

Index

Introduction

Nested particle filters

 Nested filtering

 Nested particle filter (NPF)

Gradient-based exploration of the parameter space

State-space model

We are interested in systems can be represented by **Markov state-space dynamical models**:

$$\begin{aligned} \text{(state)} \quad \mathbf{x}_t &= \mathbf{f}(\mathbf{x}_{t-1}, \boldsymbol{\theta}) + \mathbf{v}_t, \\ \text{(observation)} \quad \mathbf{y}_t &= \mathbf{g}(\mathbf{x}_t, \boldsymbol{\theta}) + \mathbf{r}_t, \end{aligned}$$

In terms of a set of relevant probability density functions (pdfs):

$$\boldsymbol{\theta} \sim p(\boldsymbol{\theta}) \text{ and } \mathbf{x}_0 \sim p(\mathbf{x}_0) \quad (1)$$

$$\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}, \boldsymbol{\theta}) \quad (2)$$

$$\mathbf{y}_t \sim p(\mathbf{y}_t | \mathbf{x}_t, \boldsymbol{\theta}) \quad (3)$$

State-space model

We are interested in systems can be represented by **Markov state-space dynamical models**:

$$\begin{aligned} \text{(state)} \quad \mathbf{x}_t &= \mathbf{f}(\mathbf{x}_{t-1}, \boldsymbol{\theta}) + \mathbf{v}_t, \\ \text{(observation)} \quad \mathbf{y}_t &= \mathbf{g}(\mathbf{x}_t, \boldsymbol{\theta}) + \mathbf{r}_t, \end{aligned}$$

In terms of a set of **relevant probability density functions (pdfs)**:

$$\boldsymbol{\theta} \sim p(\boldsymbol{\theta}) \text{ and } \mathbf{x}_0 \sim p(\mathbf{x}_0) \quad (1)$$

$$\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}, \boldsymbol{\theta}) \quad (2)$$

$$\mathbf{y}_t \sim p(\mathbf{y}_t | \mathbf{x}_t, \boldsymbol{\theta}) \quad (3)$$

Goal

→ We want to **approximate the joint posterior distribution of θ and \mathbf{x}_t** , i.e., $p(\mathbf{x}_t, \theta | \mathbf{y}_{1:t})$.

→ For a long sequence of observations, i.e., **online**.

State-of-the-art methods

Methods for Bayesian inference of both θ and \mathbf{x}_t :

- **particle Markov chain Monte Carlo (PMCMC)**¹
- **sequential Monte Carlo square (SMC²)**²
- **nested particle filters (NPFs)**³

- They can **quantify the uncertainty** or estimation error.
- They can be applied to a **broad class of models**.
- They provide **theoretical guarantees**.
- Both PMCMC and SMC² are **batch techniques**, while the NPF is a **recursive method**.

¹Andrieu, Doucet, and Holenstein 2010.

²Chopin, Jacob, and Papaspiliopoulos 2013.

³Crisan and Míguez 2018.

State-of-the-art methods

Methods for Bayesian inference of both θ and \mathbf{x}_t :

- **particle Markov chain Monte Carlo (PMCMC)**¹
- **sequential Monte Carlo square (SMC²)**²
- **nested particle filters (NPFs)**³

→ They can **quantify the uncertainty** or estimation error.

→ They can be applied to a **broad class of models**.

→ They provide **theoretical guarantees**.

→ Both PMCMC and SMC² are **batch techniques**, while the NPF is a **recursive method**.

¹Andrieu, Doucet, and Holenstein 2010.

²Chopin, Jacob, and Papaspiliopoulos 2013.

³Crisan and Míguez 2018.

State-of-the-art methods

Methods for Bayesian inference of both θ and \mathbf{x}_t :

- **particle Markov chain Monte Carlo (PMCMC)**¹
 - **sequential Monte Carlo square (SMC²)**²
 - **nested particle filters (NPFs)**³
- They can **quantify the uncertainty** or estimation error.
- They can be applied to a **broad class of models**.
- They provide **theoretical guarantees**.
- Both PMCMC and SMC² are **batch techniques**, while the NPF is a **recursive method**.

¹Andrieu, Doucet, and Holenstein 2010.

²Chopin, Jacob, and Papaspiliopoulos 2013.

³Crisan and Míguez 2018.

