

SMC for Recursive Parameter Estimation in State-Space Models

Arnaud Doucet
Departments of Statistics & Computer Science
University of British Columbia

Problem Statement

- $\{X_n\}_{n \geq 1}$ latent/hidden Markov process with

$$X_1 \sim \mu_\theta(\cdot) \text{ and } X_n | (X_{n-1} = x) \sim f_\theta(\cdot | x).$$

- $\{Y_n\}_{n \geq 1}$ observation process such that observations are conditionally independent given $\{X_n\}_{n \geq 1}$ and

$$Y_n | (X_n = x) \sim g_\theta(\cdot | x).$$

- **Objectives:** Assume the observations available correspond to $\theta = \theta^*$, obtain a recursive algorithm to estimate θ^* .

- *Linear Gaussian state-space model*

$$\begin{aligned}X_1 &\sim \mathcal{N}(0, 1), \quad X_n = \alpha X_{n-1} + \sigma_v V_n, \\Y_n &= X_n + \sigma_w W_n\end{aligned}$$

where $V_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$, $W_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. In this case, we have $\theta = (\alpha, \sigma_v, \sigma_w)$.

- *Stochastic volatility model*

$$\begin{aligned}X_1 &\sim \mathcal{N}(0, 1), \quad X_n = \alpha X_{n-1} + \sigma_v V_n, \\Y_n &= \beta \exp(X_n/2) W_n\end{aligned}$$

where $V_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$, $W_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. In this case, we have $\theta = (\alpha, \sigma_v, \beta)$.

Approaches to Recursive Parameter Estimation

- Bayesian approaches where θ is an unknown random parameter with a prior $p(\theta)$. In this case, inference relies on the sequence of distributions $p(\theta | y_{1:n})$.
- Point estimation based on recursive Maximum Likelihood and pseudo-likelihood approaches.

- In a Bayesian framework, θ is an unknown random parameter with a prior $p(\theta)$.
- At time n , inference relies on

$$p(\theta | y_{1:n}) = \int p(x_{1:n}, \theta | y_{1:n}) dx_{1:n}$$

where

$$p(x_{1:n}, \theta | y_{1:n}) \propto p(y_{1:n} | x_{1:n}, \theta) p(x_{1:n} | \theta) p(\theta).$$

- We know the sequence of distributions $p(x_{1:n}, \theta | y_{1:n})$ up to a normalizing constant so we can use SMC methods.

Preliminary Warning

- We have

$$p(\theta | y_{1:n}) = \frac{p(y_{1:n} | \theta) p(\theta)}{p(y_{1:n})}$$

- We have seen previously that, even for a fixed value θ , the SMC estimate $\hat{p}(y_{1:n} | \theta)$ of $p(y_{1:n} | \theta)$ is under favourable mixing assumptions such that

$$\frac{\mathbb{V}[\hat{p}(y_{1:n} | \theta)]}{p(y_{1:n} | \theta)^2} \leq C \frac{n}{N};$$

i.e. the performance degrade linearly with the time index n .

- Intuitively, estimating the whole posterior $p(\theta | y_{1:n})$ is obviously more difficult than estimating $p(y_{1:n} | \theta)$ for a specific value of θ . Hence the SMC algorithms targetting $p(\theta | y_{1:n})$ might not enjoy very good convergence properties... Indeed this is unfortunately the case.

- Numerous SMC schemes have been proposed to address this problem.
- I will only discuss schemes providing asymptotically consistent estimates of $p(x_{1:n}, \theta | y_{1:n})$, hence of $p(\theta | y_{1:n})$; i.e. for n fixed we have convergence for $N \rightarrow \infty$.
- Approaches introducing some artificial random walk dynamics on the parameter/making fixed-lag approximations do not satisfy this property.

Naive SMC Scheme for Parameter Estimation

- Sample $(X_1^{(i)}, \theta_0^{(i)}) \sim q(\cdot, \cdot | y_1)$ and

$$W_1^{(i)} \propto \frac{p(\theta_0^{(i)}) \mu_{\theta_0^{(i)}}(X_1^{(i)}) g_{\theta_0^{(i)}}(y_1 | X_1^{(i)})}{q(X_1^{(i)}, \theta_0^{(i)} | y_1)}.$$

- Resample $\{(X_1^{(i)}, \theta_0^{(i)}), W_1^{(i)}\}$ to obtain particles $\{X_1^{(i)}, \theta_1^{(i)}\}$

- At time $n \geq 2$, sample $X_n^{(i)} \sim q_{\theta_{n-1}^{(i)}}(\cdot | y_n, X_{n-1}^{(i)})$ and

$$W_n^{(i)} \propto \frac{f_{\theta_{n-1}^{(i)}}(X_n^{(i)} | X_{n-1}^{(i)}) g_{\theta_{n-1}^{(i)}}(y_n | X_n^{(i)})}{q_{\theta_{n-1}^{(i)}}(X_n^{(i)} | y_n, X_{n-1}^{(i)})}.$$

