

Sequential Monte Carlo methods

Lecture 8 – Path space view of the particle filter

Fredrik Lindsten, Linköping University 2025-02-04

Aim: Introduce the path space view of the particle filter, explain the path degeneracy problem and briefly mention the low-variance resampling methods.

Outline:

- 1. Path space view of the particle filter
- 2. Path degeneracy
- 3. Mitigating the path degeneracy problem
 - a. Effective samples size (ESS)
 - b. Low variance resampling
- 4. Parameter inference in SSMs

Reminder – the bootstrap particle filter

Algorithm 1 Bootstrap particle filter (for i = 1, ..., N)

- 1. Initialization (t = 0):
- (a) Sample $\mathbf{x}_0^i \sim p(\mathbf{x}_0)$.
- (b) Set initial weights: $w_0^i = 1/N$.
- 2. for t = 1 to T do
- (a) **Resample:** sample ancestor indices $a_t^i \sim C(\{w_{t-1}^j\}_{j=1}^N)$.
- (b) **Propagate:** sample $x_t^i \sim p(x_t | x_{t-1}^{a_t^i})$ and set $x_{0:t}^i = (x_{0:t-1}^{a_t^i}, x_t^i)$.
- (c) Weight: compute $\widetilde{w}_t^i = p(y_t | x_t^i)$ and normalize $w_t^i = \widetilde{w}_t^i / \sum_{j=1}^N \widetilde{w}_t^j$.

The ancestor indices $\{a_t^i\}_{i=1}^N$ allow us to keep track of exactly what happens in each resampling step.

Note the bookkeeping added to the propagation step 2b.

Example evolution of three particles for t = 0, 1, 2.



Example evolution of three particles for t = 0, 1, 2.



The **ancestral path** of x_2^1 , i.e. $x_{0:2}^1$, is shown as the thick line.

Bookkeeping – ancestor indices

At time t = 1, particle x_0^2 is resampled twice and particle x_0^3 is resampled once (whereas particle x_0^1 is not resampled). Hence, at time t = 1, the **ancestor indices** are

$$a_1^1 = 2, a_1^2 = 2$$
 and $a_1^3 = 3$.

Similarly, at time t = 2, the **ancestor indices** are given by

$$a_2^1 = 2, a_2^2 = 3$$
 and $a_2^3 = 3$.

The **ancestral path** of x_2^1 , i.e. $x_{0:2}^1$, is shown as a thick line. It is defined recursively from the ancestor indices

$$x_{0:2}^{1} = (x_{0}^{a_{1}^{2}}, x_{1}^{a_{2}^{1}}, x_{2}^{1}) = (x_{0}^{a_{1}^{2}}, x_{1}^{2}, x_{2}^{1}) = (x_{0}^{2}, x_{1}^{2}, x_{2}^{1}).$$

Algorithm 2 joint filtering bootstrap PF (for i = 1, ..., N)

- 1. Initialization (t = 0):
- (a) Sample $x_0^i \sim p(x_0)$.
- (b) Set initial weights: $w_0^i = 1/N$.
- 2. for t = 1 to T do
- (a) **Resample:** sample ancestor indices $a_t^i \sim C(\{w_{t-1}^j\}_{j=1}^N)$.
- (b) **Propagate:** sample $x_t^i \sim p(x_t | x_{t-1}^{a_t^i})$ and set $x_{0:t}^i = (x_{0:t-1}^{a_t^i}, x_t^i)$.

(c) Weight: compute $\widetilde{w}_t^i = p(y_t | x_t^i)$ and normalize $w_t^j = \widetilde{w}_t^j / \sum_{j=1}^N \widetilde{w}_t^j$.

Bootstrap PF targeting the joint filtering PDF

It can be shown that Algorithm 2 targets the joint filtering pdf $p(\mathbf{x}_{0:t} | y_{1:t}) = p(\mathbf{x}_{0:t-1} | y_{1:t-1}) \frac{p(\mathbf{x}_t | \mathbf{x}_{t-1})p(y_t | \mathbf{x}_t)}{p(y_t) | y_{1:t-1})}.$

It resamples entire trajectories $x_{0:t}^{i}$, not just individual states x_{t}^{i} .

