

Maximum Likelihood Parameter Estimation in State-Space Models

Arnaud Doucet
Department of Statistics, Oxford University

University College London

4th October 2012

- Let $\{X_t\}_{t \geq 1}$ be a latent/hidden \mathcal{X} -valued Markov process with

$$X_1 \sim \mu(\cdot) \text{ and } X_t | (X_{t-1} = x) \sim f(\cdot | x).$$

- Let $\{Y_t\}_{t \geq 1}$ be an \mathcal{Y} -valued Markov observation process such that

$$Y_t | (X_t = x) \sim g(\cdot | x).$$

- Particle filters estimate $\{p(x_{1:t} | y_{1:t})\}_{t \geq 1}$ on-line but only estimates of $\{p(x_t | y_{1:t})\}_{t \geq 1}$ and $\{p(y_{1:t})\}_{t \geq 1}$ are reliable.
- Particle smoothing methods allow us to obtain reliable estimates of $\{p(x_t | y_{1:T})\}_{t=1}^T$.

State-Space Models with Unknown Parameters

- In most scenarios of interest, the state-space model contains an unknown static parameter $\theta \in \Theta$ so that

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

- The observations $\{Y_t\}_{t \geq 1}$ are conditionally independent given $\{X_t\}_{t \geq 1}$ and θ

$$Y_t | (X_t = x_t) \sim g_\theta(\cdot | x).$$

- **Aim:** We would like to infer θ either on-line or off-line.

- **Stochastic Volatility model**

$$X_t = \phi X_{t-1} + \sigma V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

$$Y_t = \beta \exp(X_t/2) W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

where $\theta = (\phi, \sigma^2, \beta)$.

- **Biochemical Network model**

$$\Pr(X_{t+dt}^1 = x_t^1 + 1, X_{t+dt}^2 = x_t^2 \mid x_t^1, x_t^2) = \alpha x_t^1 dt + o(dt),$$

$$\Pr(X_{t+dt}^1 = x_t^1 - 1, X_{t+dt}^2 = x_t^2 + 1 \mid x_t^1, x_t^2) = \beta x_t^1 x_t^2 dt + o(dt),$$

$$\Pr(X_{t+dt}^1 = x_t^1, X_{t+dt}^2 = x_t^2 - 1 \mid x_t^1, x_t^2) = \gamma x_t^2 dt + o(dt),$$

with

$$Y_k = X_{k\Delta T}^1 + W_k \text{ with } W_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$$

where $\theta = (\alpha, \beta, \gamma)$.

Parameter Inference in State-Space Models

- Online Bayesian parameter inference.
- Offline Maximum Likelihood parameter inference.
- Online Maximum Likelihood parameter inference.

Bayesian Parameter Inference in State-Space Models

- Set a prior $p(\theta)$ on θ so inference relies now on

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

where

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

with

$$p_{\theta}(x_{1:t}, y_{1:t}) = \mu_{\theta}(x_1) \prod_{k=2}^t f_{\theta}(x_k | x_{k-1}) \prod_{k=1}^t g_{\theta}(y_k | x_k)$$

- We have

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

- Standard and more sophisticated particle methods to sample from $\{p(\theta, x_{1:t} | y_{1:t})\}_{t \geq 1}$ are ALL unreliable.

At time $t = 1$

- $(\bar{\theta}_1^{(i)}, \bar{X}_1^{(i)}) \sim p(\theta) \mu_\theta(x_1)$ then
$$\bar{p}(\theta, x_1 | y_1) = \sum_{i=1}^N W_1^{(i)} \delta_{(\bar{\theta}_1^{(i)}, \bar{X}_1^{(i)})}(\theta, x_1), \quad W_1^{(i)} \propto g_{\bar{\theta}_1^{(i)}}(y_1 | \bar{X}_1^{(i)}).$$
- $(\theta_1^{(i)}, X_1^{(i)}) \sim \bar{p}(\theta, x_1 | y_1)$
and $\hat{p}(\theta, x_1 | y_1) = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_1^{(i)}, X_1^{(i)}}(\theta, x_1).$

