Machine Learning - Waseda University SMC

AD

July 2011

• $\{X_k\}_{k\geq 1}$ hidden $\mathcal X$ -valued Markov process with

$$X_{1}\sim\mu\left(x_{1}\right) \text{ and } X_{k}|\left(X_{k-1}=x_{k-1}\right)\sim f\left(\left.x_{k}\right|x_{k-1}\right).$$

• $\{Y_k\}_{k\geq 1}$ observed \mathcal{Y} -valued process with observations conditionally independent given $\{X_k\}_{k\geq 1}$ with

$$Y_k|(X_k=x_k)\sim g(y_k|x_k).$$

• Main Objective: Estimate $\{X_k\}_{k\geq 1}$ given $\{Y_k\}_{k\geq 1}$ online/offline.

Inference in State-Space Models

• Given observations $y_{1:n} := (y_1, y_2, \dots, y_n)$, inference about $X_{1:n} := (X_1, \dots, X_n)$ relies on the posterior

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

where

$$p(x_{1:n}, y_{1:n}) = \underbrace{\mu(x_1) \prod_{k=2}^{n} f(x_k | x_{k-1}) \prod_{k=1}^{n} g(y_k | x_k)}_{p(x_{1:n})}$$

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

- We want to compute $p(x_{1:n}|y_{1:n})$ and $p(y_{1:n})$ sequentially in time n.
- For non-linear non-Gaussian models, numerical approximations are required.

Monte Carlo Methods

• Assume you can generate $X_{1:n}^{(i)} \sim p(x_{1:n}|y_{1:n})$ where i = 1, ..., N then MC approximation is

$$\widehat{p}(x_{1:n}|y_{1:n}) = \frac{1}{N} \sum_{i=1}^{N} \delta_{X_{1:n}^{(i)}}(x_{1:n})$$

• Integration is straightforward

$$\int \varphi_{n}\left(x_{1:n}\right)\widehat{p}\left(x_{1:n}\right|y_{1:n}\right)dx_{1:n}=\frac{1}{N}\sum_{i=1}^{N}\varphi_{n}\left(X_{1:n}^{\left(i\right)}\right).$$

• Marginalisation is straightforward

$$\widehat{p}(x_{k}|y_{1:n}) = \int \widehat{p}(x_{k}|y_{1:n}) dx_{1:k-1} dx_{k+1:n} = \frac{1}{N} \sum_{i=1}^{N} \delta_{X_{k}^{(i)}}(x_{k})$$

• **Problem**: Sampling from $p(x_{1:n}|y_{1:n})$ is impossible in general cases.

- Divide and conquer strategy: Break the problem of sampling from $p(x_{1:n}|y_{1:n})$ into a collection of simpler subproblems. First approximate $p(x_1|y_1)$ at time 1, then $p(x_{1:2}|y_{1:2})$ at time 2 and so on.
- Each target distribution is approximated by a cloud of random samples termed *particles* evolving according to *importance sampling* and *resampling* steps.

Importance Sampling

• Assume you have at time n-1

$$\widehat{\rho}(x_{1:n-1}|y_{1:n-1}) = \frac{1}{N} \sum_{i=1}^{N} \delta_{X_{1:n-1}^{(i)}}(x_{1:n-1}).$$

• By sampling $\widetilde{X}_{n}^{(i)} \sim f\left(x_{n} | X_{n-1}^{(i)}\right)$ and setting $\widetilde{X}_{1:n}^{(i)} = \left(X_{1:n-1}^{(i)}, \widetilde{X}_{n}^{(i)}\right)$ then

$$\widehat{p}(x_{1:n}|y_{1:n-1}) = \frac{1}{N} \sum_{i=1}^{N} \delta_{\widetilde{X}_{1:n}^{(i)}}(x_{1:n}).$$

Our target at time n is

$$p(x_{1:n}|y_{1:n}) = \frac{g(y_n|x_n) p(x_{1:n}|y_{1:n-1})}{\int g(y_n|x_n) p(x_{1:n}|y_{1:n-1}) dx_n}$$

so by substituting $\widehat{p}\left(\left.x_{1:n}\right|\left.y_{1:n-1}\right)$ to $p\left(\left.x_{1:n}\right|\left.y_{1:n-1}\right)$ we obtain

$$\widetilde{p}(x_{1:n}|y_{1:n}) = \sum_{i=1}^{N} W_n^{(i)} \delta_{\widetilde{X}_{1:n}^{(i)}}(x_{1:n}), \ W_n^{(i)} \propto g(y_n|\widetilde{X}_{1:n}^{(i)}).$$

