

Hess-MC²: Sequential Monte Carlo Squared using Hessian Information and Second Order Proposals

Joshua Murphy*, Conor Rosato[†], Andrew Millard*, Lee Devlin*, Paul Horridge* and Simon Maskell*

* Department of Electrical Engineering and Electronics, University of Liverpool, United Kingdom

[†] Department of Pharmacology and Therapeutics, University of Liverpool, United Kingdom

Email: {joshua.murphy, cmrosa, andrew.millard, ljdevlin, p.horridge, smaskell}@liverpool.ac.uk

Abstract—When performing Bayesian inference using Sequential Monte Carlo (SMC) methods, two considerations arise: the accuracy of the posterior approximation and computational efficiency. To address computational demands, Sequential Monte Carlo Squared (SMC²) is well-suited for high-performance computing (HPC) environments. The design of the proposal distribution within SMC² can improve accuracy and exploration of the posterior as poor proposals may lead to high variance in importance weights and particle degeneracy. The Metropolis-Adjusted Langevin Algorithm (MALA) uses gradient information so that particles preferentially explore regions of higher probability. In this paper, we extend this idea by incorporating second-order information, specifically the Hessian of the log-target. While second-order proposals have been explored previously in particle Markov Chain Monte Carlo (p-MCMC) methods, we are the first to introduce them within the SMC² framework. Second-order proposals not only use the gradient (first-order derivative), but also the curvature (second-order derivative) of the target distribution. Experimental results on synthetic models highlight the benefits of our approach in terms of step-size selection and posterior approximation accuracy when compared to other proposals.

Index Terms—Bayesian inference, Parameter estimation, Sequential Monte Carlo, Differentiable particle filters

I. INTRODUCTION

Bayesian inference offers a framework for uncertainty quantification, but remains computationally challenging in complex or high-dimensional models. Sequential Monte Carlo (SMC) methods are a flexible class of algorithms for approximating

posterior distributions, especially in non-linear non-Gaussian problem settings. However, two key considerations must be balanced in SMC algorithms: computational efficiency and accuracy of the posterior approximation. Sequential Monte Carlo Squared (SMC²) is one such method which is well-suited to estimating the dynamic states and parameters of state space models (SSM) [1]. SMC² consists of two layers: an inner particle filter (PF) [2] layer which tracks the dynamic states and an SMC sampler [3] which estimates the parameters. Recent work in [4] introduced a framework for deploying SMC² on distributed memory architectures, which addresses the computational efficiency problem. This parallelizability represents a significant advantage over many other Bayesian inference algorithms. In [4], only a random walk (RW) proposal was considered but RW can converge slowly in high-dimensional parameter spaces. To use gradient-based proposals the PF needs to be differentiated which can be problematic due to some of the discrete operations being non-differentiable [5]. Devising methods for differentiating PFs is therefore an active area of research [6]–[9]. In the context of parameter estimation for SSMs, first-order (FO) gradient information has been incorporated into particle-Markov Chain Monte Carlo (p-MCMC) methods [10], [11], which improved proposal efficiency. Second-order (SO) information which incorporates Hessian matrices has been included in general MCMC [12] and p-MCMC which has further improved inference quality [13]. Hessian matrices capture both the slope (via gradients) and the local curvature of the posterior, allowing for more informed and efficient proposals.

In the context of SMC², recent work to increase accuracy has focused on improving the proposal, including [14], in which a differentiable Common Random Number (CRN) PF [9] is employed to obtain the gradient of the log-likelihood with respect to the parameters. The resulting gradients were then used within the Metropolis-Adjusted Langevin Algorithm (MALA) proposal [15] to adaptively guide particles toward regions of higher posterior density using local information about the target. Building on this idea, we propose the use of SO information in the form of the Hessian matrix of the log-target density. Using automatic differentiation in PyTorch [16], we compute SO derivatives to construct curvature-aware proposals that improve the concentration and diversity of

particles in high-probability regions. Choosing an appropriate step-size can be challenging for RW, FO, and SO proposals. Our analysis shows that accuracy is more sensitive to step-size in the case of RW proposals, exhibiting much higher variance compared to FO and SO methods.

