mathbf{z}) \\
\mathbf{z}^{(i)} &= \arg \max_{\mathbf{z}} f(\mathbf{z} | \boldsymbol{\alpha}_{t-1}^{(j_i)}) g_t(\mathbf{y}_t | \mathbf{z}) f(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k_i)} | \mathbf{z})
\end{aligned}$$

for respectively the forward particle filter, backward particle filter, and smoother. The downside is a $\mathcal{O}(r^2 n_{\max} N_S)$ or $\mathcal{O}(r^2 n_{\max} N)$ computational complexity at each time point as we have to evaluate $\mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{X}_t$ for each particle or particle pair. Instead we can make the approximation once at each time point at respectively $\sum_{i=1}^N w_{t-1}^{(j)} \boldsymbol{\alpha}_{t-1}^{(j)}$, $\sum_{i=1}^N \tilde{w}_{t+1}^{(j)} \tilde{\boldsymbol{\alpha}}_{t+1}^{(j)}$, and

$$(\mathbf{Q}^{-1} + \mathbf{F}^\top \mathbf{Q}^{-1} \mathbf{F})^{-1} \left(\mathbf{Q}^{-1} \mathbf{F} \sum_{i=1}^N w_{t-1}^{(j)} \boldsymbol{\alpha}_{t-1}^{(j)} + \mathbf{F}^\top \mathbf{Q}^{-1} \sum_{i=1}^N \tilde{w}_{t+1}^{(j)} \tilde{\boldsymbol{\alpha}}_{t+1}^{(j)} \right)$$

which will reduce the computational complexity at each time point to $\mathcal{O}(r n_{\max} N_S + r p n_{\max})$ or $\mathcal{O}(r n_{\max} N + r p n_{\max})$.

3 Higher Order State Space Models

Now, we will consider the case where $p > r$. Currently this is not supported in the package but may be in the future. An example is a second order vector autoregression, VAR(2), with $2r = p$ and

$$\mathbf{F} = \begin{pmatrix} \mathbf{F}_1 & \mathbf{F}_2 \\ \mathbf{I}_r & \mathbf{0} \end{pmatrix}, \quad \mathbf{F}_i \in \mathbb{R}^{r \times r}$$

$$\boldsymbol{\alpha}_t = (\boldsymbol{\xi}_t^\top \quad \boldsymbol{\xi}_{t-1}^\top)^\top$$

Here

$$f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}) = \delta_{\mathbf{K}\mathbf{F}\boldsymbol{\alpha}_{t-1}}(\mathbf{K}\boldsymbol{\alpha}_t) \phi(\mathbf{R}^+ \boldsymbol{\alpha}_t | \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}, \mathbf{Q}), \quad \mathbf{K} = \begin{pmatrix} \mathbf{0} \\ \mathbf{I}_{p-r} \end{pmatrix}$$

$$= \delta_{\boldsymbol{\xi}_{t-1}}(\mathbf{K}\boldsymbol{\alpha}_t) \phi(\boldsymbol{\xi}_t | \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}, \mathbf{Q})$$

the second equality follows in a second order vector autoregression. This is easily implemented in the forward particle filter by sampling $\boldsymbol{\xi}_t | \boldsymbol{\alpha}_{t-1}$ and setting the remaining $p - r$ variables of $\boldsymbol{\alpha}_t$ to $\mathbf{K}\mathbf{F}\boldsymbol{\alpha}_{t-1}$ (the last $p - r$ rows of $\mathbf{F}\boldsymbol{\alpha}_{t-1}$). It is not as easy for the backward filter. To see this, consider the artificial transition density in Equation (7)

$$\tilde{p}(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t+1}) = \frac{f(\boldsymbol{\alpha}_{t+1} | \boldsymbol{\alpha}_t) \gamma_t(\boldsymbol{\alpha}_t)}{\gamma_{t+1}(\boldsymbol{\alpha}_{t+1})}$$

$$= \frac{\delta_{\mathbf{K}\mathbf{F}\boldsymbol{\alpha}_t}(\mathbf{K}\boldsymbol{\alpha}_{t+1}) \phi(\mathbf{R}^+ \boldsymbol{\alpha}_{t+1} | \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_t, \mathbf{Q}) \gamma_t(\boldsymbol{\alpha}_t)}{\gamma_{t+1}(\boldsymbol{\alpha}_{t+1})}$$

where the last equality holds in this example. This distribution is obviously degenerate as it only has mass at the point $\boldsymbol{\xi}_t = \mathbf{K}\mathbf{F}\boldsymbol{\alpha}_t$ in the VAR(2) model.

3.1 Backward Particle Filter

What we can do in the backward particle filter in a VAR(o) model with $o > 1$ is to sample $\boldsymbol{\xi}_t$ at time t given $\tilde{\boldsymbol{\alpha}}_{t+o}^{(k)}$. Thus, we start Algorithm 3 by sampling $\boldsymbol{\alpha}_{d+o}$. Further, we change

the artificial prior distribution density in Equation (6) to

$$\begin{aligned}\gamma_t(\boldsymbol{\alpha}_t) &= \phi\left(\boldsymbol{\alpha}_t \mid \overleftarrow{\mathbf{m}}_t, \overleftarrow{\mathbf{P}}_t\right) \\ \gamma_t(\boldsymbol{\xi}_t) &= \phi\left(\boldsymbol{\alpha}_t \mid \mathbf{R}^+ \overleftarrow{\mathbf{m}}_t, \mathbf{R}^+ \overleftarrow{\mathbf{P}}_t \mathbf{R}^{+\top}\right) \\ \overleftarrow{\mathbf{m}}_t &= \mathbf{F}^t \boldsymbol{\mu}_0 \\ \overleftarrow{\mathbf{P}}_t &= \begin{cases} \mathbf{Q}_0 & t = 0 \\ \mathbf{F} \overleftarrow{\mathbf{P}}_{t-1} \mathbf{F}^\top + \mathbf{R} \mathbf{Q} \mathbf{R}^\top & t > 0 \end{cases}\end{aligned}$$

Next, we find the conditional distribution $\tilde{\boldsymbol{\alpha}}_{t+o}^{(k)} \mid \boldsymbol{\xi}_t$. To do so, we first find the joint distribution of $(\boldsymbol{\alpha}_{t+o}^\top, \boldsymbol{\xi}_t^\top)^\top$. We can find that

$$\begin{aligned}\boldsymbol{\xi}_k &= \mathbf{R}^+ \mathbf{F}^k \boldsymbol{\alpha}_0 + \sum_{i=1}^k \mathbf{G}(k-i) \boldsymbol{\epsilon}_i \\ \mathbf{G}(j) &= \begin{cases} \mathbf{I}_r & j = 0 \\ \sum_{i=1}^{\min(j,o)} \mathbf{F}_i \mathbf{G}(j-i) & j > 0 \end{cases}\end{aligned}$$