- Resample $\{(X_{1:n}^{(i)}, \theta_{n-1}^{(i)}), W_n^{(i)}\}$ to obtain particles $\{X_{1:n}^{(i)}, \theta_n^{(i)}\}$

- This is just a standard SMC scheme...
- We have

$$\hat{p}(x_{1:n}, \theta | y_{1:n}) = \sum_{i=1}^N W_n^{(i)} \delta_{(x_{1:n}, \theta_n^{(i)})} (x_{1:n}, \theta).$$

- In particular, we have

$$\hat{p}(\theta | y_{1:n}) = \sum_{i=1}^N W_n^{(i)} \delta_{\theta_n^{(i)}}(\theta)$$

where $\theta_n^{(i)}$ correspond to the particles having been sampled at time 1 which have survived to the resampling steps at time $1, 2, \dots, n$.

- This algorithm provides an asymptotically consistent estimate of the targets under very weak assumptions....
- ... and yes it is a very bad algorithm. We only sample particles in the Θ space at time 1; this is followed by successive resampling steps.
- After a few time steps, we have

$$\hat{p}(\theta | y_{1:n}) = \delta_{\bar{\theta}}(\theta)$$

where $\theta_n^{(i)} = \bar{\theta}$ for $i \in \{1, \dots, N\}$. This is somewhat similar to the problem we faced before when there was no unknown parameter but we were interested in estimating $p(x_1 | y_{1:n})$... but the problem is even worse as, because of the lack of ergodicity, this error propagate itself.

- Theoretically, it means that we do not have a uniform convergence result for $\hat{p}(\theta | y_{1:n})$; only the following very weak result

$$\mathbb{E} \left[\left| \int \varphi(\theta) (\hat{p}(d\theta | y_{1:n}) - p(d\theta | y_{1:n})) \right|^p \right]^{1/p} \leq \frac{c(n)}{\sqrt{N}}$$

where $c(n)$ increases over time.

How to improve performance?

- We can use all the advanced methods discussed previously: auxiliary method, resample-move, block sampling.
- Resample move is especially attractive in this context: it consists in adding at time n an MCMC move $K_n(x'_{1:n}, \theta' | x_{1:n}, \theta)$ of invariant distribution $p(x_{1:n}, \theta | y_{1:n})$. To keep the algorithm on-line, we can only update a fixed number of variables; say here θ only.
- For example, we could use a Gibbs step

$$K_n(x'_{1:n}, \theta' | x_{1:n}, \theta) = \delta_{x_{1:n}}(x'_{1:n}) p(\theta' | y_{1:n}, x_{1:n})$$

Resample Move SMC for Parameter Estimation

- Sample $(X_1^{(i)}, \theta_0^{(i)}) \sim q(\cdot, \cdot | y_1)$ and

$$W_1^{(i)} \propto \frac{p(\theta_0^{(i)}) \mu_{\theta_0^{(i)}}(X_1^{(i)}) g_{\theta_0^{(i)}}(y_1 | X_1^{(i)})}{q(X_1^{(i)}, \theta_0^{(i)} | y_1)}.$$

- Resample $\left\{ (X_1^{(i)}, \theta_0^{(i)}), W_1^{(i)} \right\}$ to obtain particles $\left\{ X_1^{(i)}, \bar{\theta}_1^{(i)} \right\}$.

- Sample $\theta_1^{(i)} \sim p(\cdot | y_1, X_1^{(i)})$.

- At time $n \geq 2$, sample $X_n^{(i)} \sim q_{\theta_{n-1}^{(i)}}(\cdot | y_n, X_{n-1}^{(i)})$ and

$$W_n^{(i)} \propto \frac{f_{\theta_{n-1}^{(i)}}(X_n^{(i)} | X_{n-1}^{(i)}) g_{\theta_{n-1}^{(i)}}(y_n | X_n^{(i)})}{q_{\theta_{n-1}^{(i)}}(X_n^{(i)} | y_n, X_{n-1}^{(i)})}.$$

- Resample $\left\{ (X_{1:n}^{(i)}, \theta_{n-1}^{(i)}), W_n^{(i)} \right\}$ to obtain particles $\left\{ X_{1:n}^{(i)}, \bar{\theta}_n^{(i)} \right\}$.

- Sample $\theta_n^{(i)} \sim p(\cdot | y_{1:n}, X_{1:n}^{(i)})$.

- At first glance, this algorithm seems difficult to implement as it requires storing the paths $\{X_{1:n}^{(i)}\}$ so memory requirements increase.
- However, in many practical applications, we have

$$p(\theta | y_{1:n}, x_{1:n}) = p(\theta | s_n(x_{1:n}, y_{1:n}))$$

i.e. it depends only on a set of sufficient statistics $s_n(x_{1:n}, y_{1:n})$ of fixed dimension.

Example: Linear Gaussian state-space model

- We have

$$\begin{aligned}X_1 &\sim \mathcal{N}(0, 1), \quad X_n = \alpha X_{n-1} + \sigma_v V_n, \\Y_n &= X_n + \sigma_w W_n\end{aligned}$$

where $V_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$, $W_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$.