Resampling

• We have a "weighted" approximation $\widetilde{p}(x_{1:n}|y_{1:n})$ of $p(x_{1:n}|y_{1:n})$

$$\widetilde{p}(x_{1:n}|y_{1:n}) = \sum_{i=1}^{N} W_n^{(i)} \delta_{\widetilde{X}_{1:n}^{(i)}}(x_{1:n}).$$

• To obtain N samples $X_{1:n}^{(i)}$ approximately distributed according to $p(x_{1:n}|y_{1:n})$, we just resample

$$X_{1:n}^{(i)} \sim \widetilde{p}\left(x_{1:n} | y_{1:n}\right)$$

to obtain

$$\widehat{p}(x_{1:n}|y_{1:n}) = \frac{1}{N} \sum_{i=1}^{N} \delta_{X_{1:n}^{(i)}}(x_{1:n}).$$

• Particles with high weights are copied multiples times, particles with low weights die.

Bootstrap Filter (Gordon, Salmond & Smith, 1993)

At time n = 1

• Sample $\widetilde{X}_{1}^{\left(i
ight)}\sim\mu\left(x_{1}
ight)$ then

$$\widetilde{p}(x_{1}|y_{1}) = \sum_{i=1}^{N} W_{1}^{(i)} \delta_{\widetilde{X}_{1}^{(i)}}(x_{1}), \quad W_{1}^{(i)} \propto g\left(y_{1}|\widetilde{X}_{1}^{(i)}\right).$$

• Resample $X_1^{(i)} \sim \widetilde{p}(x_1 | y_1)$ to obtain $\widehat{p}(x_1 | y_1) = \frac{1}{N} \sum_{i=1}^N \delta_{X_1^{(i)}}(x_1)$.

At time $n \ge 2$

• Sample
$$\widetilde{X}_n^{(i)} \sim f\left(x_n | X_{n-1}^{(i)}\right)$$
, set $\widetilde{X}_{1:n}^{(i)} = \left(X_{1:n-1}^{(i)}, \widetilde{X}_n^{(i)}\right)$ and

$$\widetilde{p}(x_{1:n}|y_{1:n}) = \sum_{i=1}^{N} W_n^{(i)} \delta_{\widetilde{X}_{1:n}^{(i)}}(x_{1:n}), \ W_n^{(i)} \propto g(y_n|\widetilde{X}_n^{(i)}).$$

• Resample $X_{1:n}^{(i)} \sim \widetilde{p}(x_{1:n}|y_{1:n})$ to obtain $\widehat{p}(x_{1:n}|y_{1:n}) = \frac{1}{N} \sum_{i=1}^{N} \delta_{X_{1:n}^{(i)}}(x_{1:n}).$

SMC Output

• At time *n*, we get

$$\widehat{p}(x_{1:n}|y_{1:n}) = \frac{1}{N} \sum_{i=1}^{N} \delta_{X_{1:n}^{(i)}}(x_{1:n}).$$

The marginal likelihood estimate is given by

$$\widehat{p}(y_{1:n}) = \prod_{k=1}^{n} \widehat{p}(y_k | y_{1:k-1}) = \prod_{k=1}^{n} \left(\frac{1}{N} \sum_{i=1}^{N} g\left(y_k | \widetilde{X}_k^{(i)}\right) \right)$$

- Computational complexity is $\mathcal{O}(N)$ and memory requirements $\mathcal{O}(nN)$.
- If we are only interested in $p(x_n | y_{1:n})$ or $p(s_n(x_{1:n}) | y_{1:n})$ where $s_n(x_{1:n}) = \Psi_n(x_n, s_{n-1}(x_{1:n-1}))$ is fixed-dimensional then memory requirements $\mathcal{O}(N)$.