Thus,

$$\begin{aligned}\mathbf{E}(\boldsymbol{\xi}_t) &= \mathbf{R}^+ \mathbf{F}^t \boldsymbol{\mu}_0 = \mathbf{R}^+ \mathbf{m}_t \\ \mathbf{E}(\boldsymbol{\alpha}_t) &= \mathbf{m}_t \\ \text{Var}(\boldsymbol{\xi}_t) &= \sum_{i=1}^t \mathbf{G}(t-i) \mathbf{Q} \mathbf{G}(t-i)^\top + \mathbf{R}^+ \mathbf{F}^t \mathbf{Q}_0 \mathbf{F}^{t\top} \mathbf{R}^{+\top} \\ &= \mathbf{R}^+ \mathbf{P}_t \mathbf{R}^{+\top} \\ \text{Cov}(\boldsymbol{\xi}_k, \boldsymbol{\xi}_l) &= \sum_{i=1}^l \mathbf{G}(k-i) \mathbf{Q} \mathbf{G}(l-i)^\top + \mathbf{R}^+ \mathbf{F}^k \mathbf{Q}_0 \mathbf{F}^{l\top} \mathbf{R}^{+\top}, \quad k > l \\ \text{Var}(\boldsymbol{\alpha}_t) &= \mathbf{P}_t\end{aligned}$$

with which we can find that the joint distribution is

$$\begin{pmatrix} \boldsymbol{\alpha}_{t+o} \\ \boldsymbol{\xi}_t \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} \mathbf{m}_{t+o} \\ \mathbf{R}^+ \mathbf{m}_t \end{pmatrix}, \begin{pmatrix} \mathbf{P}_{t+o} & \text{Cov}(\boldsymbol{\alpha}_{t+o}, \boldsymbol{\xi}_t) \\ \text{Cov}(\boldsymbol{\xi}_t, \boldsymbol{\alpha}_{t+o}) & \mathbf{R}^+ \mathbf{P}_t \mathbf{R}^{+\top} \end{pmatrix}\right)$$

We are now able to compute

$$\begin{aligned}\mathbf{E}(\boldsymbol{\alpha}_{t+o} \mid \boldsymbol{\xi}_t) &= \boldsymbol{\mu}_{\boldsymbol{\alpha}_{t+o} \mid \boldsymbol{\xi}_t} = \mathbf{m}_{t+o} + \text{Cov}(\boldsymbol{\alpha}_{t+o}, \boldsymbol{\xi}_t) \mathbf{R}^\top \mathbf{P}_t^{-1} \mathbf{R} (\boldsymbol{\xi}_t - \mathbf{R}^+ \mathbf{m}_t) \\ &= \mathbf{v}_t + \mathbf{V}_t \boldsymbol{\xi}_t \\ \mathbf{V}_t &= \text{Cov}(\boldsymbol{\alpha}_{t+o}, \boldsymbol{\xi}_t) \mathbf{R}^\top \mathbf{P}_t^{-1} \mathbf{R} \\ \mathbf{v}_t &= \mathbf{m}_{t+o} - \mathbf{V}_t \mathbf{R}^+ \mathbf{m}_t\end{aligned}$$

$$\text{Var}(\boldsymbol{\alpha}_{t+o} \mid \boldsymbol{\xi}_t) = \boldsymbol{\Sigma}_{\boldsymbol{\alpha}_{t+o} \mid \boldsymbol{\xi}_t} = \mathbf{P}_{t+o} - \text{Cov}(\boldsymbol{\alpha}_{t+o}, \boldsymbol{\xi}_t) \mathbf{R}^\top \mathbf{P}_t^{-1} \mathbf{R} \text{Cov}(\boldsymbol{\xi}_t, \boldsymbol{\alpha}_{t+o})$$

Having found this conditional distribution then we can apply similar arguments as we did before and find the following mean and covariance matrix in the proposal distribution

$$\begin{aligned}\overleftarrow{\boldsymbol{\mu}}(\tilde{\boldsymbol{\alpha}}_{t+o}^{(k)}, \mathbf{z}) &= \overleftarrow{\boldsymbol{\Sigma}}_t(\mathbf{z}) \left(\mathbf{R}^\top \mathbf{P}_t^{-1} \mathbf{m}_t + \mathbf{V}_t^\top \boldsymbol{\Sigma}_{\boldsymbol{\alpha}_{t+o} \mid \boldsymbol{\xi}_t}^{-1} (\tilde{\boldsymbol{\alpha}}_{t+o}^{(k)} - \mathbf{v}_t) + \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{u}_t(\mathbf{z}) \right) \\ \overleftarrow{\boldsymbol{\Sigma}}_t(\mathbf{z}) &= \left(\mathbf{R}^\top \mathbf{P}_t^{-1} \mathbf{R} + \mathbf{V}_t^\top \boldsymbol{\Sigma}_{\boldsymbol{\alpha}_{t+o} \mid \boldsymbol{\xi}_t}^{-1} \mathbf{V}_t + \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{X}_t \right)^{-1}\end{aligned}$$

which replaces the mean and covariance matrix in Equation (21) and (22), $\mathbf{z}_t \in \mathbb{R}^r$ is the value of $\boldsymbol{\xi}_t$ at which we make the Taylor expansion, and \mathbf{G}_t and \mathbf{u}_t are defined in terms of $\boldsymbol{\xi}_t$.

3.2 Combining/Smoothing

Similarly, in Algorithm 1 we sample pairs of particles $(\boldsymbol{\alpha}_{t-1}^{(j_i)}, \tilde{\boldsymbol{\alpha}}_{t+o}^{(k_i)})$, and sample $\boldsymbol{\xi}_t$ given each pair. We can replace the mean and covariance matrix in the proposal distribution in Equation (13) and (12) with the approximation

$$\begin{aligned} \overleftrightarrow{\boldsymbol{\Sigma}}_t(\mathbf{z}) &= \left(\mathbf{Q}^{-1} + \mathbf{V}_t^\top \boldsymbol{\Sigma}_{\alpha_{t+o}|\boldsymbol{\xi}_t}^{-1} \mathbf{V}_t + \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{X}_t \right)^{-1} \\ \overleftrightarrow{\boldsymbol{\mu}}_t(\boldsymbol{\alpha}_{t-1}^{(j_i)}, \tilde{\boldsymbol{\alpha}}_{t+o}^{(k_i)}, \mathbf{z}) &= \\ &\boldsymbol{\Sigma}_t(\mathbf{z}) \left(\mathbf{Q}^{-1} \mathbf{R}^+ \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j_i)} + \mathbf{V}_t^\top \boldsymbol{\Sigma}_{\alpha_{t+o}|\boldsymbol{\xi}_t}^{-1} (\tilde{\boldsymbol{\alpha}}_{t+o}^{(k_i)} - \mathbf{v}_t) + \mathbf{X}_t^\top \mathbf{G}_t(\mathbf{z}) \mathbf{X}_t \right) \end{aligned}$$

where \mathbf{z} is the value of $\boldsymbol{\xi}_t$ that we make the Taylor expansion at.