At time $t \geq 2$

- Set $\bar{\theta}_t^{(i)} = \theta_{t-1}^{(i)}, \bar{X}_t^{(i)} \sim f_{\bar{\theta}_t^{(i)}}(x_t | X_{t-1}^{(i)})$ and $\bar{X}_{1:t}^{(i)} = (X_{1:t-1}^{(i)}, \bar{X}_t^{(i)})$
$$\bar{p}(\theta, x_{1:t} | y_1) = \sum_{i=1}^N W_t^{(i)} \delta_{(\bar{\theta}_t^{(i)}, \bar{X}_{1:t}^{(i)})}(\theta, x_{1:t}), \quad W_t^{(i)} \propto g_{\bar{\theta}_t^{(i)}}(y_t | \bar{X}_t^{(i)}).$$
- $(\theta_t^{(i)}, X_{1:t}^{(i)}) \sim \bar{p}(\theta, x_{1:t} | y_{1:t})$ and
$$\hat{p}(\theta, x_{1:t} | y_{1:t}) = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_t^{(i)}, X_{1:t}^{(i)}}(\theta, x_{1:t}).$$

Online Bayesian Parameter Inference

- Provide consistent estimates but remarkably inefficient (Chopin, 2002). Particles $\{\theta_1^{(i)}\}$ in Θ space only sampled at time 1: degeneracy problem!
- Consider the extended state $Z_t = (X_t, \theta_t)$ then

$$\begin{aligned}v(z_1) &= p(\theta_1) \mu_{\theta_1}(x_1), \\f(z_t | z_{t-1}) &= \delta_{\theta_{t-1}}(\theta_t) f_{\theta_t}(x_t | x_{t-1}), \quad g(y_t | z_t) = g_{\theta_t}(y_t | x_t); \end{aligned}$$

i.e. $\theta_t = \theta_1$ for any t with θ_1 from the prior. Exponential stability assumption on $\{p(z_t | y_{1:t})\}_{t \geq 1}$ cannot be satisfied.

- Use MCMC steps on θ so as to jitter $\{\theta_t^{(i)}\}$; e.g. Andrieu, De Freitas & D. (1999); Fearnhead (2002); Gilks & Berzuini (2001); Carvalho et al. (2010).
- When $p(\theta | y_{1:t}, x_{1:t}) = p(\theta | s_t(x_{1:t}, y_{1:t}))$ where $s_t(x_{1:t}, y_{1:t})$ is a fixed-dimensional vector, “elegant” but still implicitly relies on $\hat{p}(x_{1:t} | y_{1:t})$ so degeneracy will creep in.

- At time $t - 1$, we have

$$\hat{p}(\theta, x_{1:t-1} | y_{1:t-1}) = \frac{1}{N} \sum_{i=1}^N \delta_{(\theta_{t-1}^{(i)}, X_{1:t-1}^{(i)})}(\theta, x_{1:t-1}),$$

- Set $\bar{\theta}_t^{(i)} = \theta_{t-1}^{(i)}$, sample $\bar{X}_t^{(i)} \sim f_{\bar{\theta}_t^{(i)}}(\cdot | X_{t-1}^{(i)})$, set $\bar{X}_{1:t}^{(i)} = (X_{1:t-1}^{(i)}, \bar{X}_t^{(i)})$ and

$$\begin{aligned} \bar{p}(\theta, x_{1:t} | y_{1:t}) &= \sum_{i=1}^N W_t^{(i)} \delta_{(\bar{\theta}_t^{(i)}, \bar{X}_{1:t}^{(i)})}(\theta, x_{1:t}), \\ W_t^{(i)} &\propto g_{\bar{\theta}_t^{(i)}}(y_t | \bar{X}_t^{(i)}). \end{aligned}$$

- Resample $(\theta_t^{(i)}, X_{1:t}^{(i)}) \sim \bar{p}(\theta, x_{1:t} | y_{1:t})$ then sample $\theta_t^{(i)} \sim p(\theta | y_{1:t}, X_{1:t}^{(i)})$ to obtain $\hat{p}(\theta, x_{1:t} | y_{1:t}) = \frac{1}{N} \sum_{i=1}^N \delta_{(\theta_t^{(i)}, X_{1:t}^{(i)})}(\theta, x_{1:t})$.