SMC on Path-Space - figures by Olivier Cappe



AD ()

Machine Learning - Waseda UniversitySMC



Figure: $p(x_1|y_1)$, $p(x_2|y_{1:2})$ and $\widehat{\mathbb{E}}[X_1|y_1]$, $\widehat{\mathbb{E}}[X_2|y_{1:2}]$ (top) and particle approximation of $p(x_{1:2}|y_{1:2})$ (bottom)



Figure: $p(x_k | y_{1:k})$ and $\widehat{\mathbb{E}}[X_k | y_{1:k}]$ for k = 1, 2, 3 (top) and particle approximation of $p(x_{1:3} | y_{1:3})$ (bottom)



Figure: $p(x_k | y_{1:k})$ and $\widehat{\mathbb{E}}[X_k | y_{1:k}]$ for k = 1, ..., 10 (top) and particle approximation of $p(x_{1:10} | y_{1:10})$ (bottom)



Figure: $p(x_k | y_{1:k})$ and $\widehat{\mathbb{E}}[X_k | y_{1:k}]$ for k = 1, ..., 24 (top) and particle approximation of $p(x_{1:24} | y_{1:24})$ (bottom)

Illustration of the Degeneracy Problem

Degeneracy problem. For any N and any k, there exists n(k, N) such that for any n ≥ n(k, N)

$$\widehat{p}\left(\left.x_{1:k}\right|y_{1:n}\right) = \delta_{X_{1:k}^{*}}\left(x_{1:k}\right).$$

 $\widehat{p}(x_{1:n}|y_{1:n})$ is an unreliable approximation of $p(x_{1:n}|y_{1:n})$ as $n \nearrow$.



Figure: Exact calculation of $\frac{1}{n}\mathbb{E}\left[\sum_{k=1}^{n} X_{k} | y_{1:n}\right]$ via Kalman (blue) vs SMC estimate (red) for N = 1000. As *n* increases, the SMC estimate deteriorates.

Convergence Results

- Numerous precise convergence results are available for SMC methods (Del Moral, 2004).
- Let $\varphi_n: \mathcal{X}^n \to \mathbb{R}$ and consider

$$\overline{\varphi}_{n} = \int \varphi_{n}(x_{1:n}) p(x_{1:n}|y_{1:n}) dx_{1:n},$$

$$\widehat{\varphi}_{n} = \int \varphi_{n}(x_{1:n}) \widehat{p}(x_{1:n}|y_{1:n}) dx_{1:n} = \frac{1}{N} \sum_{i=1}^{N} \varphi_{n}\left(X_{1:n}^{(i)}\right)$$

• Under very weak assumptions, we have for any p > 0

$$\mathbb{E}\left[\left|\widehat{\varphi}_{n}-\overline{\varphi}_{n}\right|^{p}\right]^{1/p} \leq \frac{C_{n}}{\sqrt{N}}$$

and

$$\lim_{N\to\infty}\sqrt{N}\left(\widehat{\varphi}_n-\overline{\varphi}_n\right)\Rightarrow\mathcal{N}\left(0,\sigma_n^2\right).$$

Very weak results: C_n and σ²_n can increase with n and will for a path-dependent φ_n (x_{1:n}) as the degeneracy problem suggests!

Stronger Convergence Results

• Exponentially stability assumption. For any x_1, x'_1

$$\frac{1}{2} \int \left| p\left(\left. x_{n} \right| y_{2:n}, X_{1} = x_{1} \right) - p\left(\left. x_{n} \right| y_{2:n}, X_{1} = x_{1}' \right) \right| \, dx_{n} \leq \alpha^{n} \, \, \text{for} \, \left. \left| \alpha \right| < 1.$$

• Marginal distribution. For $\varphi_{n}\left(x_{1:n}\right) = \varphi\left(x_{n}\right)$,

$$\mathbb{E}\left[\left|\widehat{\varphi}_{n}-\overline{\varphi}_{n}\right|^{p}\right]^{1/p} \leq \frac{C}{\sqrt{N}},\\ \lim_{N\to\infty}\sqrt{N}\left(\widehat{\varphi}_{n}-\overline{\varphi}_{n}\right) \Rightarrow \mathcal{N}\left(0,\sigma_{n}^{2}\right) \text{ where } \sigma_{n}^{2} \leq D,$$

where C and D typically exponential in $\dim(X_n)$.

• Marginal likelihood.

 $\lim_{N\to\infty}\sqrt{N}\left(\log\widehat{p}\left(y_{1:n}\right)-\log p\left(y_{1:n}\right)\right)\Rightarrow\mathcal{N}\left(0,\overline{\sigma}_{n}^{2}\right) \text{ with } \overline{\sigma}_{n}^{2}\leq A n.$

• Resampling is necessary. Without resampling, we have

$$\log \hat{p}(y_{1:n}) = \log \frac{1}{N} \sum_{i=1}^{N} \prod_{k=1}^{n} g\left(y_{k} | \widetilde{X}_{k}^{(i)}\right) \text{ which has a variance}$$

increasing exponentially with n even for trivial examples.