Index

Introduction

Nested particle filters

 Nested filtering

 Nested particle filter (NPF)

Gradient-based exploration of the parameter space

Nested filtering

We aim at computing the joint posterior pdf $p(\boldsymbol{\theta}, \mathbf{x}_t | \mathbf{y}_{1:t})$, as

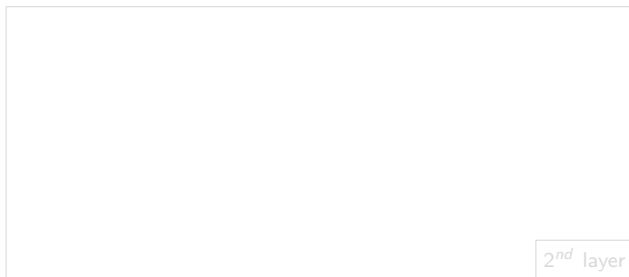
$$p(\mathbf{x}_t, \boldsymbol{\theta} | \mathbf{y}_{1:t}) = \underbrace{p(\mathbf{x}_t | \boldsymbol{\theta}, \mathbf{y}_{1:t})}_{2^{nd} \text{ layer}} \underbrace{p(\boldsymbol{\theta} | \mathbf{y}_{1:t})}_{1^{st} \text{ layer}}$$

Nested filtering

At every time step t :

$$\underbrace{p(\theta|y_{1:t-1})}_{\text{Pred. pdf of } \theta}$$

1st layer



$$p(y_t|\theta, y_{1:t-1})$$

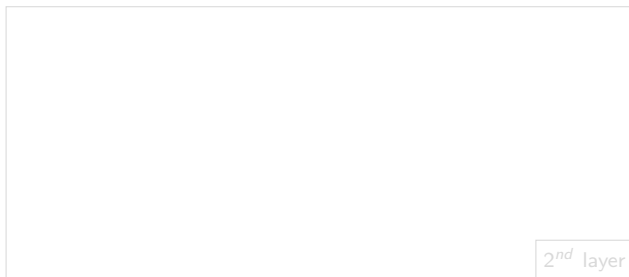
$$\underbrace{p(\theta|y_{1:t})}_{\text{Post. pdf of } \theta} \propto p(y_t|\theta, y_{1:t-1})p(\theta|y_{1:t-1})$$

Nested filtering

At every time step t :

$$\underbrace{p(\theta|y_{1:t-1})}_{\text{Pred. pdf of } \theta}$$

1st layer



2nd layer

$$p(y_t|\theta, y_{1:t-1}) = \int p(y_t|x_t, \theta) p(x_t|\theta, y_{1:t-1}) dx_t$$

$$\underbrace{p(\theta|y_{1:t})}_{\text{Post. pdf of } \theta} \propto p(y_t|\theta, y_{1:t-1}) p(\theta|y_{1:t-1})$$

Nested filtering

At every time step t :

$$\underbrace{p(\theta|y_{1:t-1})}_{\text{Pred. pdf of } \theta}$$

1st layer

Classical filtering problem (given θ)

Predictive pdf of x : $p(x_t|y_{1:t-1}, \theta)$

Likelihood: $p(y_t|x_t, \theta)$

Posterior pdf of x : $p(x_t|y_{1:t}, \theta)$

2nd layer

$$p(y_t|\theta, y_{1:t-1}) = \int p(y_t|x_t, \theta)p(x_t|\theta, y_{1:t-1})dx_t$$

$$\underbrace{p(\theta|y_{1:t})}_{\text{Post. pdf of } \theta} \propto p(y_t|\theta, y_{1:t-1})p(\theta|y_{1:t-1})$$

Nested filtering

At every time step t :

$$\underbrace{p(\theta|y_{1:t-1})}_{\text{Pred. pdf of } \theta}$$

1st layer

Filtering (given θ)

Predictive pdf of x : $p(x_t|y_{1:t-1}, \theta)$

Likelihood: $p(y_t|x_t, \theta)$

Posterior pdf of x : $p(x_t|y_{1:t}, \theta)$

2nd layer

$$p(y_t|\theta, y_{1:t-1}) = \int p(y_t|x_t, \theta)p(x_t|\theta, y_{1:t-1})dx_t$$

$$\underbrace{p(\theta|y_{1:t})}_{\text{Post. pdf of } \theta} \propto p(y_t|\theta, y_{1:t-1})p(\theta|y_{1:t-1})$$

Jittering

- NPF \longrightarrow jittering: $\bar{\theta}_t^i \sim \kappa_N(d\theta|\theta')$, where

$$\kappa_N(d\theta|\theta') = (1 - \epsilon_N)\delta_{\theta'}(\theta) + \epsilon_N\kappa(d\theta|\theta')$$