A Toy Example

- Linear Gaussian state-space model

$$X_t = \theta X_{t-1} + \sigma_V V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

$$Y_t = X_t + \sigma_W W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1).$$

- We set $p(\theta) \propto \mathbf{1}_{(-1,1)}(\theta)$ so

$$p(\theta | y_{1:t}, x_{1:t}) \propto \mathcal{N}(\theta; m_t, \sigma_t^2) \mathbf{1}_{(-1,1)}(\theta)$$

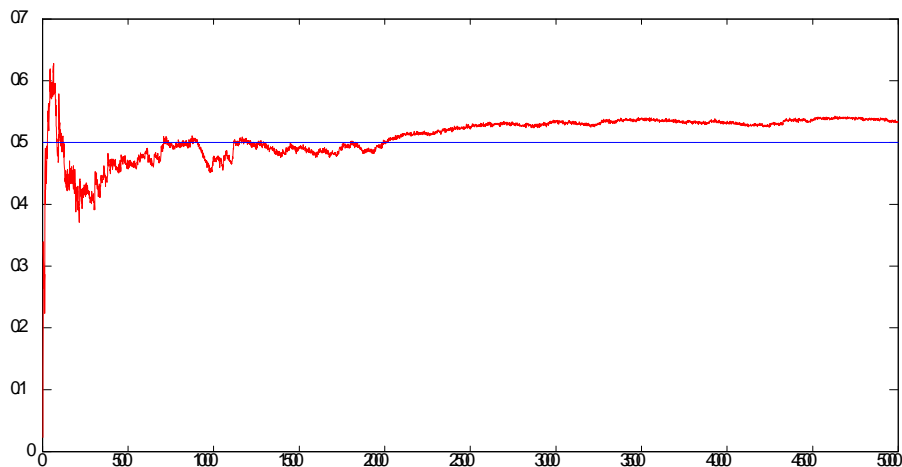
where

$$\sigma_t^2 = S_{2,t}^{-1}, \quad m_t = S_{2,t}^{-1} S_{1,t}$$

with

$$S_{1,t} = \sum_{k=2}^t x_{k-1} x_k, \quad S_{2,t} = \sum_{k=2}^t x_{k-1}^2$$

Illustration of the Degeneracy Problem



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

Another Toy Example

- Linear Gaussian state-space model

$$X_t = \rho X_{t-1} + V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

$$Y_t = X_t + \sigma W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1).$$

- We set $\rho \sim \mathcal{U}_{(-1,1)}$ and $\sigma^2 \sim \mathcal{IG}(1, 1)$.
- We use particle filter with perfect adaptation and Gibbs moves with $N = 10000$; particle learning (Andrieu, D. & De Freitas, 1999; Carvalho et al., 2010)
- 50 runs of the particle method vs ground truth obtained using Kalman filter on states and grid on parameters.

Another Illustration of Degeneracy for Particle Learning

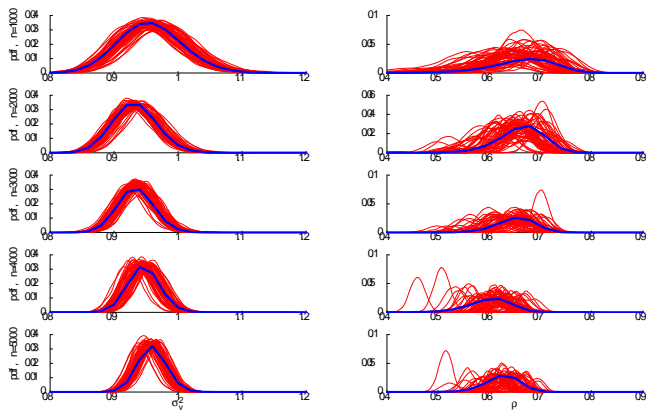


Figure: Estimates of $p(\rho | y_{1:t})$ and $p(\sigma^2 | y_{1:t})$ over 50 runs (red) vs ground truth (blue) for $t = 10^3, 2.10^3, \dots, 5.10^3$ for $N = 10^4$ (Kantas et al., 2012)

- For fixed θ , $\mathbb{V} [\widehat{p}_\theta (y_{1:t}) / p_\theta (y_{1:t})]$ is in $\mathcal{O} (t/N)$.
- In a Bayesian context, $p (\theta | y_{1:t}) \propto p_\theta (y_{1:t}) p (\theta)$ so we implicitly need to compute $p_\theta (y_{1:t})$ at each particle location $\theta^{(i)}$.
- It appears impossible to obtain uniformly in time stable estimates of $\{p (\theta | y_{1:t})\}_{t \geq 1}$ for a fixed N .
- However for a given time horizon T , we can use PF to sample efficiently from $p (\theta | y_{1:T})$; see Lecture 3.