Improving the Sampling Step

- Boostrap filter. Very inefficient for vague prior/peaky likelihood; e.g. $p(x_{n-1}|y_{1:n-1}) = \mathcal{N}(x_{n-1}; m, \sigma^2)$, $f(x_n|x_{n-1}) = \mathcal{N}(x_n; x_{n-1}, \sigma_v^2)$ and $g(y_n|x_n) = \mathcal{N}(y_n; x_n, \sigma_w^2)$.
- Optimal proposal/Perfect adaptation. Resample
 - $W_n \propto p(y_n | x_{n-1})$, sample $p(x_n | y_n, x_{n-1}) \propto g(y_n | x_n) f(x_n | x_{n-1})$.



Various standard improvements

• **Approximate optimal proposal**. Design analytical approximation via EKF, UKF $\hat{p}(x_n | y_n, x_{n-1})$ of $p(x_n | y_n, x_{n-1})$. Sample $\hat{p}(x_n | y_n, x_{n-1})$ and set

$$W_n \propto \frac{g\left(y_n \mid x_n\right) f\left(x_n \mid x_{n-1}\right)}{\widehat{p}\left(x_n \mid y_n, x_{n-1}\right)};$$

see also Auxiliary Particle Filters (Pitt & Shephard, 1999)

• Resample Move (Gilks & Berzuini, 1999). After the resampling step, you have $X_{1:n}^{(i)} = X_{1:n}^{(j)}$ for $i \neq j$. To add diversity among particles, use an MCMC kernel $X_{1:n}^{\prime(i)} \sim K_n \left(x_{1:n} | X_{1:n}^{(i)} \right)$ where

$$p(x'_{1:n}|y_{1:n}) = \int p(x_{1:n}|y_{1:n}) K_n(x'_{1:n}|x_{1:n}) dx_{1:n}$$

Here K_n does not have to be ergodic.

Improving the Resampling Step

• Resample N times $X_{1:n}^{(i)} \sim \widetilde{p}(x_{1:n}|y_{1:n}) = \sum_{i=1}^{N} W_n^{(i)} \delta_{\widetilde{X}_{1:n}^{(i)}}(x_{1:n})$ to obtain $\widehat{p}(x_{1:n}|y_{1:n})$ is called *multinomial resampling* as

$$\widehat{p}(x_{1:n}|y_{1:n}) = \frac{1}{N} \sum_{i=1}^{N} \delta_{X_{1:n}^{(i)}}(x_{1:n}) = \sum_{i=1}^{N} \frac{N_n^{(i)}}{N} \delta_{\widetilde{X}_{1:n}^{(i)}}(x_{1:n})$$

where
$$\left\{N_{n}^{(i)}\right\}$$
 follow a multinomial with $\mathbb{E}\left[N_{n}^{(i)}\right] = NW_{n}^{(i)}$, $\mathbb{V}\left[N_{n}^{(1)}\right] = NW_{n}^{(i)}\left(1 - W_{n}^{(i)}\right)$.

• Better resampling steps can be designed with $\mathbb{E}\left[N_n^{(i)}\right] = NW_n^{(i)}$ but smaller $\mathbb{V}\left[N_n^{(i)}\right]$; e.g. stratified resampling (Kitagawa, 1996).

Online Bayesian Parameter Estimation

Assume we have

$$\begin{split} X_n | \left(X_{n-1} = x_{n-1} \right) &\sim f_\theta \left(\left. x_n \right| x_{n-1} \right), \\ Y_n | \left(X_n = x_n \right) &\sim g_\theta \left(\left. y_n \right| x_n \right), \end{split}$$

where θ is an *unknown* static parameter with prior $p(\theta)$. • Given data $y_{1:n}$, inference relies on

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

where

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

• SMC methods apply as it is a standard model with extended state $Z_n = (X_n, \theta_n)$ where

$$f\left(\left.z_{n}\right|z_{n-1}\right) = \underbrace{\delta_{\theta_{n-1}}\left(\theta_{n}\right)}_{\text{practical problems}} f_{\theta_{n}}\left(\left.x_{n}\right|x_{n-1}\right), \ g\left(\left.y_{n}\right|z_{n}\right) = g_{\theta}\left(\left.y_{n}\right|x_{n}\right).$$