4 Log-Likelihood Evaluation and Parameter Estimation

In this section, I show an example of parameter estimation in the first order random walk using a Monte Carlo EM-algorithm. Then I cover the general vector autoregression model and how one can estimate the fixed effects. See Kantas et al. (2015); Del Moral et al. (2010); Schön et al. (2011) for a general coverage of parameter estimation with particle filters. Firstly though, I will remark that we can approximate the log-likelihood for a particular value of $\boldsymbol{\theta} = \{\mathbf{Q}, \mathbf{Q}_0, \boldsymbol{\mu}_0, \mathbf{F}\}$ as described in Doucet and Johansen (2009, p. 5) and Malik and Pitt (2011, p. 193) using the forward particle filter shown in Algorithm 2. Details are omitted here for the sake of brevity.

The formulas for parameter estimation for the first order random are particularly simple. We need to estimate \mathbf{Q} and \mathbf{a}_0 elements of $\boldsymbol{\varphi} = \{\mathbf{Q}, \mathbf{Q}_0, \boldsymbol{\mu}_0\}$. We do this by running Algorithm 1 for the current $\boldsymbol{\varphi}$. This yields the following quantities from the E-step

$$\begin{aligned} \mathbf{t}_t^{(\boldsymbol{\varphi})} &\approx \sum_{i=1}^{N_s} \hat{\boldsymbol{\alpha}}_t^{(i)} \hat{w}_t^{(i)} \\ \mathbf{T}_t^{(\boldsymbol{\varphi})} &\approx \sum_{i=1}^{N_s} \left(\hat{\boldsymbol{\alpha}}_t^{(i)} - \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j_{it})} \right) \left(\hat{\boldsymbol{\alpha}}_t^{(i)} - \mathbf{F} \boldsymbol{\alpha}_{t-1}^{(j_{it})} \right)^\top \hat{w}_t^{(i)} \end{aligned} \tag{23}$$

where we have extended the notation in Algorithm 1 such that superscript j_{it} is the index from forward cloud at time $t-1$ matching with i th smoothed particle at time t . Then we carry out the M-step by updating $\boldsymbol{\mu}_0$ and \mathbf{Q} given the summary statistics above

$$\boldsymbol{\mu}_0 = \mathbf{t}_1^{(\boldsymbol{\varphi})} \quad \mathbf{Q} = \frac{1}{d-1} \sum_{t=2}^d \mathbf{R}^+ \mathbf{T}_t^{(\boldsymbol{\varphi})} \mathbf{R}^{+\top} \tag{24}$$

We then take another iteration of the EM-algorithm with the new $\boldsymbol{\mu}_0$ and \mathbf{Q} and repeat till a convergence criteria is satisfied.

4.1 Vector Autoregression Models

We start by defining the following matrices to cover estimation in general vector autoregression models for the latent space variable

$$\begin{aligned}\mathbf{N} &= \left(\hat{\boldsymbol{\alpha}}_2^{(1)}, \dots, \hat{\boldsymbol{\alpha}}_2^{(N_s)}, \hat{\boldsymbol{\alpha}}_3^{(1)}, \dots, \hat{\boldsymbol{\alpha}}_3^{(N_s)}, \hat{\boldsymbol{\alpha}}_4^{(1)}, \dots, \hat{\boldsymbol{\alpha}}_d^{(N_s)} \right)^\top \mathbf{R}^{+\top} \\ \mathbf{M} &= \left(\boldsymbol{\alpha}_1^{(j_{12})}, \dots, \boldsymbol{\alpha}_1^{(j_{N_s 2})}, \boldsymbol{\alpha}_2^{(j_{13})}, \dots, \boldsymbol{\alpha}_2^{(j_{N_s 3})}, \boldsymbol{\alpha}_3^{(j_{14})}, \dots, \boldsymbol{\alpha}_{d-1}^{(j_{N_s d})} \right)^\top \\ \mathbf{W} &= \text{diag} \left(\hat{w}_2^{(1)}, \dots, \hat{w}_2^{(N_s)}, \hat{w}_3^{(1)}, \dots, \hat{w}_3^{(N_s)}, \hat{w}_4^{(1)}, \dots, \hat{w}_d^{(N_s)} \right)\end{aligned}$$

where $\text{diag}(\cdot)$ is a diagonal matrix. We suppress the dependence above on the result of the E-step in a given iteration of the EM-algorithm to ease the notation. The goal is to estimate \mathbf{F} and \mathbf{Q} in Equation (1). We can find that the M-step maximizers are

$$\hat{\mathbf{F}}^\top \mathbf{R}^{+\top} = (\mathbf{M}^\top \mathbf{W} \mathbf{M})^{-1} \mathbf{M}^\top \mathbf{W} \mathbf{N} \quad (25)$$

$$\hat{\mathbf{Q}} = \frac{1}{d-1} \left(\mathbf{N} - \mathbf{R}^+ \hat{\mathbf{F}} \mathbf{M} \right)^\top \mathbf{W} \left(\mathbf{N} - \mathbf{R}^+ \hat{\mathbf{F}} \mathbf{M} \right) \quad (26)$$

which are the typical vector autoregression estimators with weights. Equation (25) and (26) can easily be computed in parallel using QR decompositions as in the `bam` function in the `mgcv` package with a low memory footprint (see Wood et al., 2014). This is currently implemented. Though, the gains from a parallel implementation may be small as the computational complexity is independent of the number of observations. Consequently, the computation involved here is often fast relative to other part of the Monte Carlo EM-algorithm as the dimension of the state vector is small.

4.2 Restricted Vector Autoregression Models

Suppose that we want to restrict some of the parameters of \mathbf{F} and \mathbf{Q} . Let

$$\begin{aligned}(s_1, s_2, \dots, s_r)^\top &= \mathbf{J} \boldsymbol{\psi} \\ (o_{21}, o_{31}, \dots, o_{r1}, o_{32}, \dots, o_{r2}, o_{43}, \dots, o_{r, r-1})^\top &= \mathbf{K} \boldsymbol{\phi}\end{aligned}$$

Then we can restrict the model such that

$$\begin{aligned}\sigma_i &= \exp(s_i) & \rho_{ij} &= \frac{2}{1 + \exp(-o_{ij})} - 1 \\ \text{vec}(\mathbf{R}^+ \mathbf{F}) &= \mathbf{G} \boldsymbol{\theta} & \mathbf{Q} &= \mathbf{V} \mathbf{C} \mathbf{V} \\ \mathbf{V} &= \begin{pmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \ddots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & \sigma_r \end{pmatrix} & \mathbf{C} &= \begin{pmatrix} 1 & \rho_{21} & \cdots & \rho_{r1} \\ \rho_{21} & 1 & \ddots & \rho_{r2} \\ \vdots & \ddots & \ddots & \vdots \\ \rho_{r1} & \cdots & \rho_{r, r-1} & 1 \end{pmatrix}\end{aligned}$$

and where $\text{vec}(\cdot)$ is the vectorization function which stacks the columns of a matrix from left to right. E.g.,

$$\text{vec}(\mathbf{A}) = (a_{11}, a_{21}, a_{31}, a_{12}, a_{22}, a_{32}, a_{13}, a_{23}, a_{33})^\top$$