Likelihood Function Estimation

- Let $y_{1:T}$ being given, the log-(marginal) likelihood is given by

$$\ell(\theta) = \log p_{\theta}(y_{1:T}).$$

- For any $\theta \in \Theta$, one can estimate $\ell(\theta)$ using particle methods, variance $\mathcal{O}(T/N)$.
- Direct maximization of $\ell(\theta)$ difficult as estimate $\hat{\ell}(\theta)$ is not a smooth function of θ even for fixed random seed.
- For $\dim(X_t) = 1$, we can obtain smooth estimate of log-likelihood function by using a smoothed resampling step (e.g. Pitt, 2011); i.e. piecewise linear approximation of $\Pr(X_t < x | y_{1:t})$.
- For $\dim(X_t) > 1$, we can obtain estimates of $\ell(\theta)$ highly positively correlated for neighbouring values in Θ (e.g. Lee, 2008).

- To maximise $\ell(\theta)$ w.r.t θ , use at iteration $k + 1$

$$\theta_{k+1} = \theta_k + \gamma_k \nabla \ell(\theta)|_{\theta=\theta_k}$$

where $\nabla \ell(\theta)|_{\theta=\theta_k}$ is the so-called score vector.

- $\nabla \ell(\theta)|_{\theta=\theta_k}$ can be estimated using finite differences but more efficiently using Fisher's identity

$$\nabla \ell(\theta) = \int \nabla \log p_{\theta}(x_{1:T}, y_{1:T}) p_{\theta}(x_{1:T} | y_{1:T}) dx_{1:T}$$

where

$$\begin{aligned} \nabla \log p_{\theta}(x_{1:T}, y_{1:T}) &= \nabla \log \mu_{\theta}(x_1) \\ &+ \sum_{t=2}^T \nabla \log f_{\theta}(x_t | x_{t-1}) + \sum_{t=1}^T \nabla \log g_{\theta}(y_t | x_t). \end{aligned}$$

Particle Calculation of the Score Vector

- We have

$$\begin{aligned} \nabla \ell(\theta) &= \int \{ \nabla \log \mu_{\theta}(x_1) + \nabla \log g_{\theta}(y_1 | x_1) \} p_{\theta}(x_1 | y_{1:T}) dx_1 \\ + \sum_{t=2}^T \int \{ \nabla \log f_{\theta}(x_t | x_{t-1}) + \nabla \log g_{\theta}(y_t | x_t) \} p_{\theta}(x_{t-1}, x_t | y_{1:T}) dx_{t-1} dx_t \end{aligned}$$

- To approximate $\nabla \ell(\theta)$, we just need particle approximations of $\{p_{\theta}(x_{t-1}, x_t | y_{1:T})\}_{t=2}^T$.
- All the particle smoothing methods detailed before can be applied.
- Similar “smoothed additive functionals” have to be computed when implementing the Expectation-Maximization.

Comparison Direct Method vs FB

- We want to estimate

$$\bar{\varphi}_T = \sum_{t=1}^T \int \varphi(x_{t-1}, x_t, y_t) p(x_{t-1}, x_t | y_{1:T}) dx_{t-1} dx_t.$$

Method	Direct	FB
# particles	N	N
cost	$\mathcal{O}(TN)$	$\mathcal{O}(TN^2), \mathcal{O}(TN)$
Var.	$\mathcal{O}(T^2/N)$	$\mathcal{O}(T/N)$
Bias	$\mathcal{O}(T/N)$	$\mathcal{O}(T/N)$
MSE=Bias ² +Var	$\mathcal{O}(T^2/N)$	$\mathcal{O}(T^2/N^2)$

- “Fast” implementations FB of computational complexity $\mathcal{O}(NT)$ outperform direct approach as MSE is $\mathcal{O}(T^2/N^2)$ whereas it is $\mathcal{O}(T^2/N)$ for direct SMC.
- “Naive” implementations FB and TF have MSE of same order as direct method for fixed computational complexity but MSE is bias dominated for FB/TF whereas it is variance dominated for Direct SMC.