- Assume for sake of simplicity that only α is unknown with $p(\alpha) = \mathcal{U}_{[-1,1]}(\alpha)$.
- It is easy to check that

$$p(\alpha | y_{1:n}, x_{1:n}) \propto \mathcal{N}(\alpha; m_n, \sigma_n^2) \mathbf{1}_{[-1,1]}(\alpha)$$

where

$$\sigma_n^2 = \left(\sum_{k=1}^{n-1} x_k^2 \right)^{-1}, \quad m_n = \sigma_n^2 \left(\sum_{k=2}^n x_{k-1} x_k \right).$$

- In practice, we only need to store $\sum_{k=2}^n x_{k-1} x_k$ and $\sum_{k=1}^{n-1} x_k^2$ instead of $x_{1:n}$.

Resample Move SMC with Sufficient Statistics for Parameter Estimation

- $(X_1^{(i)}, \theta_0^{(i)}) \sim q(\cdot, \cdot | y_1)$ and $W_1^{(i)} \propto \frac{p(\theta_0^{(i)}) \mu_{\theta_0^{(i)}}(X_1^{(i)}) g_{\theta_0^{(i)}}(y_1 | X_1^{(i)})}{q(X_1^{(i)}, \theta_0^{(i)} | y_1)}$.
- Resample $\left\{ (X_1^{(i)}, \theta_0^{(i)}), W_1^{(i)} \right\}$ to obtain $\left\{ X_1^{(i)}, s_1(X_1^{(i)}, y_1), \bar{\theta}_1^{(i)} \right\}$.
- $\theta_1^{(i)} \sim p(\cdot | s_1(X_1^{(i)}, y_1))$.
- At time $n \geq 2$, $X_n^{(i)} \sim q_{\theta_{n-1}^{(i)}}(\cdot | y_n, X_{n-1}^{(i)})$ and $W_n^{(i)} \propto \frac{f_{\theta_{n-1}^{(i)}}(X_n^{(i)} | X_{n-1}^{(i)}) g_{\theta_{n-1}^{(i)}}(y_n | X_n^{(i)})}{q_{\theta_{n-1}^{(i)}}(X_n^{(i)} | y_n, X_{n-1}^{(i)})}$.
- Resample $\left\{ (X_n^{(i)}, s_n(X_{1:n}^{(i)}, y_{1:n}), \theta_{n-1}^{(i)}), W_n^{(i)} \right\}$ to obtain $\left\{ X_n^{(i)}, s_n(X_{1:n}^{(i)}, y_{1:n}), \bar{\theta}_n^{(i)} \right\}$.
- Sample $\theta_n^{(i)} \sim p(\cdot | s_n(X_{1:n}^{(i)}, y_{1:n}))$.

- This algorithm appears elegant.
- This algorithm and some variations have already appeared several times in the literature (Andrieu, De Freitas & D., 1999), (Fearnhead, 2002), (Storvik, 2002), (Johannes & Polson, 2007).
- This algorithm suffers from very severe limitations and is not robust as, once more, it relies implicitly on the SMC approximation of a sequence of distributions $p(x_{1:n}|y_{1:n})$ of increasing dimension; the pitfalls of this approach were first discussed in (Andrieu, De Freitas & D., 1999), see also (Andrieu, D. & Tadic, 2005).

Illustration of the degeneracy phenomenon

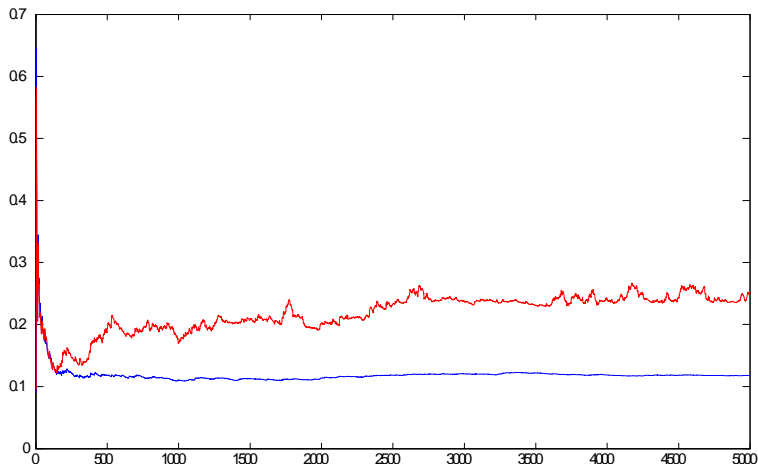


Figure: Sufficient statistics computed exactly through the Kalman smoother (blue) and the SMC method (red).