- $0 < \epsilon_N \leq \frac{1}{\sqrt{N}}$
 - $\kappa(d\theta|\theta')$ is an arbitrary Markov kernel with mean θ' and finite variance, e.g., $\kappa(d\theta|\theta') = \mathcal{N}(\theta|\theta', \tilde{\sigma}^2 I)$, with $\tilde{\sigma}^2 < \infty$.
- Guarantees convergence to the true posterior when $N \longrightarrow \infty$

Take-aways

Advantages:

- Only framework that is *online* and *Bayesian* on θ
- Applicable to general parametric state-space models
- Asymptotic convergence guarantees

Limitations:

- Coverage speed might be slow (depends on the jittering (hyper)parameters)
- This problem gets worse as the dimension of θ increases

Take-aways

Advantages:

- Only framework that is *online* and *Bayesian* on θ
- Applicable to general parametric state-space models
- Asymptotic convergence guarantees

Limitations:

- Coverage speed might be slow (depends on the jittering (hyper)parameters)
- This problem gets worse as the dimension of θ increases

Langevin nested particle filter (LNPF)

Iteratively, move the parameters towards areas of higher probability with the Unadjusted Langevin algorithm (ULA):

$$\boldsymbol{\theta}_{t,k+1}^l = \boldsymbol{\theta}_{t,k}^l + \gamma_k \cdot \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta} \mid \mathbf{y}_{1:t}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t,k}^l} + \sqrt{2\gamma_k} \mathbf{v}_k, \quad (4)$$

where $\mathbf{v}_k \sim \mathcal{N}(0, \mathbf{I}_{d_{\boldsymbol{\theta}}})$, $\gamma_k > 0$ is a step size sequence.

We replace jittering by ULA:

$$\text{(jittering)} \quad \bar{\boldsymbol{\theta}}_t^i \sim \kappa_{\text{LNPF}}(d\boldsymbol{\theta}|\boldsymbol{\theta}') = (1 - \epsilon_N)\delta_{\boldsymbol{\theta}'}(\boldsymbol{\theta}) + \epsilon_N \kappa_{\text{jitter}}(d\boldsymbol{\theta}|\boldsymbol{\theta}') \quad (5)$$

$$\text{(ULA)} \quad \bar{\boldsymbol{\theta}}_t^i \sim \kappa_{\text{LNPF}}(d\boldsymbol{\theta}|\boldsymbol{\theta}') = (1 - \epsilon_N)\delta_{\boldsymbol{\theta}'}(\boldsymbol{\theta}) + \epsilon_N \kappa_{\text{ULA}}(d\boldsymbol{\theta}|\boldsymbol{\theta}') \quad (6)$$

*(ULA needs to adjust number of iterations and step size)

Langevin nested particle filter (LNPF)

Iteratively, move the parameters towards areas of higher probability with the Unadjusted Langevin algorithm (ULA):

$$\boldsymbol{\theta}_{t,k+1}^l = \boldsymbol{\theta}_{t,k}^l + \gamma_k \cdot \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta} \mid \mathbf{y}_{1:t}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t,k}^l} + \sqrt{2\gamma_k} \mathbf{v}_k, \quad (4)$$

where $\mathbf{v}_k \sim \mathcal{N}(0, \mathbf{I}_{d_{\boldsymbol{\theta}}})$, $\gamma_k > 0$ is a step size sequence.

We replace jittering by ULA:

$$\text{(jittering)} \quad \bar{\boldsymbol{\theta}}_t^i \sim \kappa_{LNPF}(d\boldsymbol{\theta}|\boldsymbol{\theta}') = (1 - \epsilon_N)\delta_{\boldsymbol{\theta}'}(\boldsymbol{\theta}) + \epsilon_N \kappa_{\text{jitter}}(d\boldsymbol{\theta}|\boldsymbol{\theta}') \quad (5)$$

$$\text{(ULA)} \quad \bar{\boldsymbol{\theta}}_t^i \sim \kappa_{LNPF}(d\boldsymbol{\theta}|\boldsymbol{\theta}') = (1 - \epsilon_N)\delta_{\boldsymbol{\theta}'}(\boldsymbol{\theta}) + \epsilon_N \kappa_{ULA}(d\boldsymbol{\theta}|\boldsymbol{\theta}') \quad (6)$$

*(ULA needs to adjust number of iterations and step size)

Langevin nested particle filter (LNPF)

Iteratively, move the parameters towards areas of higher probability with the Unadjusted Langevin algorithm (ULA):

$$\boldsymbol{\theta}_{t,k+1}^l = \boldsymbol{\theta}_{t,k}^l + \gamma_k \cdot \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta} \mid \mathbf{y}_{1:t}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t,k}^l} + \sqrt{2\gamma_k} \mathbf{v}_k, \quad (4)$$

where $\mathbf{v}_k \sim \mathcal{N}(0, \mathbf{I}_{d_{\boldsymbol{\theta}}})$, $\gamma_k > 0$ is a step size sequence.