The structure of this paper is as follows: in Section II we describe a PF with the log-likelihood, FO and SO gradients outlined in Sections II-A, II-B and II-C. The SMC sampler and different variants of the proposal are shown in Sections III and III-A, respectively. Two examples are shown in Section IV with conclusions and future work described in Section V.

II. PARTICLE FILTER

State-Space Models (SSMs) can be used to represent the dependence of latent states in non-linear non-Gaussian dynamical systems. An SSM consists of a state equation and an observation equation parameterized by θ

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

$$\mathbf{y}_t | \mathbf{x}_t \sim p(\mathbf{y}_t | \mathbf{x}_t, \theta). \quad (2)$$

An SSM sequentially simulates the latent states, $\mathbf{x}_{1:T} = \{\mathbf{x}_1, \dots, \mathbf{x}_t, \dots, \mathbf{x}_T\}$, conditional on data it receives at each time increment t , $\mathbf{y}_{1:T} = \{\mathbf{y}_1, \dots, \mathbf{y}_t, \dots, \mathbf{y}_T\}$, over T timesteps. The joint distribution at timestep t is

$$p(\mathbf{x}_{1:t}, \mathbf{y}_{1:t} | \theta) = p(\mathbf{x}_1 | \theta) p(\mathbf{y}_1 | \mathbf{x}_1, \theta) \times \prod_{\tau=2}^t p(\mathbf{x}_\tau | \mathbf{x}_{\tau-1}, \theta) p(\mathbf{y}_\tau | \mathbf{x}_\tau, \theta). \quad (3)$$

A PF can be used to recursively approximate the distribution in (3) using a set of N_x particles and an importance sampling approach. As time evolves, new particles are drawn from a proposal distribution, $q(\mathbf{x}_t | \mathbf{x}_{t-1}, \theta, \mathbf{y}_t)$. These particles are weighted at each timestep according to

$$\mathbf{w}_{1:t}^j = \mathbf{w}_{1:t-1}^j \frac{p(\mathbf{y}_t | \mathbf{x}_t^j, \theta) p(\mathbf{x}_t^j | \mathbf{x}_{t-1}^j, \theta)}{q(\mathbf{x}_t^j | \mathbf{x}_{t-1}^j, \theta, \mathbf{y}_t)}, \quad (4)$$

Normalized weights $\tilde{\mathbf{w}}_{1:t}^j$ are calculated by

$$\tilde{\mathbf{w}}_{1:t}^j = \frac{\mathbf{w}_{1:t}^j}{\sum_{j'=1}^{N_x} \mathbf{w}_{1:t}^{j'}}. \quad (5)$$

As time evolves, a small number of particles may have the majority of the weight, a phenomenon known as particle degeneracy. The number of effective samples can be used to monitor the performance of the samples is given by

$$N_x^{\text{eff}} = \frac{1}{\sum_{j=1}^{N_x} (\tilde{\mathbf{w}}_t^j)^2}, \quad (6)$$

and we resample when N_x^{eff} falls below $N_x/2$. Resampling involves sampling N_x new samples with replacement from the existing set with probabilities according to their weights. Unbiased estimates of integrals of functions of the posterior with respect to (3) are realized as

$$\int p(\mathbf{x}_{1:t}, \mathbf{y}_{1:t} | \theta) f(\mathbf{x}_{1:t}) d\mathbf{x}_{1:t} = \frac{1}{N_x} \sum_{j=1}^{N_x} \mathbf{w}_{1:t}^j f(\mathbf{x}_t^j). \quad (7)$$

A. Log-Likelihood

The log-posterior of the parameters (not the states) is

$$\log \pi(\theta) = \log p(\theta) + \log p(\mathbf{y}_{1:T} | \theta), \quad (8)$$

where $\log p(\theta)$ is the log-prior and $\log p(\mathbf{y}_{1:T} | \theta)$ represents the log-likelihood. An unbiased estimate of the log-likelihood is calculated from (7) with $f(\cdot) = 1$.

B. First-Order Gradients

In accordance with (8), the gradient of the log-posterior is

$$\nabla \log \pi(\theta) = \nabla \log p(\theta) + \nabla \log p(\mathbf{y}_{1:T} | \theta). \quad (9)$$

An approximation of the gradient of the likelihood is determined as a function of the derivative of the weights:

$$\frac{d}{d\theta} \log p(\mathbf{y}_{1:t} | \theta) \approx \frac{1}{N_x} \sum_{j=1}^{N_x} \frac{d}{d\theta} \mathbf{w}_{1:t}^j \quad (10)$$

The details for deriving the gradient of the log-likelihood using the chain rule can be found in Appendix A of [14].