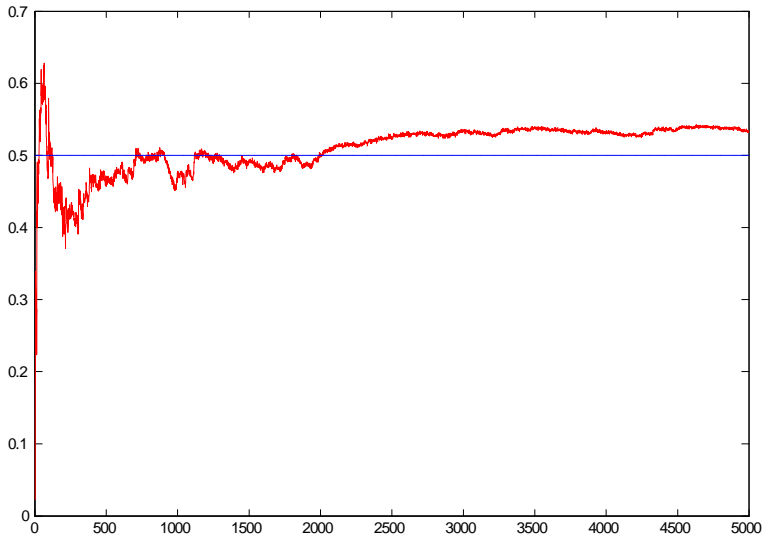


Figure: SMC approximation of $\mathbb{E}[\theta | y_{1:n}]$ for $N = 1000$ particles (red) as a function of n versus true value (blue). The algorithm converges towards a wrong value.

Additional Comments

- These algorithms provide asymptotically consistent approximations; i.e. for fixed n , the SMC approximation converges towards the true target as N increases...
- Still, it does not mean that such algorithms perform well in practice. For a fixed N and an increasing n , the error will increase; i.e. it is not possible to obtain uniform convergence results.
- You can use any advanced method you want but, as long as you rely on an SMC approximation of $p(x_{1:n}|y_{1:n})$ (or $p(s_n(x_{1:n}, y_{1:n})|y_{1:n})$), then you will face the same problem eventually for n large enough.
- For a **fixed time horizon**, and N large enough, such methods might perform reasonably well and cannot be completely ruled out. However you have to be extremely careful: determining a large enough N is difficult (see SMC project for more information and quantitative results).

- The credible intervals estimates computed via such approaches are much tighter than they should be (because of the degeneracy phenomenon) so you cannot trust them.
- You can expect these methods to perform very poorly when the dimension of the parameter space is high; say superior to 5-10.

- It is impossible to obtain an asymptotically convergent SMC algorithm to estimate $p(\theta | y_{1:n})$ which enjoys uniform convergence properties.
- At the price of a non-vanishing bias, it should be possible to obtain much better approximations of $p(\theta | y_{1:n})$ based on fixed-lag approximations. The main problem is that it is difficult to quantify the bias in practical situations.

Recursive Maximum Likelihood

- Recursive Maximum Likelihood is a fairly old and popular approach in the system identification/control community.
- We show here how to implement an SMC version of it for general state-space models.
- Under stationary assumptions (e.g. Tadic & D., 2005), we have

$$\frac{1}{n} \log p_{\theta} (Y_{1:n}) = \frac{1}{n} \sum_{k=1}^n \log p_{\theta} (Y_k | Y_{1:k-1}) \rightarrow l(\theta)$$

with

$$l(\theta) = \int \int_{\mathcal{Y} \times \mathcal{P}(\mathcal{X})} \log \left(\int g_{\theta}(y|x) \mu(x) dx \right) \lambda_{\theta, \theta^*}(dy, d\mu),$$

where $\mathcal{P}(\mathcal{X})$ is the space of probability distributions on \mathcal{X} and $\lambda_{\theta, \theta^*}(dy, d\mu) = \int \lambda_{\theta, \theta^*}(dx, dy, d\mu)$; $\lambda_{\theta, \theta^*}(dx, dy, d\mu)$ being the invariant distribution of the Markov chain $\{X_n, Y_n, p_{\theta}(x_n | Y_{1:n-1})\}_{n \geq 1}$.

Stochastic Approximation

- The set of global maxima of the averaged log-likelihood $l(\theta)$ includes θ^* .
- The function $l(\theta)$ is unknown but can be maximized using a stochastic approximation algorithm

$$\theta_n = \theta_{n-1} + \gamma_n \nabla \log p_{\theta_{1:n-1}}(Y_n | Y_{1:n-1}) \quad (1)$$

where the stepsize sequence $\{\gamma_n\}_{n \geq 1}$ is a positive non-increasing sequence.

- $p_{\theta_{1:n}}(x_n | Y_{1:n})$ denotes the filter computed using θ_{t-1} at time t and similarly for $\nabla \log p_{\theta_{1:n-1}}(Y_n | Y_{1:n-1})$.
- We typically need $\sum \gamma_n = \infty$ and $\sum \gamma_n^2 < \infty$; i.e. one selects $\gamma_n = \gamma_0 \cdot n^{-\alpha}$ where $\gamma_0 > 0$ and $0.5 < \alpha \leq 1$.
- This algorithm is a stochastic gradient algorithm and is not guaranteed to converge towards θ^* ; only to a local maximum of $l(\theta)$.
- For finite-state space hidden Markov models, this algorithm was proposed and studied by (Le Gland & Mevel, 1997).