Bootstrap PF targeting the joint filtering PDF

It can be shown that Algorithm 2 targets the joint filtering pdf $p(x_{0:t} | y_{1:t}) = p(x_{0:t-1} | y_{1:t-1}) \frac{p(x_t | x_{t-1})p(y_t | x_t)}{p(y_t) | y_{1:t-1})}.$

It resamples entire trajectories $x_{0:t}^{i}$, not just individual states x_{t}^{i} .

Resulting approximation of the joint filtering PDF

$$\widehat{p}^{N}(\mathbf{x}_{0:t} | \mathbf{y}_{1:t}) = \sum_{i=1}^{N} w_{t}^{i} \delta_{\mathbf{x}_{0:t}^{i}}(\mathbf{x}_{0:t}).$$

Bootstrap PF targeting the joint filtering PDF

It can be shown that Algorithm 2 targets the joint filtering pdf $p(x_{0:t} | y_{1:t}) = p(x_{0:t-1} | y_{1:t-1}) \frac{p(x_t | x_{t-1})p(y_t | x_t)}{p(y_t) | y_{1:t-1})}.$

It resamples entire trajectories $x_{0:t}^{i}$, not just individual states x_{t}^{i} .

Resulting approximation of the joint filtering PDF

$$\widehat{p}^{N}(\mathbf{x}_{0:t} | \mathbf{y}_{1:t}) = \sum_{i=1}^{N} w_{t}^{i} \delta_{\mathbf{x}_{0:t}^{i}}(\mathbf{x}_{0:t}).$$

Problem: While it can be shown that the estimate $\hat{p}^N(x_{0:t} | y_{1:t})$ produced by Algorithm 2 converges asymptotically as $N \to \infty$, it is in many cases **not** a good approximation of $p(x_{0:t} | y_{1:t})!$

Path degeneracy

1D Gaussian random walk, measured in Gaussian noise, T = 25.

Target the joint filtering density using a bootstrap PF (Alg. 2) with N = 30 particles.

$$\widehat{p}(\mathbf{x}_{0:25} | \mathbf{y}_{1:25}) = \sum_{i=1}^{30} w_{25}^{i} \delta_{\mathbf{x}_{0:25}^{i}}(\mathbf{x}_{0:25}).$$

ex) Path degeneracy

ex) Path degeneracy



At each point in time all particles are plotted using a black dot and each particle is connected with its ancestor using a black line.

ex) Path degeneracy



At each point in time all particles are plotted using a black dot and each particle is connected with its ancestor using a black line.



The grey dots represents $\hat{p}(x_t | y_{1:t})$ at each point in time.

The black lines represents $\hat{p}(\mathbf{x}_{0:25} | \mathbf{y}_{1:t}).$

Note that all ancestral paths $\{x_{0:25}^i\}_{i=1}^N$ share a common ancestor at time t = 6 (and consequently for all times t < 6 as well).

Note that all ancestral paths $\{x_{0:25}^i\}_{i=1}^N$ share a common ancestor at time t = 6 (and consequently for all times t < 6 as well).

Let us use the resulting particle system $\{w_{25}^i, x_{0:25}^i\}_{i=1}^N$ to compute a Monte Carlo estimate of $\mathbb{E}[x_3 | y_{1:25}]$,

$$\mathbb{E}[\mathbf{x}_3 \mid \mathbf{y}_{1:25}] \approx \sum_{i=1}^{30} w_{25}^i \mathbf{x}_3^i$$

Note that all ancestral paths $\{x_{0:25}^i\}_{i=1}^N$ share a common ancestor at time t = 6 (and consequently for all times t < 6 as well).