C. Second-Order Gradients

In accordance with (8), the negative Hessian of the log-posterior is expressed as:

$$\nabla^2 \log \pi(\theta) = -\nabla^2 \log p(\theta) - \nabla^2 \log p(\mathbf{y}_{1:T} | \theta). \quad (11)$$

An approximation to the negative Hessian of the log-likelihood is determined as a function of the derivative of the log-weights:

$$\begin{aligned} \frac{d^2}{d\theta^2} \log p(\mathbf{y}_{1:t} | \theta) &\approx \frac{1}{N_x} \frac{d}{d\theta} \left(\sum_{j=1}^{N_x} \tilde{\mathbf{w}}_{1:t}^j \right) \frac{d}{d\theta} \log \mathbf{w}_{1:t}^j \\ &+ \frac{1}{N_x} \sum_{j=1}^{N_x} \tilde{\mathbf{w}}_{1:t}^j \frac{d^2}{d\theta^2} \log \mathbf{w}_{1:t}^j, \end{aligned} \quad (12)$$

where, by (5),

$$\frac{d}{d\theta} \sum_{j=1}^{N_x} \tilde{\mathbf{w}}_{1:t}^j = \sum_{j=1}^{N_x} \frac{d}{d\theta} \left(\frac{\mathbf{w}_{1:t}^j}{\sum_{i=1}^{N_x} \mathbf{w}_{1:t}^i} \right) \quad (13)$$

$$= \sum_{j=1}^{N_x} \tilde{\mathbf{w}}_{1:t}^j \left(\frac{d}{d\theta} \log \mathbf{w}_{1:t}^j - \frac{d}{d\theta} \log p(\mathbf{y}_{1:t} | \theta) \right). \quad (14)$$

Following the logic of [14] Appendix A, we can find an estimate of the second derivative of the log weights. Taking the second term in (12),

$$\frac{d^2}{d\theta^2} \log \mathbf{w}_{1:t}^j = \frac{d^2}{d\theta^2} \log \mathbf{w}_{1:t-1}^j + \frac{d^2}{d\theta^2} \log p(\mathbf{y}_t | \mathbf{x}_t^j). \quad (15)$$

Let

$$L(\mathbf{x}_t^j, \theta, \mathbf{y}_t) \triangleq \log p(\mathbf{y}_t | \mathbf{x}_t^j, \theta), \quad (16)$$

where the likelihood is Gaussian with a variance that is independent of \mathbf{x}_t^j , such that

$$L(\mathbf{x}_t^j, \theta, \mathbf{y}_t) \triangleq \log \mathcal{N}(\mathbf{y}_t; h(\mathbf{x}_t^j, \theta), R(\theta)). \quad (17)$$

The SO derivative of the likelihood L in (16) is given by

$$\begin{aligned} \frac{d^2}{d\theta^2} L(\mathbf{x}_t^j, \theta, \mathbf{y}_t) &= \frac{d}{d\theta} \left(\frac{\partial}{\partial h} \log \mathcal{N}(\mathbf{y}_t; h, R) \right) \cdot \frac{dh}{d\theta} \\ &+ \frac{\partial}{\partial h} \log \mathcal{N}(\mathbf{y}_t; h, R) \cdot \frac{d^2 h}{d\theta^2} \\ &+ \frac{d}{d\theta} \left(\frac{\partial}{\partial R} \log \mathcal{N}(\mathbf{y}_t; h, R) \right) \cdot \frac{dR}{d\theta} \\ &+ \frac{\partial}{\partial R} \log \mathcal{N}(\mathbf{y}_t; h, R) \cdot \frac{d^2 R}{d\theta^2}, \end{aligned} \quad (18)$$

for

$$\frac{dh}{d\theta} = \frac{\partial h}{\partial \mathbf{x}_t^j} \frac{d\mathbf{x}_t^j}{d\theta} + \frac{\partial h}{\partial \theta}, \quad (19)$$