SMC Approximation

- We need to approximate $\nabla \log p_\theta (Y_n | Y_{1:n-1})$.
- The first approach consists of using

$$\nabla \log p_\theta (Y_n | Y_{1:n-1}) = \nabla \log p_\theta (Y_{1:n}) - \nabla \log p_\theta (Y_{1:n-1})$$

where Fisher's identity yields

$$\nabla \log p_\theta (Y_{1:n}) = \int \nabla \log p_\theta (x_{1:n}, Y_{1:n}) \cdot p_\theta (x_{1:n} | Y_{1:n}) dx_{1:n}$$

with

$$\begin{aligned} \nabla \log p_\theta (x_{1:n}, Y_{1:n}) &= \nabla \log \mu_\theta (x_1) + \nabla \log g_\theta (Y_1 | x_1) \\ &\quad + \sum_{k=2}^n \nabla \log f_\theta (x_k | x_{k-1}) + \nabla \log g_\theta (Y_k | x_k). \end{aligned}$$

SMC Implementation of Fisher's identity

- Given your favourite SMC approximation $\widehat{p}_\theta (x_{1:n} | Y_{1:n})$ of $p_\theta (x_{1:n} | Y_{1:n})$; say

$$\widehat{p}_\theta (x_{1:n} | Y_{1:n}) = \sum_{i=1}^N W_n^{(i)} \delta_{X_{1:n}^{(i)}} (x_{1:n})$$

then we can compute an estimate

$$\widehat{\nabla \log p_\theta} (Y_{1:n}) = \sum_{i=1}^N W_n^{(i)} \nabla \log p_\theta (X_{1:n}^{(i)}, Y_{1:n}).$$

- This estimate can be easily computed recursively using

$$\begin{aligned} \nabla \log p_\theta (x_{1:n}, Y_{1:n}) &= \nabla \log p_\theta (x_{1:n-1}, Y_{1:n-1}) \\ &\quad + \nabla \log f_\theta (x_n | x_{n-1}) + \nabla \log g_\theta (Y_n | x_n). \end{aligned}$$

- We obtain

$\widehat{\nabla \log p_\theta} (Y_n | Y_{1:n-1}) = \widehat{\nabla \log p_\theta} (Y_{1:n}) - \widehat{\nabla \log p_\theta} (Y_{1:n-1})$ but this estimate has poor properties as, once more, it relies implicitly on an approximation of the joint distribution $p_\theta (x_{1:n} | Y_{1:n}) \dots$

SMC Approximation of the Sensitivity Equations

- There is an alternative way to compute $\nabla \log p_\theta (Y_n | Y_{1:n-1})$ based on sensitivity equations.
- We have

$$\nabla \log p_\theta (Y_{n+1} | Y_{1:n}) = \int \nabla \log p_\theta (x_{n+1}, Y_{n+1} | Y_{1:n}) p(x_{n+1} | Y_{1:n+1}) dx_{n+1}$$

with

$$\begin{aligned} \nabla p_\theta (x_{n+1}, Y_{n+1} | Y_{1:n}) &= g_\theta (Y_{n+1} | x_{n+1}) \int f_\theta (x_{n+1} | x_n) p_\theta (x_n | Y_{1:n}) \\ &\times (\nabla \log g_\theta (Y_{n+1} | x_{n+1}) + \nabla \log f_\theta (x_{n+1} | x_n) + \nabla \log p_\theta (x_n | Y_{1:n})) dx_n. \end{aligned}$$

- By differentiating $\nabla \log p_\theta (x_{n+1} | Y_{1:n+1})$, we obtain

$$\begin{aligned} \nabla p_\theta (x_{n+1} | Y_{1:n+1}) &= \frac{\nabla p_\theta (x_{n+1}, Y_{n+1} | Y_{1:n})}{p_\theta (Y_{n+1} | Y_{1:n})} \\ &- p_\theta (x_{n+1} | Y_{1:n+1}) \nabla \log p_\theta (Y_{n+1} | Y_{1:n}) \end{aligned}$$

SMC Approximation of Filter Sensitivity

- To implement this recursion, we need to approximate $\nabla p_\theta(x_n | Y_{1:n})$. This is a signed measure such that

$$\int \nabla p_\theta(x_n | Y_{1:n}) dx_n = 0.$$

- A first idea to approximate $\nabla p_\theta(x_n | Y_{1:n})$ consists of using the identity

$$\nabla p_\theta(x_n | Y_{1:n}) = \int \nabla \log p_\theta(x_{1:n} | Y_{1:n}) \cdot p_\theta(x_{1:n} | Y_{1:n}) dx;$$

this would rely once more on an SMC approximation of $p_\theta(x_{1:n} | Y_{1:n})$... and it is just a convoluted way to rewrite the previous algorithm.