Let us use the resulting particle system $\{w_{25}^i, x_{0:25}^i\}_{i=1}^N$ to compute a Monte Carlo estimate of $\mathbb{E}[x_3 \mid y_{1:25}]$,

$$\mathbb{E}[\mathbf{x}_3 \,|\, \mathbf{y}_{1:25}] \approx \sum_{i=1}^{30} w_{25}^i \mathbf{x}_3^i$$

Boils down to an estimate using a single sample, since x_3^i is identical for all i = 1, ..., 30.

Path degeneracy follows as a direct consequence of resampling.

The resampling step will by construction result in that for any time s there exists a time t > s, such that the PF approximation $\hat{p}^{N}(x_{0:t} | y_{1:t})$ consists of a single particle at time s.

In the above example this happened for s = 6 and t = 25.

Mitigating path degeneracy

The impact of the path degeneracy problem can be reduced:

- 1. Do not resample at each iteration, when?
- 2. Better resampling algorithms

3. ...

Effective sample size (ESS)

The effective sample size (ESS) $N_{\rm eff}$ is a diagnostics tool that tells us when our weights are problematic in the sense that they are close to being degenerate.

$$N_{\mathrm{eff}} = rac{N}{\mathbb{E}_q\left[\{\omega(x^i)\}^2
ight]} \leq N.$$

Effective sample size (ESS)

The effective sample size (ESS) N_{eff} is a diagnostics tool that tells us when our weights are problematic in the sense that they are close to being degenerate.

$$N_{\rm eff} = \frac{N}{\mathbb{E}_q\left[\{\omega(x^i)\}^2\right]} \le N.$$

We cannot evaluate $N_{\rm eff}$ exactly, but we can compute an estimate

$$\widehat{N}_{\text{eff}} = \frac{1}{\sum_{i=1}^{N} (w^i)^2}.$$

Effective sample size (ESS)

The effective sample size (ESS) $N_{\rm eff}$ is a diagnostics tool that tells us when our weights are problematic in the sense that they are close to being degenerate.

$$N_{\rm eff} = \frac{N}{\mathbb{E}_q\left[\{\omega(x^i)\}^2\right]} \le N.$$

We cannot evaluate $N_{\rm eff}$ exactly, but we can compute an estimate

$$\widehat{N}_{\text{eff}} = \frac{1}{\sum_{i=1}^{N} (w^i)^2}.$$

ESS-adaptive resampling: When \hat{N}_{eff} falls below some threshold N_{thres} we resample the particles, otherwise we continue without resampling.

Ex. 1) Let $w^i = 1/N$ for all i = 1, ..., N (independent samples),

$$\widehat{N}_{\text{eff}} = \frac{1}{\sum_{i=1}^{N} (w^i)^2} = \frac{1}{N \times 1/N^2} = N.$$

Ex. 1) Let $w^i = 1/N$ for all i = 1, ..., N (independent samples),

$$\widehat{N}_{\text{eff}} = \frac{1}{\sum_{i=1}^{N} (w^i)^2} = \frac{1}{N \times 1/N^2} = N.$$

Ex. 2) Let $w^i = 0$ for i = 1, ..., N - 1 and $w^N = 1$ (completely degenerate),

$$\widehat{N}_{\text{eff}} = \frac{1}{\sum_{i=1}^{N} (w^i)^2} = 1.$$

Bootstrap PF with ESS-adaptive resampling

Algorithm 3 joint filtering bootstrap PF (for i = 1, ..., N)

- 1. Initialization (t = 0):
- (a) Sample $\mathbf{x}_0^i \sim p(\mathbf{x}_0)$.
- (b) Set initial weights: $w_0^i = 1/N$.
- 2. for t = 1 to T do
- (a) Compute $\widehat{N}_{\text{eff}} = \frac{1}{\sum_{i=1}^{N} (w_{t-1}^i)^2}$.
- (b) ESS-adapted resample: If $\widehat{N}_{eff} < N_{thres}$ sample ancestor indices $a_t^i \sim C(\{w_{t-1}^j\}_{j=1}^N)$ and set $w_{t-1}^i = 1/N$. If $\widehat{N}_{eff} \ge N_{thres}$ set $a_t^i = i$.
- (c) **Propagate:** sample $x_t^i \sim p(x_t | x_{t-1}^{a_t^i})$. $x_{0:t}^i = \{x_{0:t-1}^{a_t^i}, x_t^i\}$.
- (d) Weight: compute $\widetilde{w}_t^i = p(y_t | x_t^i) w_{t-1}$ and normalize $w_t^i = \widetilde{w}_t^i / \sum_{j=1}^N \widetilde{w}_t^j$.