Then,

$$\begin{aligned} \frac{d}{d\theta} \left(\frac{\partial}{\partial h} \log \mathcal{N}(\mathbf{y}_t; h, R) \right) &= \frac{\partial^2}{\partial^2 h} \log \mathcal{N}(\mathbf{y}_t; h, R) \cdot \frac{dh}{d\theta} \\ &+ \frac{\partial^2}{\partial h \partial R} \log \mathcal{N}(\mathbf{y}_t; h, R) \cdot \frac{dR}{d\theta}, \end{aligned} \quad (20)$$

$$\begin{aligned} \frac{dh^2}{d\theta^2} &= \frac{\partial^2 h}{\partial \mathbf{x}_t^j^2} \left(\frac{d\mathbf{x}_t^j}{d\theta} \right)^2 + \frac{2\partial^2 h}{\partial \theta \partial \mathbf{x}_t^j} \left(\frac{d\mathbf{x}_t^j}{d\theta} \right) + \frac{\partial h}{\partial \mathbf{x}_t^j} \left(\frac{d^2 \mathbf{x}_t^j}{d\theta^2} \right) \\ &+ \frac{\partial^2 h}{\partial \theta^2} \end{aligned} \quad (21)$$

and

$$\begin{aligned} \frac{d}{d\theta} \left(\frac{\partial}{\partial R} \log \mathcal{N}(\mathbf{y}_t; h, R) \right) &= \frac{\partial^2}{\partial h \partial R} \log \mathcal{N}(\mathbf{y}_t; h, R) \cdot \frac{dh}{d\theta} \\ &+ \frac{d^2}{d^2 R} \log \mathcal{N}(\mathbf{y}_t; h, R) \cdot \frac{dR}{d\theta}. \end{aligned} \quad (22)$$

III. SEQUENTIAL MONTE CARLO SQUARED

SMC² runs for K iterations, targeting the posterior of the parameters $\pi(\theta)$ at each iteration k . The joint distribution from all states until $k = K$ is defined to be

$$\pi(\theta_{1:K}) = \pi(\theta_K) \prod_{k=2}^K L_k(\theta_{k-1} | \theta_k), \quad (23)$$

where $L_k(\theta_{k-1} | \theta_k)$ is the L-kernel, a user-defined probability distribution. At $k = 1$, N samples are drawn from a prior distribution $q_1(\cdot)$ and weighted according to

$$\mathbf{v}_1^i = \frac{\pi(\theta_1^i)}{q_1(\theta_1^i)}. \quad (24)$$

At $k > 1$, subsequent samples are proposed based on samples from the previous iteration via a proposal distribution, $q(\theta_k^i | \theta_{k-1}^i)$. These samples are weighted according to

$$\mathbf{v}_k^i = \mathbf{v}_{k-1}^i \frac{\pi(\theta_k^i)}{\pi(\theta_{k-1}^i)} \frac{L_k(\theta_{k-1}^i | \theta_k^i)}{q_k(\theta_k^i | \theta_{k-1}^i)}. \quad (25)$$

At the SMC sampler level we employ a parallelized version of systematic resampling [17] when the effective sample size,

calculated as in (6), goes below $N/2$. Estimates of functions on the distribution are realized by

$$\mathbb{E}_\pi[f(\theta)] \approx \sum_{i=1}^N \tilde{\mathbf{v}}_k^i f(\theta_k^i), \quad (26)$$

with samples from previous iterations incorporated through recycling [18].

A. Proposals

In this paper we consider RW, FO, and SO proposals and derive L-kernels based off a change of variables (CoV). The proposals are

$$q_k(\theta_k | \theta_{k-1}) = \begin{cases} \mathcal{N}(\theta_{k-1}, \Gamma) & \text{RW} \\ \mathcal{N}(\theta_{k-1} + \frac{\Gamma}{2} G_{k-1}, \Gamma) & \text{FO} \\ \mathcal{N}(\theta_{k-1} + \frac{\Gamma}{2} H_{k-1} G_{k-1}, \Gamma H_{k-1}), & \text{SO} \end{cases} \quad (27)$$

where we let the gradient and Hessian of the log-posterior be $G_{k-1} = \nabla \log \pi(\theta_{k-1})$ and $H_{k-1} = (-\nabla^2 \log \pi(\theta_{k-1}))^{-1}$. These proposals are parameterized in terms of a scalar step size ϵ such that $\Gamma = \epsilon^2 \mathbf{I}$.