- An alternative consists of using (Poyadjis, D. & Singh, 2005)

$$\nabla p_{\theta}(x_n | Y_{1:n}) = \frac{\nabla p_{\theta}(x_n | Y_{1:n})}{p_{\theta}(x_n | Y_{1:n})} \cdot p_{\theta}(x_n | Y_{1:n});$$

that is if $\hat{p}_{\theta}(x_n | Y_{1:n}) = \sum_{i=1}^N W_n^{(i)} \delta_{X_n^{(i)}}(x_n)$ then

$$\widehat{\nabla} p_{\theta}(x_n | Y_{1:n}) = \sum_{i=1}^N W_n^{(i)} \frac{\widetilde{\nabla} p_{\theta}(X_n^{(i)} | Y_{1:n})}{\widetilde{p}_{\theta}(X_n^{(i)} | Y_{1:n})} \delta_{X_n^{(i)}}(x_n)$$

- This only relies on approximation of the marginals; the price to pay is that we now need a pointwise estimate of $\widetilde{p}_{\theta}(X_n^{(i)} | Y_{1:n})$ and $\widetilde{\nabla} p_{\theta}(X_n^{(i)} | Y_{1:n})$. The algorithm is thus in $O(N^2)$.

- At time $n - 1$, assume approximations of the filtering distribution and its derivatives of the form

$$\widehat{p}_\theta(x_{n-1} | y_{1:n-1}) = \sum_{i=1}^N W_{n-1}^{(i)} \delta_{X_{n-1}^{(i)}}(x_{n-1}),$$

$$\widehat{\nabla} p_\theta(x_{n-1} | y_{1:n-1}) = \sum_{i=1}^N W_{n-1}^{(i)} A_{n-1}^{(i)} \delta_{X_{n-1}^{(i)}}(x_{n-1}),$$

are available where $A_{n-1}^{(i)}$ is an approximation of

$$\nabla p_\theta(X_{n-1}^{(i)} | y_{1:n-1}) / p_\theta(X_{n-1}^{(i)} | y_{1:n-1}).$$

- We obtain the pointwise approximations of $p_\theta(x_n, y_n | y_{1:n-1})$, $\nabla p_\theta(x_n, y_n | y_{1:n-1})$

$$\widetilde{p}_\theta(x_n, y_n | y_{1:n-1}) = \sum_{i=1}^N W_{n-1}^{(i)} g(y_n | x_n) f(x_n | y_n, X_{n-1}^{(i)}),$$

$$\begin{aligned} \widetilde{\nabla} p_\theta(x_n, y_n | y_{1:n-1}) &= g_\theta(y_n | x_n) \sum_{i=1}^N W_{n-1}^{(i)} f_\theta(x_n | X_{n-1}^{(i)}) \\ &\quad \times \left(\nabla \log g_\theta(y_n | x_n) + \nabla \log f_\theta(x_n | X_{n-1}^{(i)}) + A_{n-1}^{(i)} \right). \end{aligned}$$

- We use a marginalized version of the APF which relies on a joint probability density

$$q_{\theta}(x_n, y_n | x_{n-1}) = q_{\theta}(x_n | y_n, x_{n-1}) q_{\theta}(y_n | x_{n-1})$$

which is an approximation of

$$p_{\theta}(x_n, y_n | x_{n-1}) = g_{\theta}(y_n | x_n) f_{\theta}(x_n | x_{n-1})$$

- We construct the marginal importance distribution

$$q_{\theta}(x_n | y_n) = \sum_{i=1}^N \widetilde{W}_n^{(i)} q_n(x_n | y_n, X_{n-1}^{(i)}),$$

$$\widetilde{W}_n^{(i)} \propto W_{n-1}^{(i)} q_{\theta}(y_n | X_{n-1}^{(i)}).$$

- Sampling from $q_{\theta}(x_n | y_n)$ includes implicitly the resampling step.

SMC for Sensitivity

- Sample $X_n^{(i)} \sim q_\theta(\cdot | y_n)$.
- Evaluate

$$w_n^{(i)} = \frac{\tilde{p}_\theta(X_n^{(i)}, y_n | y_{1:n-1})}{q_\theta(X_n^{(i)} | y_n)}, \quad a_n^{(i)} = \frac{\widehat{\nabla} p_\theta(X_n^{(i)}, y_n | y_{1:n-1})}{q_\theta(X_n^{(i)} | y_n)}$$

$$W_n^{(i)} \propto w_n^{(i)} \text{ with } \sum_{i=1}^N W_n^{(i)} = 1,$$

$$W_n^{(i)} A_n^{(i)} = \frac{a_n^{(i)}}{\sum_{j=1}^N w_n^{(j)}} - W_n^{(i)} \frac{\sum_{j=1}^N a_n^{(j)}}{\sum_{j=1}^N w_n^{(j)}},$$

- We have

$$\widehat{\nabla \log p_\theta}(Y_n | Y_{1:n-1}) = \frac{\sum_{i=1}^N a_n^{(i)}}{\sum_{i=1}^N w_n^{(i)}}.$$