- For fixed θ , $\mathbb{V}[\log \hat{p}_{\theta}(y_{1:n})]$ is in Cn/N. In a Bayesian context, the problem is even more severe as $p(\theta|y_{1:n}) \propto p_{\theta}(y_{1:n}) p(\theta)$. Exponential stability assumption cannot hold as $\theta_n = \theta_1$.
- To mitigate but NOT solve the problem, introduce MCMC steps on θ; e.g. (Andrieu, D.&D.,1999; Fearnhead, 1998, 2002; Gilks & Berzuini 1999,2001,2003; Storvik, 2002; Polson & Johannes, 2007; Vercauteren et al., 2005).
- When $p(\theta|y_{1:n}, x_{1:n}) = p(\theta|s_n(x_{1:n}, y_{1:n}))$ where $s_n(x_{1:n}, y_{1:n})$ is fixed-dimensional, this is an elegant algorithm but still relies on $\hat{p}(x_{1:n}|y_{1:n})$ so degeneracy will creep in.
- As dim (Z_n) = dim (X_n) + dim (θ), such methods are not recommended for high-dimensional θ, especially with vague priors.

Example of SMC with MCMC for Parameter Estimation

• Given at time n-1, the approximation at time n

$$\widehat{p}(\theta, x_{1:n-1}|y_{1:n-1}) = \frac{1}{N} \sum_{i=1}^{N} \delta_{\left(\theta_{n-1}^{(i)}, X_{1:n-1}^{(i)}\right)}(\theta, x_{1:n-1}).$$

• Sample
$$\widetilde{X}_n^{(i)} \sim f_{\theta_{n-1}^{(i)}}\left(x_n | X_{n-1}^{(i)}\right)$$
, set $\widetilde{X}_{1:n}^{(i)} = \left(X_{1:n-1}^{(i)}, \widetilde{X}_n^{(i)}\right)$ and

$$\widetilde{p}(\theta, x_{1:n}|y_{1:n}) = \sum_{i=1}^{N} W_n^{(i)} \delta_{\left(\theta_{n-1}^{(i)}, \widetilde{X}_{1:n}^{(i)}\right)}(x_{1:n}), \ W_n^{(i)} \propto g_{\theta_{n-1}^{(i)}}\left(y_n | \widetilde{X}_n^{(i)}\right).$$

• Resample $X_{1:n}^{(i)} \sim \widetilde{p}(x_{1:n}|y_{1:n})$ then sample $\theta_n^{(i)} \sim p\left(\theta|y_{1:n}, X_{1:n}^{(i)}\right)$ to obtain $\widehat{p}(\theta, x_{1:n}|y_{1:n}) = \frac{1}{N} \sum_{i=1}^N \delta_{\left(\theta_n^{(i)}, X_{1:n}^{(i)}\right)}(\theta, x_{1:n}).$

Illustration of the Degeneracy Problem



SMC estimate of $\mathbb{E} \left[\theta | y_{1:n} \right]$, as *n* increases the degeneracy creeps in.

AD ()

Machine Learning - Waseda UniversitySMC

• Given data $y_{1:n}$, inference relies on

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

where

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

- For a given θ , SMC can estimate both $p_{\theta}(x_{1:n}|y_{1:n})$ and $p_{\theta}(y_{1:n})$.
- Is it possible to use SMC within MCMC to sample from p (θ, x_{1:n}| y_{1:n})?

Metropolis-Hastings (MH) Sampler

• To sample from a target $\pi(z)$, the MH sampler generates a Markov chain $\{Z^{(i)}\}$ according to the following mechanism. Given $Z^{(i-1)}$, propose a candidate $Z^* \sim q\left(z^* | Z^{(i-1)}\right)$ and with probability

$$\alpha\left(Z^{(i-1)}, Z^*\right) = \min\left(1, \frac{\pi\left(Z^*\right)q\left(Z^{(i-1)}\right|Z^*\right)}{\pi\left(Z^{(i-1)}\right)q\left(Z^*\right|Z^{(i-1)}\right)}\right)$$

set $Z^{(i)} = Z^*$, otherwise $Z^{(i)} = Z^{(i-1)}$.

• It can be easily shown that

$$\pi(z') = \int \pi(z) K(z'|z) dz$$

where K(z'|z) is the transition kernel of the MH and under weak assumptions $Z^{(i)} \sim \pi(z)$ as $i \to \infty$.