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

$\mathbf{G} \in \mathbb{R}^{rp \times g}$ is a known matrix with $g \leq rp$ and we assume that it has full column rank. Similarly, $\mathbf{J} \in \mathbb{R}^{r \times l}$ with $l \leq r$ and $\mathbf{K} \in \mathbb{R}^{r(r-1)/2 \times k}$ with $k \leq r(r-1)/2$. Both are known and have full column rank. We assume that \mathbf{G} is such that \mathbf{F} is non-singular for some $\boldsymbol{\theta}$ since Equation (9) used. Further, we assume that \mathbf{J} and \mathbf{K} are such that \mathbf{Q} is a positive definite matrix for some $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$ pair. \mathbf{V} is a diagonal matrix containing the standard deviations and \mathbf{C} is the correlation matrix.

We cannot jointly maximize $\boldsymbol{\theta}$, $\boldsymbol{\psi}$, and $\boldsymbol{\phi}$ analytically but we can maximize $\boldsymbol{\theta}$ analytically conditional on $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$. Hence, we can employ a Monte Carlo expectation conditional maximization algorithm in which we take two so-called conditional maximization steps (see Meng and Rubin, 1993, for the, non-Monte Carlo, expectation maximization algorithm). The first conditional maximization step is

$$\boldsymbol{\theta}^{(i+1)} = \mathbf{G}^+ \left(\mathbf{Q}^{(i)} \otimes (\mathbf{M}^\top \mathbf{W} \mathbf{M})^{-1} \right) \mathbf{G}^{+\top} \mathbf{G}^\top \text{vec} \left(\mathbf{M}^\top \mathbf{W} \mathbf{N} \mathbf{Q}^{-(i)} \right) \quad (27)$$

where \otimes is the Kronecker product and $\mathbf{Q}^{-(i)}$ is the inverse of $\mathbf{Q}^{(i)}$. Equation (27) is easily computed with the QR decomposition we make to compute for Equation (25). Having obtained the new $\boldsymbol{\theta}^{(i+1)}$, we update \mathbf{F} and denote the new estimate $\widehat{\mathbf{F}}^{(i+1)}$. The second conditional maximization step which updates $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$ is

$$\mathbf{Z} = \left(\mathbf{N} - \mathbf{R} + \widehat{\mathbf{F}}^{(i+1)} \mathbf{M} \right)^\top \mathbf{W} \left(\mathbf{N} - \mathbf{R} + \widehat{\mathbf{F}}^{(i+1)} \mathbf{M} \right)$$

$$\boldsymbol{\psi}^{(i+1)}, \boldsymbol{\phi}^{(i+1)} = \arg \max_{\boldsymbol{\psi}, \boldsymbol{\phi}} -(d-1) \log |\mathbf{Q}(\boldsymbol{\psi}, \boldsymbol{\phi})| - \text{tr} \left(\mathbf{Q}(\boldsymbol{\psi}, \boldsymbol{\phi})^{-1} \mathbf{Z} \right)$$

which can be done numerically. We have made \mathbf{Q} 's dependence on $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$ explicit to emphasize which factors are affected. \mathbf{C} may not be a valid correlation matrix for all $\boldsymbol{\phi} \in \mathbb{R}^k$ for some choices of \mathbf{K} . Thus, the numerical optimization algorithm is constrained to valid correlation matrices. This completes the two conditional maximization steps. The next E-step is then performed using $\boldsymbol{\theta}^{(i+1)}$, $\boldsymbol{\psi}^{(i+1)}$, $\boldsymbol{\phi}^{(i+1)}$. Meng and Rubin (1993, see the discussion) comments that it may be beneficial to perform an E-step between each conditional maximization step when the E-step is relatively cheap. This is not the case here since all the above computations are independent of the number observations, n_{\max} . Thus, if we have a moderately large number of observations at each time point relative to the dimension of the state vector, then we will use most of the computation time performing the E-step.

4.3 Estimating Fixed Effect Coefficients

Next, we turn to estimating the fixed effect coefficients, $\boldsymbol{\omega}$, in Equation (1). If we assume that observations, y_{it} s, are from an exponential family then it is easy to show that the M-step estimator amounts to generalized linear model with N_s observations for each y_{it}

which differ only by an offset term and a weight. The offset term comes from the $\mathbf{x}_{it}^\top \mathbf{R} + \hat{\boldsymbol{\alpha}}_j^{(t)}$ term in Equation (1) for each of the $j = 1, \dots, N_s$ smoothed particles. The corresponding weights are the smoothed weights, $\hat{w}_j^{(t)}$. The problem can be solved in parallel using QR decompositions as in Section 4.1. This is what is done in the current implementation.

Currently, I only take one iteration of the iteratively re-weighted least squares. I gather I have to repeat till convergence though... This is however not nice computationally and the difference in the estimate from one M-step iteration to the next is very minor when you only take one iteratively re-weighted least square iteration...

5 Other Filter and Smoother 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 their article. The non-auxiliary version in Briers et al. (2009) is more feasible as it only requires evaluation of f in the smoothing part of the generalized two-filter smoother (see Equation (46) in their 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. (2009) and Del Moral et al. (2010) they may still be useful as the computational cost in the smoothing step is independent of the number of observations, n_{\max} . Further, the computational cost can be reduced to $\mathcal{O}(N \log(N))$ run times with approximations like 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 approximations as a function of the unknown parameters. Though, I am not sure how well this method scale with higher state dimensions, p .

Kantas et al. (2015) show empirically that it may be worth just using a forward filter. However, the example is with an univariate outcome ($n_{\max} = 1$ not to be confused with the number of periods d). In the problems shown in this vignette, the computational complexity 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_{\max} . 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.