We replace jittering by ULA:

$$\text{(jittering)} \quad \bar{\boldsymbol{\theta}}_t^i \sim \kappa_{NPF}(d\boldsymbol{\theta}|\boldsymbol{\theta}') = (1 - \epsilon_N)\delta_{\boldsymbol{\theta}'}(\boldsymbol{\theta}) + \epsilon_N \kappa_{\text{jitter}}(d\boldsymbol{\theta}|\boldsymbol{\theta}') \quad (5)$$

$$\text{(ULA)} \quad \bar{\boldsymbol{\theta}}_t^i \sim \kappa_{LNPF}(d\boldsymbol{\theta}|\boldsymbol{\theta}') = (1 - \epsilon_N)\delta_{\boldsymbol{\theta}'}(\boldsymbol{\theta}) + \epsilon_N \kappa_{ULA}(d\boldsymbol{\theta}|\boldsymbol{\theta}') \quad (6)$$

*(ULA needs to adjust number of iterations and step size)

Challenges: intractable score

The gradient of interest includes the score, that is **intractable** and its computation is **not recursive**.

$$\nabla_{\theta} \log p(\theta \mid \mathbf{y}_{1:t}) = \nabla_{\theta} \log p(\mathbf{y}_{1:t} \mid \theta) + \nabla_{\theta} \log p(\theta) \quad (7)$$

→ Using Fisher's identity we have two approximations⁵

$$\mathcal{O}(N): \quad \nabla_{\theta} \log p(\mathbf{y}_{1:t} \mid \theta) = \int \nabla_{\theta} \log p(\mathbf{y}_{1:t}, \mathbf{x}_{1:t} \mid \theta) p(\mathbf{x}_{1:t} \mid \mathbf{y}_{1:t}, \theta) d\mathbf{x}_{1:t}$$

$$\mathcal{O}(N^2): \quad \nabla_{\theta} \log p(\mathbf{y}_{1:t} \mid \theta) = \int \nabla_{\theta} \log p(\mathbf{y}_{1:t}, \mathbf{x}_t \mid \theta) p(\mathbf{x}_t \mid \mathbf{y}_{1:t}, \theta) d\mathbf{x}_t$$

*(This is described for a fixed θ)

⁵Poyiadjis, Doucet, and Singh 2011.

Challenges: intractable score

The gradient of interest includes the score, that is **intractable** and its computation is **not recursive**.

$$\nabla_{\theta} \log p(\theta \mid \mathbf{y}_{1:t}) = \nabla_{\theta} \log p(\mathbf{y}_{1:t} \mid \theta) + \nabla_{\theta} \log p(\theta) \quad (7)$$

→ Using Fisher's identity we have two approximations⁵

$$\mathcal{O}(N): \quad \nabla_{\theta} \log p(\mathbf{y}_{1:t} \mid \theta) = \int \nabla_{\theta} \log p(\mathbf{y}_{1:t}, \mathbf{x}_{1:t} \mid \theta) p(\mathbf{x}_{1:t} \mid \mathbf{y}_{1:t}, \theta) d\mathbf{x}_{1:t}$$

$$\mathcal{O}(N^2): \quad \nabla_{\theta} \log p(\mathbf{y}_{1:t} \mid \theta) = \int \nabla_{\theta} \log p(\mathbf{y}_{1:t}, \mathbf{x}_t \mid \theta) p(\mathbf{x}_t \mid \mathbf{y}_{1:t}, \theta) d\mathbf{x}_t$$

*(This is described for a fixed θ)

⁵Poyiadjis, Doucet, and Singh 2011.

Challenges: intractable score

The gradient of interest includes the score, that is **intractable** and its computation is **not recursive**.

$$\nabla_{\theta} \log p(\theta \mid \mathbf{y}_{1:t}) = \nabla_{\theta} \log p(\mathbf{y}_{1:t} \mid \theta) + \nabla_{\theta} \log p(\theta) \quad (7)$$

→ Using Fisher's identity we have two approximations⁵

$$\mathcal{O}(N): \quad \nabla_{\theta} \log p(\mathbf{y}_{1:t} \mid \theta) = \int \nabla_{\theta} \log p(\mathbf{y}_{1:t}, \mathbf{x}_{1:t} \mid \theta) p(\mathbf{x}_{1:t} \mid \mathbf{y}_{1:t}, \theta) d\mathbf{x}_{1:t}$$

$$\mathcal{O}(N^2): \quad \nabla_{\theta} \log p(\mathbf{y}_{1:t} \mid \theta) = \int \nabla_{\theta} \log p(\mathbf{y}_{1:t}, \mathbf{x}_t \mid \theta) p(\mathbf{x}_t \mid \mathbf{y}_{1:t}, \theta) d\mathbf{x}_t$$

*(This is described for a fixed θ)

⁵Poyiadjis, Doucet, and Singh 2011.