1) *Random-Walk*: The proposal and L-kernel are evaluated in (25) as

$$q_k(\theta_k | \theta_{k-1}) = \mathcal{N}(\theta_k; \theta_{k-1}, \epsilon^2 \mathbf{I}), \quad (28)$$

$$L_k(\theta_{k-1} | \theta_k) = \mathcal{N}(\theta_{k-1}; \theta_k, \epsilon^2 \mathbf{I}). \quad (29)$$

2) *First-Order*: We can rewrite the FO proposal in (27) as

$$\theta_k = \theta_{k-1} + \frac{\epsilon^2}{2} G_{k-1} + \epsilon \mathbf{p}_{k-1}, \quad (30)$$

and can then propose samples according to

$$\mathbf{p}_{k-1} \sim \mathcal{N}(0, \mathbf{I}), \quad (31)$$

$$\mathbf{p}_{k-0.5} = \frac{\epsilon}{2} G_{k-1} + \mathbf{p}_{k-1}, \quad (32)$$

$$\theta_k = \theta_{k-1} + \epsilon \mathbf{p}_{k-0.5}, \quad (33)$$

$$\mathbf{p}_k = \frac{\epsilon}{2} G_k + \mathbf{p}_{k-0.5}, \quad (34)$$

which resembles the leapfrog integrator. As in [19] leapfrog is considered to be a function, f_{LF} which transforms θ_{k-1} to θ_k so the proposal can be rewritten through a CoV

$$q_k(\theta_k | \theta_{k-1}) = q_k^p(f_{LF}(\theta_{k-1}, \mathbf{p}_{k-1}) | \theta_{k-1}), \quad (35)$$

$$= q_k^p(\mathbf{p}_{k-1} | \theta_{k-1}) \left| \frac{df_{LF}(\theta_{k-1}, \mathbf{p}_{k-1})}{d\mathbf{p}_{k-1}} \right|^{-1} \quad (36)$$

$$= \mathcal{N}(\mathbf{p}_{k-1}; 0, \mathbf{I}) \left| \frac{df_{LF}(\theta_{k-1}, \mathbf{p}_{k-1})}{d\mathbf{p}_{k-1}} \right|^{-1}. \quad (37)$$

The L-kernel can likewise be written in terms of f_{LF} as leapfrog is reversible so that the negative of the final momentum \mathbf{p}_k will bring a sample from θ_k to θ_{k-1}

$$\begin{aligned} L_k(\theta_{k-1} | \theta_k) &= L_k^p(-\mathbf{p}_k | \theta_k) \left| \frac{df_{LF}(\theta_k, -\mathbf{p}_k)}{d\mathbf{p}_k} \right|^{-1} \\ &= \mathcal{N}(-\mathbf{p}_k; 0, \mathbf{I}) \left| \frac{df_{LF}(\theta_k, -\mathbf{p}_k)}{d\mathbf{p}_k} \right|^{-1}. \end{aligned} \quad (38)$$

3) *Second-Order*: As in Section III-A2, we write the SO proposal in leapfrog form

$$\mathbf{p}_{k-1} \sim \mathcal{N}(0, H_{k-1}^{-1}), \quad (39)$$

$$\mathbf{p}_{k-0.5} = \frac{\epsilon}{2} G_{k-1} + \mathbf{p}_{k-1}, \quad (40)$$

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \epsilon H_{k-1} \mathbf{p}_{k-0.5}, \quad (41)$$

$$\mathbf{p}_k = \frac{\epsilon}{2} G_k + \mathbf{p}_{k-0.5}, \quad (42)$$

and through a CoV can evaluate the proposal and L-kernel as

$$q_k(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1}) = \mathcal{N}(\mathbf{p}_{k-1}; 0, H_{k-1}^{-1}) \left| \frac{df_{LF}(\boldsymbol{\theta}_{k-1}, \mathbf{p}_{k-1})}{d\mathbf{p}_{k-1}} \right|^{-1}, \quad (43)$$

$$L_k(\boldsymbol{\theta}_{k-1} | \boldsymbol{\theta}_k) = \mathcal{N}(-\mathbf{p}_k; 0, H_{k-1}^{-1}) \left| \frac{df_{LF}(\boldsymbol{\theta}_k, -\mathbf{p}_k)}{d\mathbf{p}_k} \right|^{-1}. \quad (44)$$

There are cases where H_{k-1} is not a positive semi-definite matrix which is likely to occur far from posterior modes where the amount of information is limited. While there are heuristics to handle non positive semi-definite estimates of the Hessian [20], we choose to revert to the FO so as not to rely on insufficient information to effectively inform moves in the parameter space.