Multinomial resampling

Multinomial resampling introduced during lecture 4

$$a^i \sim \mathcal{C}(\{w^j\}_{j=1}^N), \qquad \mathbb{P}\left(a^i = j\right) = w^j.$$

Multinomial resampling

Multinomial resampling introduced during lecture 4

$$a^i \sim \mathcal{C}(\{w^j\}_{j=1}^N), \qquad \mathbb{P}\left(a^i = j\right) = w^j.$$



Blue circular disc – weights $\{w^i\}_{i=1}^8$.

Solid arrows – selected particles $\{x^i\}_{i=1}^8$.

Alternative implementations of resampling



Divide the circle into strata (grey dashed lines).

Stratified resampling randomly selects 1 sample from each strata.

Alternative implementations of resampling



Divide the circle into strata (grey dashed lines).

Stratified resampling randomly selects 1 sample from each strata.

Systematic resampling randomly generates 1 offset and then it picks one sample from each strata using this offset. N.B. There are some (pathological) cases when it can perform worse than simple multinomial sampling.

Alternatively, create [*Nwⁱ*] **deterministic** copies of particle *i*.

We then sample the remaining $N - \sum_{i=1}^{N} \lfloor Nw^i \rfloor$ particles by applying any resampling scheme (e.g. stratified) to the **residuals**, i.e. we use weights

$$\frac{Nw^{i} - \lfloor Nw^{i} \rfloor}{N - \sum_{j=1}^{N} \lfloor Nw^{j} \rfloor}$$

for the resampling.

Path degeneracy can mitigated by **backward simulation** (results in particle **smoothers**).



Fredrik Lindsten and Thomas B. Schön. Backward simulation methods for Monte Carlo statistical inference. Foundations and Trends in Machine Learning, 6(1):1-143, 2013.

Path degeneracy can mitigated by **backward simulation** (results in particle **smoothers**).



Fredrik Lindsten and Thomas B. Schön. Backward simulation methods for Monte Carlo statistical inference. Foundations and Trends in Machine Learning, 6(1):1-143, 2013.

Particle MCMC algorithms (lectures 11–14) can be used to tackle the path degeneracy issue and solve the state smoothing problem (jointly with inferring unknown model parameters).



Christophe Andrieu, Arnaud Doucet and Roman Holenstein. Particle Markov chain Monte Carlo methods. Journal of the Royal Statistical Society: Series B, 72:269-342, 2010. In estimating the fixed-lag smoothing density $p(x_{t-l+1:t} | y_{1:t})$ for some small l > 1 we can make use of

$$\widehat{\rho}(\mathbf{x}_{t-l+1:t} | \mathbf{y}_{1:t}) = \sum_{i=1}^{N} w_{t}^{i} \delta_{\mathbf{x}_{t-l+1:t}^{i}}(\mathbf{x}_{t-l+1:t}),$$

where the particle system comes from a particle filter targeting the joint filtering density.

If *l* is taken too large we activate the path degeneracy problem to such a degree that it will not work.