SMC for Recursive Maximum Likelihood

- Sample $X_n^{(i)} \sim q_{\theta_{n-1}}(\cdot | y_n)$.
- Evaluate

$$w_n^{(i)} = \frac{\tilde{p}_{\theta_{n-1}}(X_n^{(i)}, y_n | y_{1:n-1})}{q_{\theta_{n-1}}(X_n^{(i)} | y_n)}, \quad a_n^{(i)} = \frac{\tilde{\nabla} p_{\theta_{n-1}}(X_n^{(i)}, y_n | y_{1:n-1})}{q_{\theta_{n-1}}(X_n^{(i)} | y_n)}$$

$$W_n^{(i)} \propto w_n^{(i)} \text{ with } \sum_{i=1}^N W_n^{(i)} = 1,$$

$$W_n^{(i)} A_n^{(i)} = \frac{a_n^{(i)}}{\sum_{j=1}^N w_n^{(j)}} - W_n^{(i)} \frac{\sum_{j=1}^N a_n^{(j)}}{\sum_{j=1}^N w_n^{(j)}},$$

- Update the parameter

$$\theta_n = \theta_{n-1} + \gamma_n \frac{\sum_{i=1}^N a_n^{(i)}}{\sum_{i=1}^N w_n^{(i)}}.$$

- This algorithm is perhaps not very elegant but simple.
- This algorithm only relies on the SMC approximation of the marginals $p(x_n | y_{1:n})$.
- Under standard mixing assumptions, we can establish uniform convergence results for $\widehat{\nabla} p_\theta(x_n | y_{1:n})$.
- There is no accumulation of errors over time contrary to the SMC approaches discussed earlier.
- It has been used successfully for high-dimensional parameter estimation problems arising in robotics and bioinformatics.
- The observed information matrix can be computed similarly.

Limitations of this approach

- It is in $O(N^2)$ although fast methods can be used to speed it up.
- It requires scaling the step-size sequence appropriately for multidimensional parameters.
- It is only useful for large datasets.

Alternative to Stochastic Gradient

- In a batch context, the EM algorithm is a very popular alternative to gradient-type approaches.
- It is possible to derive an online version of the EM.
- However, once more, this algorithm would rely on an SMC approximation of $p_{\theta}(x_{1:n} | y_{1:n})$.
- A simple fixed-lag approximation can be used to mitigate this problem but not asymptotically consistent (good course project though).

Pseudo-likelihood Approaches

- Instead of trying to maximize the likelihood, we introduce a pseudo-likelihood.
- Assuming a stationary state-space model, we have

$$p_{\theta}(x_k, y_k) = \pi_{\theta}(x_{kL+1}) g_{\theta}(y_{kL+1} | x_{kL+1}) \prod_{i=kL+2}^{(k+1)L} f_{\theta}(x_i | x_{i-1}) g_{\theta}(y_i | x_i) .$$

- The log pseudo-likelihood for m blocks of observations is given by

$$l_L(\theta, Y_{0:m-1}) := \sum_{k=0}^{m-1} \log p_{\theta}(Y_k) , \quad (2)$$

- Compared to the true likelihood, essentially ignores the dependence between data blocks.

- Under ergodicity assumptions, we have

$$\lim_{m \rightarrow \infty} \frac{1}{m} l_L(\theta, Y_{0:m-1}) =: l_L(\theta) ,$$

where

$$l_L(\theta) := \int_{Y^L} \log(p_\theta(y)) p_{\theta^*}(y) dy.$$

- It can be shown that the set of parameters maximizing $l_L(\theta)$ includes the true parameter. This follows from the fact that maximizing $l_L(\theta)$ is equivalent to minimizing the following Kullback-Leibler divergence

$$K_L(\theta, \theta^*) = l_L(\theta^*) - l_L(\theta) \geq 0 .$$

On-line EM algorithm

- To introduce the on-line EM, we first present an “ideal” batch EM algorithm to minimize $K_L(\theta, \theta^*)$ with respect to θ or equivalently to maximize $l_L(\theta)$.
- At iteration $k + 1$, given an estimate θ_k of θ^* , we update our estimate via

$$\theta_{k+1} = \arg \max_{\theta \in \Theta} Q(\theta, \theta_k) ,$$

where

$$Q(\theta, \theta_k) = \int_{\mathcal{X}^L \times \mathcal{Y}^L} \log(p_\theta(x, y)) p_{\theta_k}(x|y) p_{\theta^*}(y) dx dy .$$

- Now for any $\theta \in \Theta$

$$\begin{aligned} Q(\theta_{k+1}, \theta_k) - Q(\theta_k, \theta_k) &= K_L(\theta_k, \theta^*) - K_L(\theta_{k+1}, \theta^*) \\ &\quad + \int_{\mathcal{X}^L \times \mathcal{Y}^L} \log\left(\frac{p_{\theta_{k+1}}(x|y)}{p_{\theta_k}(x|y)}\right) p_{\theta_k}(x|y) p_{\theta^*}(y) dx dy \end{aligned}$$

so an iteration of this “ideal” EM algorithm decreases the value of $K_L(\theta_k, \theta^*)$.