6 Generalized Two-Filter Smoother

The $\mathcal{O}(N^2)$ smother from Briers et al. (2009) 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 (30) comes from the generalized two-filter formula. To motivate the smoother, we use

$$p(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_t) = \tilde{p}(\mathbf{y}_{t:d}) \frac{\tilde{p}(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d})}{\gamma_t(\boldsymbol{\alpha}_t)}$$

to generalize the two-filter formula in Kitagawa (1994) as follows

$$\begin{aligned}
p(\boldsymbol{\alpha}_t | \mathbf{y}_{1:d}) &= \frac{p(\boldsymbol{\alpha}_t | \mathbf{y}_{1:t-1}) p(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_t)}{p(\mathbf{y}_{t:d} | \mathbf{y}_{1:t-1})} \\
&\propto p(\boldsymbol{\alpha}_t | \mathbf{y}_{1:t-1}) p(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_t) \\
&= p(\boldsymbol{\alpha}_t | \mathbf{y}_{1:t-1}) \tilde{p}(\mathbf{y}_{t:d}) \frac{\tilde{p}(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d})}{\gamma_t(\boldsymbol{\alpha}_t)} \\
&\propto p(\boldsymbol{\alpha}_t | \mathbf{y}_{1:t-1}) \frac{\tilde{p}(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d})}{\gamma_t(\boldsymbol{\alpha}_t)} \\
&= \tilde{p}(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d}) \frac{[\int p(\boldsymbol{\alpha}_{t-1} | \mathbf{y}_{1:t-1}) f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}) d\boldsymbol{\alpha}_{t-1}]}{\gamma_t(\boldsymbol{\alpha}_t)} \\
&\propto \sum_{i=1}^N \tilde{w}_t^{(i)} \delta_{\tilde{\boldsymbol{\alpha}}_t^{(i)}}(\boldsymbol{\alpha}_t) \frac{[\sum_{j=1}^N w_{t-1}^{(j)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)})]}{\gamma_t(\tilde{\boldsymbol{\alpha}}_t^{(i)})}
\end{aligned} \tag{28}$$

where \propto means approximately proportional. Similar arguments leads to

$$\begin{aligned}
p(\boldsymbol{\alpha}_{t-1:t} | \mathbf{y}_{1:d}) &\propto p(\boldsymbol{\alpha}_{t-1:t} | \mathbf{y}_{1:t-1}) p(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_{t-1:t}) \\
&= f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}) p(\boldsymbol{\alpha}_{t-1} | \mathbf{y}_{1:t-1}) p(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_t) \\
&\propto f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}) p(\boldsymbol{\alpha}_{t-1} | \mathbf{y}_{1:t-1}) \frac{\tilde{p}(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d})}{\gamma_t(\boldsymbol{\alpha}_t)} \\
&\propto \sum_{i=1}^N \sum_{j=1}^N \tilde{w}_t^{(i)} \delta_{\tilde{\boldsymbol{\alpha}}_t^{(i)}}(\boldsymbol{\alpha}_t) \frac{[\sum_{k=1}^N w_{t-1}^{(k)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(k)})]}{\gamma_t(\tilde{\boldsymbol{\alpha}}_t^{(i)})} \\
&\quad \cdot \frac{w_{t-1}^{(j)} \delta_{\boldsymbol{\alpha}_{t-1}^{(j)}}(\boldsymbol{\alpha}_{t-1}) f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)})}{[\sum_{k=1}^N w_{t-1}^{(k)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(k)})]} \\
&= \sum_{i=1}^N \sum_{j=1}^N \hat{w}_t^{(i,j)} \delta_{\tilde{\boldsymbol{\alpha}}_t^{(i)}}(\boldsymbol{\alpha}_t) \delta_{\boldsymbol{\alpha}_{t-1}^{(j)}}(\boldsymbol{\alpha}_{t-1})
\end{aligned}$$

where

$$\hat{w}_t^{(i,j)} = \hat{w}_t^{(i)} \frac{w_{t-1}^{(j)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)})}{[\sum_{j=1}^N w_{t-1}^{(j)} f(\tilde{\boldsymbol{\alpha}}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)})]} \tag{29}$$

We need the latter for the Monte Carlo EM-algorithm.

Algorithm 4 $\mathcal{O}(N^2)$ generalized two-filter smoother using the method suggested by Briers et al. (2009).

Input:

- $\mathbf{Q}, \mathbf{Q}_0, \mathbf{a}_0, \mathbf{X}_1, \dots, \mathbf{X}_d, \mathbf{y}_1, \dots, \mathbf{y}_d, R_1, \dots, R_d, \boldsymbol{\omega}$
- 1: **procedure** FILTER FORWARD
 - 2: Run a forward particle filter to get 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 $\tilde{p}(\boldsymbol{\alpha}_t | \mathbf{y}_{t:d})$ 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

$$\hat{w}_t^{(i)} \propto \tilde{w}_t^{(i)} \frac{\left[\sum_{j=1}^N w_{t-1}^{(j)} f\left(\tilde{\boldsymbol{\alpha}}_t^{(i)} \mid \boldsymbol{\alpha}_{t-1}^{(j)}\right) \right]}{\gamma_t\left(\tilde{\boldsymbol{\alpha}}_t^{(i)}\right)} \quad (30)$$

With the result above, we can show the arguments behind the smoother from Fearnhead et al. (2010). Similar to Equation (28), we find that

$$\begin{aligned} p(\boldsymbol{\alpha}_t | \mathbf{y}_{1:d}) &\propto p(\boldsymbol{\alpha}_t | \mathbf{y}_{1:t-1}) p(\mathbf{y}_{t:d} | \boldsymbol{\alpha}_t) \\ &= p(\boldsymbol{\alpha}_t | \mathbf{y}_{1:t-1}) g_t(\mathbf{y}_t | \boldsymbol{\alpha}_t) p(\mathbf{y}_{t+1:d} | \boldsymbol{\alpha}_t) \\ &= \int f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}) p(\boldsymbol{\alpha}_{t-1} | \mathbf{y}_{1:t-1}) d\boldsymbol{\alpha}_{t-1} g_t(\mathbf{y}_t | \boldsymbol{\alpha}_t) \\ &\quad \cdot \int f(\boldsymbol{\alpha}_{t+1} | \boldsymbol{\alpha}_t) p(\mathbf{y}_{t+1:d} | \boldsymbol{\alpha}_{t+1}) d\boldsymbol{\alpha}_{t+1} \\ &\propto \int f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}) p(\boldsymbol{\alpha}_{t-1} | \mathbf{y}_{1:t-1}) d\boldsymbol{\alpha}_{t-1} g_t(\mathbf{y}_t | \boldsymbol{\alpha}_t) \\ &\quad \cdot \int f(\boldsymbol{\alpha}_{t+1} | \boldsymbol{\alpha}_t) \frac{\tilde{p}(\boldsymbol{\alpha}_{t+1} | \mathbf{y}_{t+1:d})}{\gamma_{t+1}(\boldsymbol{\alpha}_{t+1})} d\boldsymbol{\alpha}_{t+1} \\ &\propto \sum_{j=1}^N \sum_{k=1}^N f(\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1}^{(j)}) w_{t-1}^{(j)} g_t(\mathbf{y}_t | \boldsymbol{\alpha}_t) f(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)} | \boldsymbol{\alpha}_t) \frac{\tilde{w}_{t+1}^{(k)}}{\gamma_{t+1}(\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)})} \end{aligned}$$

Thus, we can sample $\boldsymbol{\alpha}_t$ from a proposal distribution given the time $t-1$ forward filter particle, $\boldsymbol{\alpha}_{t-1}^{(j)}$, and time $t+1$ backward filter particle, $\tilde{\boldsymbol{\alpha}}_{t+1}^{(k)}$, for all N^2 particle pairs. Alternatively, we can sample the $t-1$ and $t+1$ particles independently which yields Algorithm 1.