Experimental Results

- Consider a linear Gaussian model

$$X_t = \phi X_{t-1} + \sigma_v V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

$$Y_t = cX_t + \sigma_w W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1).$$

- We simulate 10,000 observations and compute particle estimates of

$$\int \varphi_T(x_{1:T}) p(x_{1:T} | y_{1:T}) dx_{1:T}$$

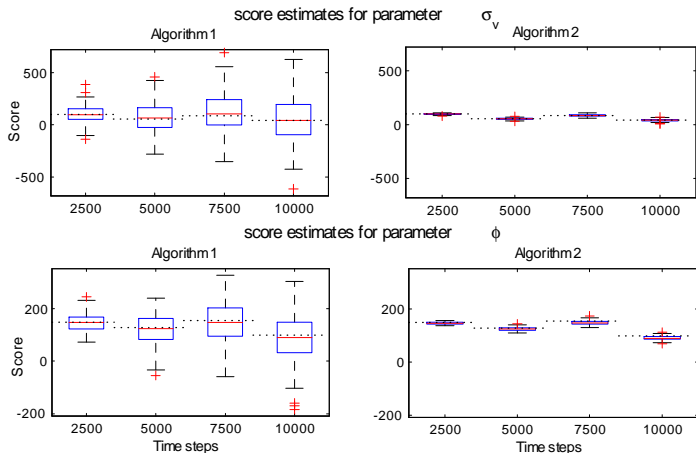
for 4 different additive functionals

$\varphi_t(x_{1:t}) = \varphi_{t-1}(x_{1:t-1}) + \varphi(x_{t-1}, x_t, y_t)$ including

$\varphi^1(x_{t-1}, x_t, y_t) = x_{t-1}x_t$, $\varphi^2(x_{t-1}, x_t, y_t) = x_t^2$. [Ground truth can be computed using Kalman smoother.]

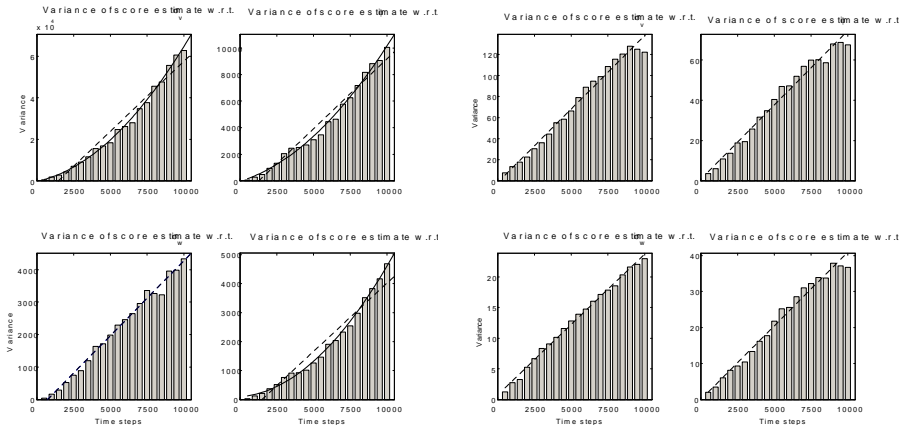
- We use 100 replications on the same dataset to estimate the empirical variance.

Boxplots of Direct vs FB Estimates



Direct (left) vs FB (right)

Empirical Variance for Direct vs FB Estimates



Direct (left) vs FB (right); the vertical scale is different

- *Recursive maximum likelihood* (Titterton, 1984; LeGland & Mevel, 1997) proceeds as follows

$$\theta_{t+1} = \theta_t + \gamma_t \nabla \log p_{\theta_{1:t}}(y_t | y_{1:t-1})$$

where $p_{\theta_{1:t}}(y_t | y_{1:t-1})$ is computed using θ_k at time k and $\sum_t \gamma_t = \infty$, $\sum_t \gamma_t^2 < \infty$. Under regularity conditions, this converges towards a local maximum of the (average) log-likelihood.