- In practice for the models which we will consider, it is necessary to compute a set of sufficient statistics $\Phi(\theta_k, \theta^*)$ at time k in order to compute Q .
- In practice, $Q(\theta, \theta_{k-1})$ cannot be computed as the expectations appearing in the expression for $\Phi(\theta_k, \theta^*)$ are with respect to a measure dependent on the unknown parameter value θ^* .
- Thanks to the ergodicity and stationarity assumptions, the observations $\{Y_k\}$ provide us with samples from $p_{\theta^*}(y)$ which can be used for the purpose of Monte Carlo integration

$$\hat{\Phi}_k = (1 - \gamma_k) \hat{\Phi}_{k-1} + \gamma_k \mathbb{E}_{\theta_{k-1}}(\Psi(X_k, Y_k) | Y_k), \quad (3)$$

where $\mathbb{E}_{\theta_{k-1}}(\phi(X_k) | Y_k)$ denotes the expectation of ϕ with respect to $p_{\theta_{k-1}}(x_k | Y_k)$.

- We then substitute $\hat{\Phi}_k$ for $\Phi(\theta_k, \theta^*)$ and obtain $\theta_k = \Lambda(\hat{\Phi}_k)$.
- If θ_k was constant and $\gamma_k = k^{-1}$ then $\hat{\Phi}_k$ would simply compute the arithmetic average of $\{\mathbb{E}_{\theta_{k-1}}(\Psi(X_k, Y_k) | Y_k)\}$, and converge towards $\Phi(\theta_k, \theta^*)$ by ergodicity.

- To summarize, the vector of sufficient statistics $\hat{\Phi}_{-1}$ is arbitrarily initialized and the on-line EM algorithm proceeds as follows for the data block indexed by $k \geq 0$.

- E-step

$$\hat{\Phi}_k = (1 - \gamma_k)\hat{\Phi}_{k-1} + \gamma_k \mathbb{E}_{\theta_{k-1}} (\Psi (X_k, Y_k) | Y_k) .$$

- M-step

$$\theta_k = \Lambda(\hat{\Phi}_k) .$$

- In scenarios where $\mathbb{E}_{\theta_k} (\Psi (X_k, Y_k) | Y_k)$ does not admit an analytical expression, a further Monte Carlo approximation can be used.

- Assume that a good approximation $q_{\theta_{k-1}}(x_k | Y_k)$ of $p_{\theta_{k-1}}(x_k | Y_k)$ is available, and that it is easy to sample from $q_{\theta_{k-1}}(x_k | Y_k)$.

- E-step

$$X_k^{(i)} \sim q_{\theta_{k-1}}(\cdot | Y_k) \text{ for } i = 1, \dots, N,$$

$$\hat{\Phi}_k = (1 - \gamma_k) \hat{\Phi}_{k-1} + \gamma_k \sum_{i=1}^N W_k^{(i)} \Psi(X_k^{(i)}, Y_k),$$

where

$$W_k^{(i)} \propto \frac{p_{\theta_{k-1}}(X_k^{(i)}, Y_k)}{q_{\theta_{k-1}}(X_k^{(i)} | Y_k)}, \quad \sum_{i=1}^N W_k^{(i)} = 1.$$

- M-step

$$\theta_k = \Lambda(\hat{\Phi}_k).$$

- If it is possible to sample from $p_{\theta_{k-1}}(x_k | Y_k)$ exactly then it is not necessary to have a large N , $N = 1$ is sufficient. Indeed it is only necessary to produce estimates of $\mathbb{E}_{\theta_{k-1}}(\Psi(X_k, Y_k) | Y_k)$.

- Note that as such the algorithm above leads to asymptotically biased estimates, but that this can be easily corrected by considering instead the following recursion for the estimation of the conditional expectation

$$\hat{F}_k = (1 - \gamma_k)\hat{F}_{k-1} + \gamma_k \frac{1}{N} \sum_{i=1}^N \frac{p_{\theta_{k-1}}(\mathbf{X}_k^{(i)}, \mathbf{Y}_k)}{q_{\theta_{k-1}}(\mathbf{X}_k^{(i)} | \mathbf{Y}_k)} \Psi(\mathbf{X}_k^{(i)}, \mathbf{Y}_k) ,$$

$$\hat{N}_k = (1 - \gamma_k)\hat{N}_{k-1} + \gamma_k \frac{1}{N} \sum_{i=1}^N \frac{p_{\theta_{k-1}}(\mathbf{X}_k^{(i)}, \mathbf{Y}_k)}{q_{\theta_{k-1}}(\mathbf{X}_k^{(i)} | \mathbf{Y}_k)} ,$$

and let $\hat{\Phi}_k = \hat{F}_k / \hat{N}_k$.

- SMC techniques can also be used to approximate this expectation. We stress here on the fact that in the situation where SMC methods are used in this context, the path degeneracy issue is easily dealt with since L is fixed, and very often of small dimension.