IV. EXAMPLES

For both examples considered, we generate data for 50 different seeds and perform parameter inference for each dataset and average the results. The algorithmic setup consists of using $N_x = 500$ particles in each PF and $N = 32$ samples in the SMC sampler over $K = 15$ iterations. We use the parallelized SMC² framework outlined in [9], [14] which enables multiple instances of the computationally intensive inner PF to be run in parallel. For each proposal we take 20 step sizes in a range that maintains numerical stability. We evaluate performance in terms of root mean square error (RMSE) between the true and the mean parameter estimates obtained from the SMC sampler. The relative runtime (RR) of FO and SO proposals are calculated with respect to RW. The analysis was performed on a distributed memory cluster equipped with two Xeon Gold 6138 CPUs, providing 384GB of memory and 40 cores. The code can be found here ¹.

A. Linear Gaussian State Space Model

We consider a Linear Gaussian State Space (LGSS) model outlined in [9], [21] with state and observation equations

$$\mathbf{x}_t | \mathbf{x}_{t-1} \sim \mathcal{N}(\mathbf{x}_t; \mu \mathbf{x}_{t-1}, \phi^2), \quad (45)$$

$$\mathbf{y}_t | \mathbf{x}_t \sim \mathcal{N}(\mathbf{y}_t; \mathbf{x}_t, \sigma^2), \quad (46)$$

where $\boldsymbol{\theta} = \{\mu, \phi, \sigma\} = \{0.75, 1, 1\}$. The PF uses the “optimal” proposal which can be derived from (46)

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}, \boldsymbol{\theta}, \mathbf{y}_t) = \mathcal{N}(\mathbf{x}_t; \rho^2 [\sigma^{-2} \mathbf{y}_t + \phi^{-2} \mu \mathbf{x}_{t-1}], \rho^2), \quad (47)$$

with $\rho^{-2} = \phi^{-2} + \sigma^{-2}$. The PF weights are updated using

$$\mathbf{w}_t^j = \mathbf{w}_{t-1}^j \mathcal{N}(\mathbf{y}_t; \mu \mathbf{x}_t, \rho^2). \quad (48)$$

¹<https://github.com/j-j-murphy/SMC-Squared-Langevin>

Each random seed is run for $T = 500$. We use $\mathcal{U}(0, 1)$, $\mathcal{U}(0, 2)$ and $\mathcal{U}(0, 2)$ as priors over the parameters. Figure 1 demonstrates how there is far less variability in RMSE across step size when using FO and SO compared to RW. Table I shows a significant benefit in the RMSE when using gradient-based proposals. While runtime for both FO and SO is higher, this is mitigated by less need for manual step-size tuning. We note that none of the 20 RW step-sizes leads to a better RMSE than SO. In the context of this model, the benefit of SO relative to FO is arguably insufficient to warrant the relatively large increase in runtime.

TABLE I: LGSS Results for median ϵ by RMSE.

Prop.- ϵ	E[μ]	E[ϕ]	E[σ]	RMSE	RR
RW-0.7	0.606	1.22	0.926	56.9×10^{-3}	1.00
FO-0.03	0.717	1.06	0.958	1.93×10^{-3}	2.45
SO-1.55	0.734	1.04	0.955	1.50×10^{-3}	12.2

B. Susceptible-Infected-Recovered Model

We also consider the Susceptible-Infected-Recovered (SIR) epidemiological model outlined in [22], [23]. A discrete time approximation of the SIR model is presented below

$$\mathbf{S}_t = \mathbf{S}_{t-1} - \beta \mathbf{I}_{t-1} \mathbf{S}_{t-1} + \epsilon_\beta, \quad (49)$$

$$\mathbf{I}_t = \mathbf{I}_{t-1} + \beta \mathbf{I}_{t-1} \mathbf{S}_{t-1} - \gamma \mathbf{I}_{t-1} - \epsilon_\beta + \epsilon_\gamma, \quad (50)$$