Further, we can find that

$$\begin{aligned}
p(\boldsymbol{\alpha}_{t-1:t}|\mathbf{y}_{1:d}) &= p(\boldsymbol{\alpha}_{t-1:t}|\mathbf{y}_{1:t-1}) g_t(\mathbf{y}_t|\boldsymbol{\alpha}_{t-1:t}) p(\mathbf{y}_{t+1:d}|\boldsymbol{\alpha}_{t-1:t}) \\
&\propto f(\boldsymbol{\alpha}_t|\boldsymbol{\alpha}_{t-1}) p(\boldsymbol{\alpha}_{t-1}|\mathbf{y}_{1:t-1}) g_t(\mathbf{y}_t|\boldsymbol{\alpha}_t) \\
&\quad \cdot \int f(\boldsymbol{\alpha}_{t+1}|\boldsymbol{\alpha}_t) \frac{\tilde{p}(\boldsymbol{\alpha}_{t+1}|\mathbf{y}_{t+1:d})}{\gamma_{t+1}(\boldsymbol{\alpha}_{t+1})} d\boldsymbol{\alpha}_{t+1} \\
&\approx \sum_{i=1}^{N_s} \delta_{\hat{\boldsymbol{\alpha}}_t^{(i)}}(\boldsymbol{\alpha}_t) \delta_{\boldsymbol{\alpha}_{t-1}^{(j_i)}}(\boldsymbol{\alpha}_t) f(\hat{\boldsymbol{\alpha}}_t^{(i)}|\boldsymbol{\alpha}_{t-1}^{(j_i)}) w_{t-1}^{(j_i)} g_t(\mathbf{y}_t|\hat{\boldsymbol{\alpha}}_t^{(i)}) \\
&\quad \cdot \int f(\boldsymbol{\alpha}_{t+1}|\hat{\boldsymbol{\alpha}}_t^{(i)}) \frac{\tilde{p}(\boldsymbol{\alpha}_{t+1}|\mathbf{y}_{t+1:d})}{\gamma_{t+1}(\boldsymbol{\alpha}_{t+1})} d\boldsymbol{\alpha}_{t+1} \\
&\approx \sum_{i=1}^{N_s} \hat{w}_t^{(i)} \delta_{\hat{\boldsymbol{\alpha}}_t^{(i)}}(\boldsymbol{\alpha}_t) \delta_{\boldsymbol{\alpha}_{t-1}^{(j_i)}}(\boldsymbol{\alpha}_t)
\end{aligned}$$

where superscripts j_i are used as in Algorithm 1 which implicitly dependent on t .

7 Gradient and Observed Information Matrix

An alternative to the Monte Carlo EM-algorithm is to approximate the gradient and use it to perform the maximization with a gradient descent algorithm. Moreover, one may be interested in the observed information matrix. Two methods are implemented in order to make such approximations. The first method is the method covered in Cappé et al. (2005, section 8.3 and chapter 11). It has the advantage that it uses the output from the forward particle filter. However, the variance of the estimates increase at least quadratically in time, d . An alternative is to use the method shown by Poyiadjis et al. (2011). Like the smoothing algorithm from Briers et al. (2009), this method has the disadvantage of having a computational complexity which is quadratic in the number of particles, N .

I will give a brief introduction to the two methods in this section. What is presented here closely follow Poyiadjis et al. (2011). First, we will need some notation. We denote the complete data log-likelihood by

$$\begin{aligned}
c(\mathbf{y}_{1:t}, \boldsymbol{\alpha}_{0:t}) &= \log h(\mathbf{y}_{1:t}, \boldsymbol{\alpha}_{0:t}) \\
h(\mathbf{y}_{1:t}, \boldsymbol{\alpha}_{0:t}) &= \nu(\boldsymbol{\alpha}_0) \prod_{k=1}^t g_k(\mathbf{y}_k|\boldsymbol{\alpha}_k) f(\boldsymbol{\alpha}_k|\boldsymbol{\alpha}_{k-1})
\end{aligned}$$

where ν is the density function of the state vector at time zero, all functions may implicitly depend on the unknown parameters, and the dimension of the arguments to c and h is given by the superscript of the arguments. A direct application of the results by Louis (1982) shows that the gradient of the observed data log-likelihood

$$o(\mathbf{y}_{1:t}) = \log \int h(\mathbf{y}_{1:t}, \mathbf{a}_{0:t}) d\mathbf{a}_{0:t}$$

w.r.t. the unknown parameters are

$$\begin{aligned}
\nabla o(\mathbf{y}_{1:t}) &= \frac{\partial}{\partial \boldsymbol{\theta}} \log \int h(\mathbf{y}_{1:t}, \mathbf{a}_{0:t}) d\mathbf{a}_{0:t} = \frac{\int h'(\mathbf{y}_{1:t}, \mathbf{a}_{0:t}) d\mathbf{a}_{0:t}}{\int h(\mathbf{y}_{1:t}, \mathbf{a}_{0:t}) d\mathbf{a}_{0:t}} \\
&= \int c'(\mathbf{y}_{1:t}, \mathbf{a}_{0:t}) p(\mathbf{a}_{0:t}|\mathbf{y}_{1:t}) d\mathbf{a}_{0:t}
\end{aligned} \tag{31}$$

where θ are the unknown parameters in the model, derivatives are w.r.t. θ , and $p(\mathbf{a}_{0:t}|\mathbf{y}_{1:t})$ is the conditional density function of $\mathbf{a}_{0:t}$ given $\mathbf{y}_{1:t}$. Moreover, the Hessian is

$$\begin{aligned}\nabla^2 o(\mathbf{y}_{1:t}) &= \frac{\int h''(\mathbf{y}_{1:t}, \mathbf{a}_{0:t}) d\mathbf{a}_{0:t}}{\int h(\mathbf{y}_{1:t}, \mathbf{a}_{0:t}) d\mathbf{a}_{0:t}} - \nabla o(\mathbf{y}_{1:t}) \nabla o(\mathbf{y}_{1:t})^\top \\ &= \int c''(\mathbf{y}_{1:t}, \mathbf{a}_{0:t}) p(\mathbf{a}_{0:t}|\mathbf{y}_{1:t}) d\mathbf{a}_{0:t} \\ &\quad + \int c'(\mathbf{y}_{1:t}, \mathbf{a}_{0:t}) c'(\mathbf{y}_{1:t}, \mathbf{a}_{0:t})^\top p(\mathbf{a}_{0:t}|\mathbf{y}_{1:t}) d\mathbf{a}_{0:t} - \nabla o(\mathbf{y}_{1:t}) \nabla o(\mathbf{y}_{1:t})^\top\end{aligned}\tag{32}$$

We can use that the forward particle filter yields not just an approximation of $p(\mathbf{a}_d|\mathbf{y}_{1:d})$ but the entire path $p(\mathbf{a}_{0:d}|\mathbf{y}_{1:d})$. That is, we can use the weights at time d from Equation (17) and make a discrete approximation of Equation (31) and (32) as shown in Cappé et al. (2005). However, the variance of the estimates grows at-least quadratically in d as shown by Poyiadjis et al. (2011). The issue is that for larger d , then few if not only one unique value of the initial state vector values (α_i with $0 \leq i \ll d$) are present in the discrete approximation.