- Note that

$$\nabla \log p_{\theta_{1:t}}(y_t | y_{1:t-1}) = \nabla \log p_{\theta_{1:t}}(y_{1:t}) - \nabla \log p_{\theta_{1:t-1}}(y_{1:t-1})$$

is given by the difference of two pseudo-score vectors where

$$\begin{aligned} \nabla \log p_{\theta_{1:t}}(y_{1:t}) := & \int \left(\sum_{k=2}^t \nabla \log f_{\theta}(x_k | x_{k-1}) \Big|_{\theta_k} \right. \\ & \left. + \nabla \log g_{\theta}(y_k | x_k) \Big|_{\theta_k} \right) p_{\theta_{1:t}}(x_{1:t} | y_{1:t}) dx_{1:t}. \end{aligned}$$

- Particle approximation follows

$$\theta_{t+1} = \theta_t + \gamma_t \widehat{\nabla \log p_{\theta_{1:t}}}(y_t | y_{1:t-1})$$

where

$$\widehat{\nabla \log p_{\theta_{1:t}}}(y_t | y_{1:t-1}) = \widehat{\nabla \log p_{\theta_{1:t}}}(y_{1:t}) - \widehat{\nabla \log p_{\theta_{1:t-1}}}(y_{1:t-1})$$

is given by the difference of particle estimates of pseudo-score vectors (Poyadjis, D. & Singh, 2011).

- Asymptotic variance of $\widehat{\nabla \log p_{\theta_{1:t}}}(y_t | y_{1:t-1})$ is uniformly bounded in $\mathcal{O}(1/N)$ for FB estimate whereas it is $\mathcal{O}(t/N)$ for direct particle method (Del Moral, D. & Singh, 2011). Bias is $\mathcal{O}(1/N)$ in both cases.
- **Major Problem:** If we use FB, this is not an online algorithm anymore as it requires a backward pass of order $\mathcal{O}(t)$ to approximate $\nabla \log p_{\theta_{1:t}}(y_{1:t}) \dots$

Variance of the Gradient Estimate for Direct vs FB

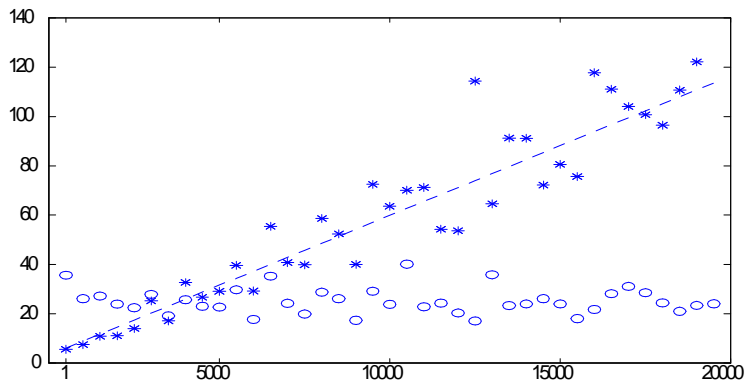


Figure: Empirical variance of the gradient estimate for standard versus FB approximations (SV model)

Online Particle ML Inference using Direct Approach

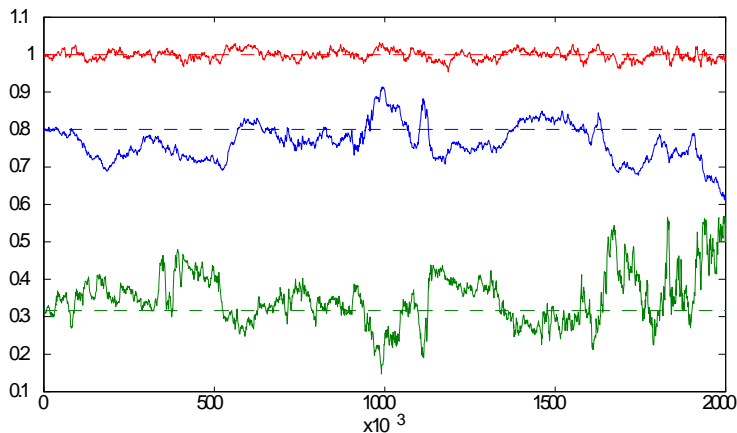


Figure: $N = 10,000$ particles, online parameter estimates for SV model.