$$\mathbf{R}_t = N_{pop} - \mathbf{S}_t - \mathbf{I}_t, \quad (51)$$

where $\boldsymbol{\theta} = \{\beta, \gamma\}$. The number of individuals in each compartment at time $t = 0$ are denoted $\mathbf{S}_0 = N_{pop} - \mathbf{I}_0$, $\mathbf{I}_0 = 1$ and $\mathbf{R}_0 = 0$ and the total population is $N_{pop} = 763$. Stochasticity is introduced by including a noise term, ϵ_θ for each parameter, which is independently drawn from $\mathcal{N}(0, 0.5)$. PF weights are updated using

$$\mathbf{w}_t^j = \mathbf{w}_{t-1}^j \mathcal{P}(\mathbf{y}_t; \mathbf{I}_t). \quad (52)$$

Each random seed is run for $T = 36$ and we use $\mathcal{U}(0, 1)$ as priors for both parameters. The results reflect those of the LGSS model. In Fig. 1, there is more variance in the RMSE estimates with RW compared to FO and more of a spread in FO than SO. However, as is also evident in Fig. 1, when looking at the best performing ϵ of RW, it performs better than that of FO. This demonstrates that it is possible to select an ϵ that makes RW competitive with gradient-based proposals but time consuming as it required evaluating 20 different step-sizes. On this smaller problem, the difference in RR and median RMSE is less pronounced for FO and SO over RW and the vast majority of the 20 RW step-sizes do not produce a better RMSE when compared to SO. In this model, the increase in runtime appears to offer commensurate improvements in RMSE to that which would be likely to be achieved using a correspondingly larger number of samples.

V. CONCLUSIONS AND FUTURE WORK

In this paper, we extend the work of [14] by computing SO gradients from a CRN-PF and incorporating them into

TABLE II: SIR Results for median ϵ by RMSE.

Prop.- ϵ	$E[\beta]$	$E[\gamma]$	RMSE	RR
RW-0.55	0.605	0.299	3.94×10^{-3}	1.00
FO-0.008	0.580	0.292	2.01×10^{-3}	1.52
SO-2.05	0.592	0.295	1.15×10^{-3}	5.26

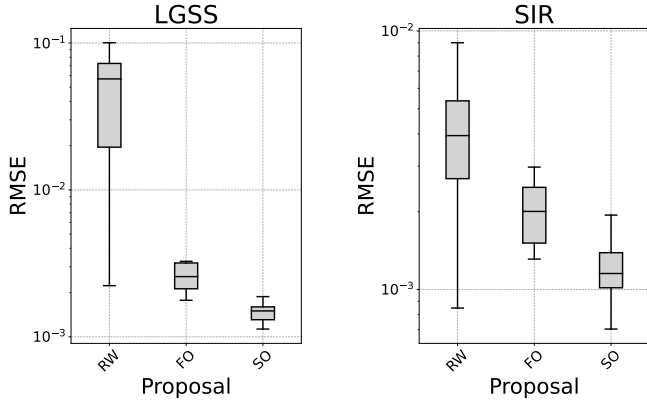


Fig. 1: RMSE distribution over 20 step-sizes by proposal.

SO proposals in the form of a Hessian matrix. We document how selecting an appropriate step-size can be time-consuming and show that SO, when compared to FO and RW, has less variability in the RMSE estimates of the parameters of an SSM. The two chosen examples are relatively simple, low-dimensional models and further work may chose to explore whether the benefits of SO proposals extend to more complex high-dimensional models.

One recent method of adapting the step-size in SMC samplers can be found in [24] which could be applicable in this work. Additional gradient-based proposals that could be employed include HMC, ChEES [25] and NUTS which was first described in [9] for particle MCMC.