Marginal Metropolis-Hastings Sampler

• Consider the following so-called marginal MH algorithm which target

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

using the proposal

$$q\left(\left(x_{1:n}^{*},\theta^{*}\right)|\left(x_{1:n},\theta\right)\right) = q\left(\theta^{*}|\theta\right)p_{\theta^{*}}\left(x_{1:n}^{*}|y_{1:n}\right).$$

• The MH acceptance probability is

$$\min \left(1, \frac{p(\theta^*, x_{1:n}^* | y_{1:n})}{p(\theta, x_{1:n} | y_{1:n})} \frac{q((x_{1:n}, \theta) | (x_{1:n}^*, \theta^*))}{q((x_{1:n}^*, \theta^*) | (x_{1:n}, \theta))} \right)$$

= min $\left(1, \frac{p_{\theta^*}(y_{1:n}) p(\theta^*)}{p_{\theta}(y_{1:n}) p(\theta)} \frac{q(\theta | \theta^*)}{q(\theta^* | \theta)} \right)$

• **Problem**: We do not know $p_{\theta}(y_{1:n})$ analytically and cannot sample from $p_{\theta}(x_{1:n}|y_{1:n})$ so this algorithm cannot be implemented.

• "Idea": Use SMC approximations of $p_{\theta}(x_{1:n}|y_{1:n})$ and $p_{\theta}(y_{1:n})$.

- At iteration *i*, given $\{\theta(i-1), X_{1:n}(i-1), \widehat{p}_{\theta^{(i-1)}}(y_{1:n})\}$ then sample $\theta^* \sim q(\theta | \theta(i-1))$, run an SMC algorithm to obtain $\widehat{p}_{\theta^*}(x_{1:n} | y_{1:n})$ and $\widehat{p}_{\theta^*}(y_{1:n})$.
- Sample $X_{1:n}^* \sim \widehat{p}_{\theta^*}(x_{1:n} | y_{1:n})$.
- With probability

$$\min\left(1,\frac{\widehat{p}_{\theta^{*}}\left(y_{1:n}\right)p\left(\theta^{*}\right)}{\widehat{p}_{\theta\left(i-1\right)}\left(y_{1:n}\right)p\left(\theta\left(i-1\right)\right)}\frac{q\left(\theta\left(i-1\right)|\theta^{*}\right)}{q\left(\theta^{*}|\theta\left(i-1\right)\right)}\right)$$

 $\begin{array}{l} \mathsf{set} \ \{\theta \left(i\right), X_{1:n} \left(i\right), \widehat{p}_{\theta^{(i)}} \left(y_{1:n}\right)\} = \{\theta^*, X_{1:n}^*, \widehat{p}_{\theta^*} \left(y_{1:n}\right)\} \ \mathsf{otherwise} \ \mathsf{set} \\ \{\theta \left(i\right), X_{1:n} \left(i\right), \widehat{p}_{\theta^{(i)}} \left(y_{1:n}\right)\} = \{\theta \left(i-1\right), X_{1:n} \left(i-1\right), \widehat{p}_{\theta^{(i-1)}} \left(y_{1:n}\right)\}. \end{array}$

Validity of the Particle Marginal MH Sampler

- This algorithm (without sampling $X_{1:n}$) was proposed as an approximate MCMC algorithm to sample from $p(\theta|y_{1:n})$ in (Fernandez-Villaverde & Rubio-Ramirez, 2007).
- Whatever being N ≥ 1, this algorithm admits exactly p (θ, x_{1:n}| y_{1:n}) as invariant distribution (Andrieu, D. & Holenstein, 2010). A particle version of the Gibbs sampler also exists.
- The higher *N*, the better the performance of the algorithm: *N* scales roughly linearly with *n*.
- Particularly useful in scenarios where X_n moderate dimensional & θ high dimensional. Admits the plug and play property (lonides et al., 2006).

Inference for Stochastic Kinetic Models

• Two species
$$X_t^1$$
 (prey) and X_t^2 (predator)

$$\begin{array}{l} \Pr\left(X_{t+dt}^{1} \!=\! x_{t}^{1} \!+\! 1, X_{t+dt}^{2} \!=\! x_{t}^{2} \,\middle|\, x_{t}^{1}, x_{t}^{2}\right) = \alpha \, x_{t}^{1} dt + o \left(dt\right), \\ \Pr\left(X_{t+dt}^{1} \!=\! x_{t}^{1} \!-\! 1, X_{t+dt}^{2} \!=\! x_{t}^{2} \!+\! 1 \!\middle|\, x_{t}^{1}, x_{t}^{2}\right) = \beta \, x_{t}^{1} \, x_{t}^{2} dt + o \left(dt\right), \\ \Pr\left(X_{t+dt}^{1} \!=\! x_{t}^{1}, X_{t+dt}^{2} \!=\! x_{t}^{2} \!-\! 1 \!\middle|\, x_{t}^{1}, x_{t}^{2}\right) = \gamma \, x_{t}^{2} dt + o \left(dt\right), \end{array}$$

observed at discrete times

$$Y_n = X_{n\Delta}^1 + W_n$$
 with $W_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, \sigma^2\right)$.