Online Particle ML Inference using FB

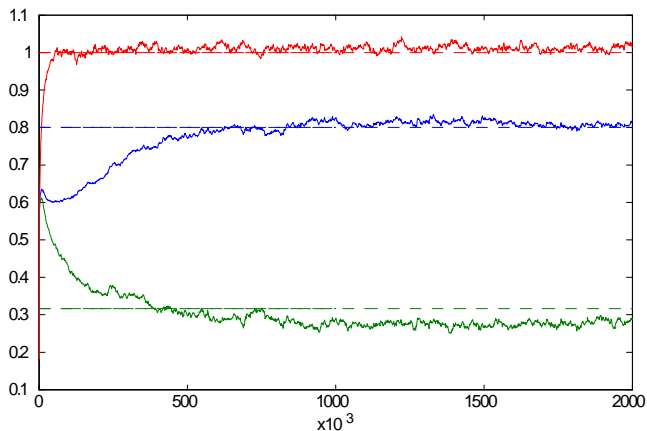


Figure: $N = 50$ particles, online parameter estimates for SV model.

Forward only Smoothing

- Dynamic programming allows us to compute in a single forward pass the FB estimates of

$$\varphi_t^\theta = \int \varphi_t(x_{1:t}) p_\theta(x_{1:t} | y_{1:t}) dx_{1:t}$$

where

$$\varphi_t(x_{1:t}) = \sum_{k=1}^t \varphi(x_{k-1}, x_k, y_k)$$

- Forward Backward (FB) decomposition states

$$p_\theta(x_{1:T} | y_{1:T}) = p_\theta(x_T | y_{1:T}) \prod_{t=1}^{T-1} p_\theta(x_t | y_{1:t}, x_{t+1})$$

where $p_\theta(x_t | y_{1:t}, x_{t+1}) = \frac{f_\theta(x_{t+1} | x_t) p_\theta(x_t | y_{1:t})}{p_\theta(x_{t+1} | y_{1:t})}$.

- Conditioned upon $y_{1:T}$, $\{X_t\}_{t=1}^T$ is a backward Markov chain of initial distribution $p(x_T | y_{1:T})$ and inhomogeneous Markov transitions $\{p_\theta(x_t | y_{1:t}, x_{t+1})\}_{t=1}^{T-1}$ independent of T .

- We have

$$\begin{aligned}\varphi_t^\theta &= \int \varphi_t(x_{1:t}) p_\theta(x_{1:t-1} | y_{1:t-1}, x_t) dx_{1:t-1} \\ &= \int \left\{ \underbrace{\int \varphi_t(x_{1:t}) p_\theta(x_{1:t-1} | y_{1:t-1}, x_t) dx_{1:t-1}}_{V_t^\theta(x_t)} \right\} p_\theta(x_t | y_{1:t}) dx_t\end{aligned}$$

- *Forward smoothing recursion*

$$V_t^\theta(x_t) = \int \left[V_{t-1}^\theta(x_{t-1}) + \varphi(x_{t-1:t}, y_t) \right] p_\theta(x_{t-1} | y_{1:t-1}, x_t) dx_{t-1}$$

- Appears implicitly in Elliott, Aggoun & Moore (1996), Ford (1998) and rediscovered a few times... Presentation follows here (Del Moral, D. & Singh, 2009).

Forward only Smoothing

- *Forward smoothing recursion*

$$V_t^\theta(x_t) = \int \left[V_{t-1}^\theta(x_{t-1}) + \varphi(x_{t-1:t}, y_t) \right] p_\theta(x_{t-1} | y_{1:t-1}, x_t) dx_{t-1}$$

- Proof is trivial

$$\begin{aligned} V_t^\theta(x_t) &= \int \varphi_t(x_{1:t}) p_\theta(x_{1:t-1} | y_{1:t-1}, x_t) dx_{1:t-1} \\ &= \int \left[\varphi_{t-1}(x_{1:t-1}) + \varphi(x_{t-1:t}, y_t) \right] p_\theta(x_{1:t-2} | y_{1:t-2}, x_{t-1}) \\ &\quad \times p_\theta(x_{t-1} | y_{1:t-1}, x_t) dx_{1:t-1} \\ &= \int \left\{ \underbrace{\int \varphi_{t-1}(x_{1:t-1}) p_\theta(x_{1:t-2} | y_{1:t-2}, x_{t-1}) dx_{1:t-2}}_{V_{t-1}^\theta(x_{t-1})} \right. \\ &\quad \left. + \varphi(x_{t-1:t}, y_t) \right\} p_\theta(x_{t-1} | y_{1:t-1}, x_t) dx_{t-1} \end{aligned}$$

- Exact implementation possible for finite state-space and linear Gaussian models.