Parameter inference in SSMs

$$\begin{aligned} X_t &= f(X_{t-1}, \theta) + V_t, \\ Y_t &= g(X_t, \theta) + E_t, \\ X_0 &\sim p(x_0 \mid \theta). \end{aligned} \qquad \begin{aligned} X_t \mid (X_{t-1} = x_{t-1}, \theta = \theta) &\sim p(x_t \mid x_{t-1}, \theta), \\ Y_t \mid (X_t = x_t, \theta = \theta) &\sim p(y_t \mid x_t, \theta), \\ X_0 &\sim p(x_0 \mid \theta). \end{aligned}$$

Two different parameter inference formulations differing in the way the unknown parameters θ are modelled:

- Maximum likelihood: *θ* modelled as deterministic.
- Bayesian: θ modelled as stochastic.

Central object - data distribution/likelihood

The data distribution can be computed by marginalizing

$$p(\mathbf{x}_{0:T}, y_{1:T} \mid \boldsymbol{\theta}) = \prod_{t=1}^{T} p(y_t \mid \mathbf{x}_t, \boldsymbol{\theta}) \prod_{t=1}^{T} p(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \boldsymbol{\theta}) p(\mathbf{x}_0 \mid \boldsymbol{\theta})$$

w.r.t. the state trajectory $x_{0:T}$

$$p(y_{1:T} \mid \boldsymbol{\theta}) = \int p(x_{0:T}, y_{1:T} \mid \boldsymbol{\theta}) dx_{0:T}.$$

Central object - data distribution/likelihood

The data distribution can be computed by marginalizing

$$p(\mathbf{x}_{0:T}, y_{1:T} \mid \boldsymbol{\theta}) = \prod_{t=1}^{T} p(y_t \mid \mathbf{x}_t, \boldsymbol{\theta}) \prod_{t=1}^{T} p(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \boldsymbol{\theta}) p(\mathbf{x}_0 \mid \boldsymbol{\theta})$$

w.r.t. the state trajectory $x_{0:T}$

$$p(y_{1:T} \mid \boldsymbol{\theta}) = \int p(\mathbf{x}_{0:T}, y_{1:T} \mid \boldsymbol{\theta}) d\mathbf{x}_{0:T}.$$

Average over all possible values for the state trajectory $x_{0:T}$.

Central object - data distribution/likelihood

The data distribution can be computed by marginalizing

$$p(\mathbf{x}_{0:T}, y_{1:T} \mid \boldsymbol{\theta}) = \prod_{t=1}^{T} p(y_t \mid \mathbf{x}_t, \boldsymbol{\theta}) \prod_{t=1}^{T} p(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \boldsymbol{\theta}) p(\mathbf{x}_0 \mid \boldsymbol{\theta})$$

w.r.t. the state trajectory $x_{0:T}$

$$p(y_{1:T} \mid \boldsymbol{\theta}) = \int p(\mathbf{x}_{0:T}, y_{1:T} \mid \boldsymbol{\theta}) d\mathbf{x}_{0:T}.$$

Average over all possible values for the state trajectory $x_{0:T}$.

Alternative way of performing the averaging:

$$p(y_{1:T} | \boldsymbol{\theta}) = \prod_{t=1}^{T} p(y_t | y_{1:t-1}, \boldsymbol{\theta}) = \prod_{t=1}^{T} \int p(y_t | x_t, \boldsymbol{\theta}) \underbrace{p(x_t | y_{1:t-1}, \boldsymbol{\theta})}_{\text{approx. by PF}} dx_t$$

Ancestral path: By starting from a particle x_{t}^{i} at time t and tracing its ancestors backwards in time via the ancestor indices we obtain $x_{0:t}^{i}$, which is the ancestral path for particle x_{t}^{i} .

Path degeneracy: The resampling step will by construction result in that for any time s there exists a time t > s such that the PF approximation $\hat{p}(x_{0:t} | y_{1:t})$ consists of a single particle at time s.

Effective sample size (ESS): An importance sampling diagnostics tool that tells us when our weights are problematic in the sense that they are close to being degenerate, i.e. it provides a way of gauging the extent of the weight degeneracy.

Backward simulation: Generates samples backwards in time. When backward sampling can be implemented it removes the path degeneracy problem (only possible in off-line situations).

Likelihood function: Deterministic function of θ obtained by inserting the available measurements into the data distribution.