Challenges: recursive approximation

→ Approximation used in the recursive maximum likelihood literature⁶, such that

$\nabla_{\theta} \log p(\mathbf{y}_{1:t} | \theta)$ is replaced by $\nabla_{\theta} \log p(\mathbf{y}_t | \mathbf{y}_{1:t-1}, \theta)$,

and it can be computed as

$$\left. \nabla_{\theta} \log p(\mathbf{y}_t | \mathbf{y}_{1:t-1}, \theta) \right|_{\theta=\theta_t} = \left. \nabla_{\theta} \log p(\mathbf{y}_{1:t} | \theta) \right|_{\theta=\theta_{1:t}} - \left. \nabla_{\theta} \log p(\mathbf{y}_{1:t-1} | \theta) \right|_{\theta=\theta_{1:t-1}}$$

*(We are assuming that process \mathbf{y}_t is ergodic)

⁶Chopin, Papaspiliopoulos, et al. 2020.

Challenges: recursive approximation

→ Approximation used in the recursive maximum likelihood literature⁶, such that

$\nabla_{\theta} \log p(\mathbf{y}_{1:t} | \theta)$ is replaced by $\nabla_{\theta} \log p(\mathbf{y}_t | \mathbf{y}_{1:t-1}, \theta)$,

and it can be computed as

$$\left. \nabla_{\theta} \log p(\mathbf{y}_t | \mathbf{y}_{1:t-1}, \theta) \right|_{\theta=\theta_t} = \left. \nabla_{\theta} \log p(\mathbf{y}_{1:t} | \theta) \right|_{\theta=\theta_{1:t}} - \left. \nabla_{\theta} \log p(\mathbf{y}_{1:t-1} | \theta) \right|_{\theta=\theta_{1:t-1}}$$

*(We are assuming that process \mathbf{y}_t is ergodic)

⁶Chopin, Papaspiliopoulos, et al. 2020.

Challenges: recursive approximation

→ Approximation used in the recursive maximum likelihood literature⁶, such that

$\nabla_{\theta} \log p(\mathbf{y}_{1:t} | \theta)$ is replaced by $\nabla_{\theta} \log p(\mathbf{y}_t | \mathbf{y}_{1:t-1}, \theta)$,

and it can be computed as

$$\left. \nabla_{\theta} \log p(\mathbf{y}_t | \mathbf{y}_{1:t-1}, \theta) \right|_{\theta=\theta_t} = \left. \nabla_{\theta} \log p(\mathbf{y}_{1:t} | \theta) \right|_{\theta=\theta_{1:t}} - \left. \nabla_{\theta} \log p(\mathbf{y}_{1:t-1} | \theta) \right|_{\theta=\theta_{1:t-1}}$$

*(We are assuming that process \mathbf{y}_t is ergodic)

⁶Chopin, Papaspiliopoulos, et al. 2020.

Conclusions

- Jittering can be very inefficient, especially with larger d_θ
- We proposed ULA updates for smarter exploration of θ space
- Challenge is approximating the score online and accurately

References

- Crisan & Míguez (2018), *Nested particle filters for online parameter estimation in discrete-time state-space Markov models*. Bernoulli, vol. 24, no. 4A, pp. 3039–3086.
- Chopin, N., Jacob, P. E., & Papaspiliopoulos, O. (2013). *SMC²: an efficient algorithm for sequential analysis of state space models*. Journal of the Royal Statistical Society Series B: Statistical Methodology, 75(3), 397-426.
- Andrieu, Christophe, Arnaud Doucet, & Roman Holenstein (2010), *Particle Markov chain Monte Carlo methods*. Journal of the Royal Statistical Society Series B: Statistical Methodology 72.3 : 269-342.
- Poyiadjis, G., Doucet, A., & Singh, S. S. (2011). *Particle approximations of the score and observed information matrix in state space models with application to parameter estimation*. Biometrika, 98(1), 65-80.
- Chopin, N., & Papaspiliopoulos, O. (2020). *An introduction to sequential Monte Carlo* (Vol. 4). Cham, Switzerland: Springer.

Thank you!

<https://sarapv.github.io/>