 We are interested in the kinetic rate constants θ = (α, β, γ) a priori distributed as (Boys et al., 2008)

$$\alpha \sim \mathcal{G}(1, 10), \quad \beta \sim \mathcal{G}(1, 0.25), \quad \gamma \sim \mathcal{G}(1, 7.5).$$

• MCMC methods require reversible jumps, PMMH requires only forward simulation.

Experimental Results



SMC Fixed-Lag Smoothing Approximation

- Direct SMC approximations of $p(x_{1:n}|y_{1:n})$ and its marginals $p(x_k|y_{1:n})$ gets poorer as $n \nearrow$.
- The fixed-lag smoothing approximation (Kitagawa & Sato, 2001) relies on

 $p(x_{1:k}|y_{1:n}) \approx p(x_{1:k}|y_{1:k+\Delta})$ for Δ large enough.

- Algorithmically: stop resampling $\left\{X_k^{(i)}\right\}$ beyond time $k + \Delta$.
- Computational cost is $\mathcal{O}(Nn)$ but non-vanishing bias as $N \to \infty$ (Olsson & al., 2006).
- Picking Δ is difficult. Δ too small results in p (x_{1:k} | y_{1:k+Δ}) being a poor approximation of p (x_{1:k} | y_{1:n}). Δ too large improves the approximation but particle degeneracy creeps in.

SMC Forward Filtering Backward Smoothing

• Forward filtering Backward smoothing (FFBS).



- SMC Implementation: For k = 1, ..., n, compute $\hat{p}(x_k | y_{1:k})$. For k = n 1, ..., 1, compute $\hat{p}(x_k | y_{1:n}) = \sum_{i=1}^{N} W_{k|n}^{(i)} \delta_{X_k^{(i)}}(x_k)$ with cost $\mathcal{O}(N^2 n)$ using $W_{k|n}^{(i)} = \sum_{j=1}^{N} W_{k+1|n}^{(j)} \frac{f(X_{k+1}^{(j)} | X_k^{(j)})}{\sum_{l=1}^{N} f(X_{k+1}^{(j)} | X_k^{(l)})}.$
- For $\varphi_n(x_{1:n}) = \sum_{k=1}^{n-1} s_k(x_k, x_{k+1})$, the SMC FFBS estimates $\{\widehat{\varphi}_n\}$ can be computed online exactly (Del Moral, D. & Singh, 2009).
- Sampling from $\hat{p}(x_{1:n}|y_{1:n})$ costs O(Nn) (Godsill, D. & West, 2004) but O(n) through rejection sampling (Douc et al., 2009).

SMC Generalized Two-Filter Smoothing

• Generalized Two-Filter smoothing (TFS)



$$\overline{p}(x_{k+1}|y_{k+1:n}) \propto p(y_{k+1:n}|x_{k+1})\overline{p}(x_{k+1}).$$

SMC Implementation: For k = 1, ..., n, compute p̂ (x_k | y_{1:k}). For k = n, ..., 1, compute p̂ (x_{k+1} | y_{k+1:n}). Combine the forward and backward filters to obtain

$$\widehat{p}(x_{k}, x_{k+1}|y_{1:n}) \propto \widehat{p}(x_{k}|y_{1:k}) \frac{f(x_{k+1}|x_{k})}{\overline{p}(x_{k+1})} \widehat{\overline{p}}(x_{k+1}|y_{k+1:n})$$

Cost \$\mathcal{O}(N^2n)\$ but \$\mathcal{O}(Nn)\$ through rejection sampling (Briers, D. & Maskell, 2008) and importance sampling (Fearnhead, Wyncoll & Tawn, 2008; Briers, D. & Singh, 2005).