As an alternative, Poyiadjis et al. (2011) develop a marginal version of Equation (31) and (32). That is,

$$\begin{aligned}\tilde{c}(\mathbf{y}_{1:t}, \alpha_t) &= \log \tilde{h}(\mathbf{y}_{1:t}, \alpha_t) \\ \tilde{h}(\mathbf{y}_{1:t}, \alpha_t) &= \begin{cases} g_t(\mathbf{y}_t|\alpha_t) \int f(\alpha_t|\mathbf{a}_{t-1}) \tilde{h}(\mathbf{y}_{1:(t-1)}, \mathbf{a}_{t-1}) d\mathbf{a}_{t-1} & t > 0 \\ \nu(\alpha_t) & t = 0 \end{cases}\end{aligned}\tag{33}$$

$$\nabla o(\mathbf{y}_{1:t}) = \int \tilde{c}'(\mathbf{y}_{1:t}, \mathbf{a}_t) p(\mathbf{a}_t|\mathbf{y}_{1:t}) d\mathbf{a}_t\tag{34}$$

$$\begin{aligned}\nabla^2 o(\mathbf{y}_{1:t}) &= \int \tilde{c}''(\mathbf{y}_{1:t}, \mathbf{a}_t) p(\mathbf{a}_t|\mathbf{y}_{1:t}) d\mathbf{a}_t \\ &\quad + \int \tilde{c}'(\mathbf{y}_{1:t}, \mathbf{a}_t) \tilde{c}'(\mathbf{y}_{1:t}, \mathbf{a}_t)^\top p(\mathbf{a}_t|\mathbf{y}_{1:t}) d\mathbf{a}_t - \nabla o(\mathbf{y}_{1:t}) \nabla o(\mathbf{y}_{1:t})^\top \\ &= \int \left(\frac{\tilde{h}''(\mathbf{y}_{1:t}, \mathbf{a}_t)}{\tilde{h}(\mathbf{y}_{1:t}, \mathbf{a}_t)} - \tilde{c}'(\mathbf{y}_{1:t}, \mathbf{a}_t) \tilde{c}'(\mathbf{y}_{1:t}, \mathbf{a}_t)^\top \right) p(\mathbf{a}_t|\mathbf{y}_{1:t}) d\mathbf{a}_t \\ &\quad + \int \tilde{c}'(\mathbf{y}_{1:t}, \mathbf{a}_t) \tilde{c}'(\mathbf{y}_{1:t}, \mathbf{a}_t)^\top p(\mathbf{a}_t|\mathbf{y}_{1:t}) d\mathbf{a}_t - \nabla o(\mathbf{y}_{1:t}) \nabla o(\mathbf{y}_{1:t})^\top\end{aligned}\tag{35}$$

While there is no analytical expression for the derivatives then one can establish a point-wise approximation recursively for $c'(\mathbf{y}_{1:t}, \mathbf{a}_t)$ and $c''(\mathbf{y}_{1:t}, \mathbf{a}_t)$ as suggested by Poyiadjis et al. (2011). To see this, let

$$s_t(\alpha_t, \alpha_{t-1}) = \log g_t(\mathbf{y}_t|\alpha_t) + \log f(\alpha_t|\alpha_{t-1})$$

Then

$$\begin{aligned}\tilde{h}'(\mathbf{y}_{1:t}, \alpha_t) &= o(\mathbf{y}_{1:t-1}) g_t(\mathbf{y}_t|\alpha_t) \int f(\alpha_t|\mathbf{a}_{t-1}) p(\mathbf{a}_{t-1}|\mathbf{y}_{1:t-1}) \\ &\quad \cdot (s'_t(\alpha_t, \mathbf{a}_{t-1}) + \tilde{c}'(\mathbf{y}_{1:(t-1)}, \mathbf{a}_{t-1})) d\mathbf{a}_{t-1}\end{aligned}\tag{36}$$

Taking the ratio of Equation (36) and (33) yields $\tilde{c}'(\mathbf{y}_{1:t}, \boldsymbol{\alpha}_t)$ in Equation (34). Moreover, for Equation (35)

$$\begin{aligned} \tilde{h}''(\mathbf{y}_{1:t}, \boldsymbol{\alpha}_t) &= o(\mathbf{y}_{1:t-1}) g_t(\mathbf{y}_t | \boldsymbol{\alpha}_t) \int f(\boldsymbol{\alpha}_t | \mathbf{a}_{t-1}) p(\mathbf{a}_{t-1} | \mathbf{y}_{1:t-1}) \\ &\quad \cdot \left((s'_t(\boldsymbol{\alpha}_t, \mathbf{a}_{t-1}) + \tilde{c}'(\mathbf{y}_{1:(t-1)}, \mathbf{a}_{t-1})) (s'_t(\boldsymbol{\alpha}_t, \mathbf{a}_{t-1}) + \tilde{c}'(\mathbf{y}_{1:(t-1)}, \mathbf{a}_{t-1}))^\top \right. \\ &\quad \left. + s''_t(\boldsymbol{\alpha}_t, \mathbf{a}_{t-1}) + \tilde{c}''(\mathbf{y}_{1:(t-1)}, \mathbf{a}_{t-1}) \right) d\mathbf{a}_{t-1} \end{aligned}$$

where we again take the ratio with (33). Unlike before, we need to evaluate two ratios, $\tilde{h}'(\mathbf{y}_{1:t}, \boldsymbol{\alpha}_t) / \tilde{h}(\mathbf{y}_{1:t}, \boldsymbol{\alpha}_t)$ and $\tilde{h}''(\mathbf{y}_{1:t}, \boldsymbol{\alpha}_t) / \tilde{h}(\mathbf{y}_{1:t}, \boldsymbol{\alpha}_t)$, which require evaluation of expressions of the form

$$\frac{\int f(\boldsymbol{\alpha}_t | \mathbf{a}_{t-1}) p(\mathbf{a}_{t-1} | \mathbf{y}_{1:t-1}) \kappa_t(\boldsymbol{\alpha}_t, \mathbf{a}_{t-1}) d\mathbf{a}_{t-1}}{\int f(\boldsymbol{\alpha}_t | \mathbf{a}_{t-1}) p(\mathbf{a}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{a}_{t-1}} \quad (37)$$

for some function κ_t . To do so, redefine the weights in Equation (17) in the forward particle filter shown in Algorithm 2 to

$$w_t^{(i)} \propto \frac{g_t(\mathbf{y}_t | \boldsymbol{\alpha}_t^{(i)}) \sum_{j=1}^N f(\boldsymbol{\alpha}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)}) w_{t-1}^{(j)}}{q\left(\boldsymbol{\alpha}_t^{(i)} \left| \left\{ (\boldsymbol{\alpha}_{t-1}^{(j)}, w_{t-1}^{(j)}) \right\}_{j=1, \dots, N}, \mathbf{y}_t \right.)\right)}$$

where we have made it explicit that the proposal distribution may depend on the previous particle cloud and assume that we use the same number of particles at time 0. Further, define the weights

$$\bar{w}_t^{(i,j)} = \frac{f(\boldsymbol{\alpha}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(j)}) w_{t-1}^{(j)}}{\sum_{k=1}^N f(\boldsymbol{\alpha}_t^{(i)} | \boldsymbol{\alpha}_{t-1}^{(k)}) w_{t-1}^{(k)}} \quad (38)$$