Particle Forward only Smoothing

- At time $t - 1$, we have $\hat{p}_\theta(x_{t-1} | y_{1:t-1}) = \frac{1}{N} \sum_{i=1}^N \delta_{X_{t-1}^{(i)}}(x_{t-1})$ and $\{\hat{V}_{t-1}^\theta(X_{t-1}^{(i)})\}_{1 \leq i \leq N}$.
- At time t , compute $\hat{p}_\theta(x_t | y_{1:t}) = \sum_{i=1}^N W_t^{(i)} \delta_{X_t^{(i)}}(x_t)$ and set

$$\begin{aligned}\hat{V}_t^\theta(X_t^{(i)}) &= \int \left\{ \hat{V}_{t-1}^\theta(x_{t-1}) + \varphi(x_{t-1}, x_t, y_t) \right\} \hat{p}_\theta(x_{t-1} | y_{1:t-1}, X_t^{(i)}) dx_{t-1} \\ &= \frac{\sum_{j=1}^N f_\theta(X_t^{(i)} | X_{t-1}^{(j)}) [\hat{V}_{t-1}^\theta(X_{t-1}^{(j)}) + \varphi(X_{t-1}^{(j)}, X_t^{(i)}, y_t)]}{\sum_{j=1}^N f_\theta(X_t^{(i)} | X_{t-1}^{(j)})},\end{aligned}$$

$$\hat{\varphi}_t^\theta = \frac{1}{N} \sum_{i=1}^N \hat{V}_t^\theta(X_t^{(i)}).$$

- This estimate is exactly the same as the Particle FB estimate, computational complexity $\mathcal{O}(N^2)$.

- At time $t - 1$, we have $\hat{p}_{\theta_{1:t-1}}(x_{t-1} | y_{1:t-1})$, $\left\{ \hat{V}_{t-1}^{\theta_{1:t-1}}(X_{t-1}^{(i)}) \right\}$,
 $\widehat{\nabla \log p_{\theta_{1:t-1}}}(y_{1:t-1}) = \int \hat{V}_{t-1}^{\theta_{1:t-1}}(x_{t-1}) \hat{p}_{\theta_{1:t-1}}(x_{t-1} | y_{1:t-1}) dx_{t-1}$
and get θ_t .

- At time t , use your favourite PF to compute $\hat{p}_{\theta_{1:t}}(x_t | y_{1:t})$ and

$$\hat{V}_t^{\theta_{1:t}}(X_t^{(i)}) = \int \left\{ \hat{V}_{t-1}^{\theta_{1:t-1}}(x_{t-1}) + \varphi(x_{t-1}, x_t, y_t) \right\} \\ \times \hat{p}_{\theta_{1:t}}(x_{t-1} | y_{1:t-1}, X_t^{(i)}) dx_{t-1},$$
$$\varphi(x_{t-1:t}, y_t) = \nabla \log f_{\theta}(x_t | x_{t-1})|_{\theta_t} + \nabla \log g_{\theta}(y_t | x_t)|_{\theta_t}$$

and

$$\widehat{\nabla \log p_{\theta_{1:t}}}(y_{1:t}) = \int \hat{V}_t^{\theta_{1:t}}(x_t) \hat{p}_{\theta_{1:t}}(x_t | y_{1:t}) dx_t$$

- Parameter update

$$\theta_{t+1} = \theta_t + \gamma_t \left(\widehat{\nabla \log p_{\theta_{1:t}}}(y_{1:t}) - \widehat{\nabla \log p_{\theta_{1:t-1}}}(y_{1:t-1}) \right)$$

- Online Bayesian parameter inference using particle methods is yet an unsolved problem.
- Particle smoothing techniques can be used to perform off-line and on-line ML parameter estimation.
- Observed information matrix can also be evaluated online in a stable manner.
- For online inference, computational complexity is $\mathcal{O}(N^2)$ at each time step and requires evaluating $f_{\theta}(x_t | x_{t-1})$.