• Exponentially stability assumption. For any x_1 , x'_1

$$\frac{1}{2} \int \left| p\left(\left| x_{n} \right| y_{2:n}, X_{1} = x_{1} \right) - p\left(\left| x_{n} \right| y_{2:n}, X_{1} = x_{1}' \right) \right| dx_{n} \leq \alpha^{n} \text{ for } |\alpha| < 1.$$

• Additive functionals. If $\varphi_n(x_{1:n}) = \sum_{k=1}^n \varphi(x_k)$, we have for the standard path-based SMC estimate (Poyiadjis, D. & Singh, 2009)

$$\lim_{N\to\infty}\sqrt{N}\left(\widehat{\varphi}_n-\overline{\varphi}_n\right)\Rightarrow \mathcal{N}\left(0,\sigma_n^2\right) \text{ where } \underline{A}n^2\leq \sigma_n^2\leq \overline{A}n^2.$$

For the FFBS and TFS estimates (Douc et al., 2009; Del Moral, D. & Singh, 2009), we have

$$\lim_{N\to\infty}\sqrt{N}\left(\widehat{\varphi}_{n}-\overline{\varphi}_{n}\right)\Rightarrow\mathcal{N}\left(0,\sigma_{n}^{2}\right) \text{ where }\sigma_{n}^{2}\leq Cn$$

• Tradeoff between computational and statistical efficiency.

• Consider a linear Gaussian model

$$\begin{split} X_1 &\sim \mathcal{N}\left(0, \frac{\sigma^2}{1 - \phi^2}\right) \text{ and } X_k = \phi X_{k-1} + \sigma_V V_k, \ V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, 1\right) \\ Y_k &= c X_k + \sigma_W W_k, \ W_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, 1\right). \end{split}$$

- We simulate 10,000 observations for $\theta = (\phi, \sigma_V, c, \sigma_W) = (0.8, 0.5, 1.0, 1.0).$
- We compute the score vector using Fisher's identity

$$\nabla \log p_{\theta}\left(y_{1:n}\right) = \int \nabla \log p_{\theta}\left(x_{1:n}, y_{1:n}\right) p_{\theta}\left(x_{1:n} \middle| y_{1:n}\right) dx_{1:n}$$

at the true value of θ and compare to its true value.

Empirical Variance for Standard vs FFBS Approximations



Standard path-based (left) vs FFBS (right); the vertical scale is different

Parameter Estimation using Gradient Ascent/EM

• Gradient ascent: To maximise $p_{\theta}(y_{1:n})$ w.r.t θ , use at iteration k+1

$$heta_{k+1} = heta_k +
abla \log p_{ heta} \left(y_{1:n}
ight) |_{ heta = heta_k}$$

where $\nabla \log p_{\theta}(y_{1:n})|_{\theta=\theta_k}$ is computed using Fisher's identity or IPA (Coquelin, Deguest & Munos, 2009) and any SMC smoothing algorithm.

• *EM algorithm*: To maximise $p_{\theta}(y_{1:n})$ w.r.t θ , the EM uses at iteration k+1

$$\theta_{k+1} = \arg \max \ Q(\theta_k, \theta).$$

where

$$Q(heta_k, heta) = \int \log p_{ heta}(x_{1:n}, y_{1:n}) \ p_{ heta_k}(x_{1:n}|y_{1:n}) dx_{1:n}$$

can be computed using any SMC smoothing algorithm.

Online Parameter Estimation using Gradient Ascent/EM

 In the online implementation (Le Gland & Mevel, 1997), update the parameter at time n+1 using

$$\theta_{n+1} = \theta_n + \gamma_{n+1} \nabla \log p_{\theta_{1:n}}(y_n | y_{1:n-1})$$

where $\sum_n \gamma_n = \infty$, $\sum_n \gamma_n^2 < \infty$ and

 $\nabla \log p_{\theta_{1:n}}(y_n | y_{1:n-1}) = \nabla \log p_{\theta_{1:n}}(y_{1:n}) - \nabla \log p_{\theta_{1:n-1}}(y_{1:n-1}).$

- An estimate of ∇ log p_{θ1:n}(y_n|y_{1:n-1}) with a time-uniform bounded variance can be computed using online SMC FFBS estimate (Del Moral, D. & Singh, 2009).
- A numerically stable SMC implementation of online EM (e.g. Cappé, 2009; Elliott, Ford & Moore, 2002) can also be implemented using online SMC FFBS estimate.
- These non-Bayesian procedures do not suffer from the degeneracy problem but require long data sets for convergence.