Now a discrete approximation of expression in Equation (37) is given by

$$\sum_{j=1}^N \bar{w}_t^{(i,j)} \kappa_t(\boldsymbol{\alpha}_t^{(i)}, \boldsymbol{\alpha}_{t-1}^{(j)})$$

Thus, the recursive formula for the gradient approximation is

$$\begin{aligned} \zeta_t^{(i)} &= \sum_{j=1}^N \bar{w}_t^{(i,j)} \left(s'_t(\boldsymbol{\alpha}_t^{(i)}, \boldsymbol{\alpha}_{t-1}^{(j)}) + \zeta_{t-1}^{(j)} \right) \\ \nabla o(\mathbf{y}_{1:t}) &\approx \sum_{i=1}^N w_t^{(i)} \zeta_t^{(i)} \end{aligned} \quad (39)$$

and for the Hessian we have

$$\begin{aligned} \mathbf{r}_t^{(i)} &= \sum_{j=1}^N \bar{w}_t^{(i,j)} \left(\left(s'_t(\boldsymbol{\alpha}_t^{(i)}, \boldsymbol{\alpha}_{t-1}^{(j)}) + \zeta_{t-1}^{(j)} \right) \left(s'_t(\boldsymbol{\alpha}_t^{(i)}, \boldsymbol{\alpha}_{t-1}^{(j)}) + \zeta_{t-1}^{(j)} \right)^\top \right. \\ &\quad \left. + s''_t(\boldsymbol{\alpha}_t^{(i)}, \boldsymbol{\alpha}_{t-1}^{(j)}) + \mathbf{r}_{t-1}^{(j)} \right) - \zeta_t^{(i)} \zeta_t^{(i)\top} \end{aligned} \quad (40)$$

such that

$$\nabla^2 o(\mathbf{y}_{1:t}) \approx \sum_{i=1}^N w_t^{(i)} \left(\boldsymbol{\zeta}_t^{(i)} \boldsymbol{\zeta}_t^{(i)\top} + \boldsymbol{\Upsilon}_t^{(i)} \right) - \nabla o(\mathbf{y}_{1:t}) \nabla o(\mathbf{y}_{1:t})^\top$$

The issue with the latter method is that the method has an $\mathcal{O}(N^2)$ computational complexity because of the sums in Equation (38), (39), and (40). A particular type of particle filters that are well suited for approximations like those in Equation (37) are the so-called independent particle filters suggest by Lin et al. (2005). The key point in these filters is that they use a proposal distribution which only depends on the observed outcome, \mathbf{y}_t , or also the previous particle cloud but not any particular particle.

An alternative to the methods in the `dynamichazard` package is the `mssm` package. It contains both the method shown in Cappé et al. (2005) and the method suggested by Poyiadjis et al. (2011) but for more general models. Moreover, the `mssm` package has an implementation of the dual k-d tree approximation method like in Klaas et al. (2006). This reduces the average-case complexity to $\mathcal{O}(N \log N)$ and thus it allows one to use substantially more particles. Lastly, the `mssm` also allows for two types of antithetic variables like those suggest by Durbin and Koopman (1997). This decreases the variance of the estimates for the same computational cost.

References

- Mark Briers, Arnaud Doucet, and Simon Maskell. Smoothing algorithms for state-space models. *Annals of the Institute of Statistical Mathematics*, 62(1):61, Jun 2009. ISSN 1572-9052. doi: 10.1007/s10463-009-0236-2. URL <https://doi.org/10.1007/s10463-009-0236-2>.
- Olivier Cappé, Eric Moulines, and Tobias Ryden. *Inference in Hidden Markov Models*. Springer-Verlag New York, 2005. ISBN 978-0-387-40264-2, 978-1-4419-2319-6.
- 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.
- J. Durbin and S. J. Koopman. Monte carlo maximum likelihood estimation for non-Gaussian state space models. *Biometrika*, 84(3):669–684, 1997. ISSN 00063444. URL <http://www.jstor.org/stable/2337587>.
- 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. The two-filter formula for smoothing and an implementation of the Gaussian-sum smoother. *Annals of the Institute of Statistical Mathematics*, 46(4):605–623, Dec 1994. ISSN 1572-9052. doi: 10.1007/BF00773470. URL <https://doi.org/10.1007/BF00773470>.
- 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.
- Ming T. Lin, Junni L. Zhang, Qiansheng Cheng, and Rong Chen. Independent particle filters. *Journal of the American Statistical Association*, 100(472):1412–1421, 2005. ISSN 01621459. URL <http://www.jstor.org/stable/27590681>.
- Thomas A. Louis. Finding the observed information matrix when using the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 44(2):226–233, 1982. ISSN 00359246. URL <http://www.jstor.org/stable/2345828>.

- Sheheryar Malik and Michael K Pitt. Particle filters for continuous likelihood evaluation and maximisation. *Journal of Econometrics*, 165(2):190–209, 2011.
- Xiao-Li Meng and Donald B. Rubin. Maximum likelihood estimation via the ECM algorithm: A general framework. *Biometrika*, 80(2):267–278, 1993. ISSN 00063444. URL <http://www.jstor.org/stable/2337198>.
- Michael K Pitt and Neil Shephard. Filtering via simulation: Auxiliary particle filters. *Journal of the American statistical association*, 94(446):590–599, 1999.
- George Poyiadjis, Arnaud Doucet, and Sumeetpal S. Singh. Particle approximations of the score and observed information matrix in state space models with application to parameter estimation. *Biometrika*, 98(1):65–80, 2011. ISSN 00063444. URL <http://www.jstor.org/stable/29777165>.
- Thomas B Schön, Adrian Wills, and Brett Ninness. System identification of nonlinear state-space models. *Automatica*, 47(1):39–49, 2011.
- Simon Wood, Yannig Goude, and Simon Shaw. Generalized additive models for large data sets. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 64(1):139–155, 2014. doi: 10.1111/rssc.12068. URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/rssc.12068>.