REFERENCES

- [1] N. Chopin, P. E. Jacob, and O. Papaspiliopoulos, “SMC²: an efficient algorithm for sequential analysis of state space models,” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 75, no. 3, pp. 397–426, 2013.
- [2] M. S. Arulampalam, S. Maskell, N. Gordon, and T. Clapp, “A tutorial on particle filters for online nonlinear/non-gaussian bayesian tracking,” *IEEE Transactions on signal processing*, vol. 50, no. 2, pp. 174–188, 2002.
- [3] P. Del Moral, A. Doucet, and A. Jasra, “Sequential monte carlo samplers,” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 68, no. 3, pp. 411–436, 2006.
- [4] C. Rosato, A. Varsi, J. Murphy, and S. Maskell, “An $O(\log_2 N)$ SMC² algorithm on distributed memory with an approx. optimal l-kernel,” in *2023 IEEE Symposium Sensor Data Fusion and International Conference on Multisensor Fusion and Integration (SDF-MFI)*, pp. 1–8, IEEE, 2023.
- [5] X. Chen and Y. Li, “An overview of differentiable particle filters for data-adaptive sequential Bayesian inference,” *Foundations of Data Science*, pp. 0–0, 2023.
- [6] C. Nemeth, P. Fearnhead, and L. Mihaylova, “Particle approximations of the score and observed information matrix for parameter estimation in state–space models with linear computational cost,” *Journal of Computational and Graphical Statistics*, vol. 25, no. 4, pp. 1138–1157, 2016.

- [7] A. Corenflos, J. Thornton, G. Deligiannidis, and A. Doucet, “Differentiable particle filtering via entropy-regularized optimal transport,” in *International Conference on Machine Learning*, pp. 2100–2111, PMLR, 2021.
- [8] P. Karkus, D. Hsu, and W. S. Lee, “Particle filter networks with application to visual localization,” in *Conference on robot learning*, pp. 169–178, PMLR, 2018.
- [9] C. Rosato, L. Devlin, V. Beraud, P. Horridge, T. B. Schön, and S. Maskell, “Efficient learning of the parameters of non-linear models using differentiable resampling in particle filters,” *IEEE Transactions on Signal Processing*, vol. 70, pp. 3676–3692, 2022.
- [10] J. Dahlin, F. Lindsten, and T. B. Schön, “Particle metropolis hastings using Langevin dynamics,” in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 6308–6312, IEEE, 2013.
- [11] C. Nemeth, C. Sherlock, and P. Fearnhead, “Particle metropolis-adjusted Langevin algorithms,” *Biometrika*, vol. 103, no. 3, pp. 701–717, 2016.
- [12] M. Girolami and B. Calderhead, “Riemann manifold Langevin and Hamiltonian monte carlo methods,” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 73, no. 2, pp. 123–214, 2011.
- [13] J. Dahlin, F. Lindsten, and T. B. Schön, “Particle metropolis–hastings using gradient and Hessian information,” *Statistics and computing*, vol. 25, no. 1, pp. 81–92, 2015.
- [14] C. Rosato, J. Murphy, A. Varsi, P. Horridge, and S. Maskell, “Enhanced SMC²: Leveraging gradient information from differentiable particle filters within Langevin proposals,” in *2024 IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems (MFI)*, pp. 1–8, IEEE, 2024.
- [15] G. O. Roberts and R. L. Tweedie, “Exponential convergence of Langevin distributions and their discrete approximations,” *Bernoulli*, pp. 341–363, 1996.
- [16] A. Paszke, “Pytorch: An imperative style, high-performance deep learning library,” *arXiv preprint arXiv:1912.01703*, 2019.
- [17] A. Varsi, S. Maskell, and P. G. Spirakis, “An $O(\log_2 N)$ fully-balanced resampling algorithm for particle filters on distributed memory architectures,” *Algorithms*, vol. 14, no. 12, pp. 342–362, 2021.
- [18] T. L. T. Nguyen, F. Septier, G. W. Peters, and Y. Delignon, “Efficient sequential Monte-Carlo samplers for Bayesian inference,” *IEEE Transactions on Signal Processing*, vol. 64, no. 5, pp. 1305–1319, 2015.
- [19] L. Devlin, M. Carter, P. Horridge, P. L. Green, and S. Maskell, “The no-u-turn sampler as a proposal distribution in a sequential Monte Carlo sampler without accept/reject,” *IEEE Signal Processing Letters*, 2024.
- [20] J. Nocedal and S. J. Wright, *Numerical optimization*. Springer, 1999.
- [21] J. Dahlin and T. B. Schön, “Getting started with particle metropolis-hastings for inference in nonlinear dynamical models,” *Journal of Statistical Software*, vol. 88, pp. 1–41, 2019.
- [22] D. M. Sheinson, J. Niemi, and W. Meiring, “Comparison of the performance of particle filter algorithms applied to tracking of a disease epidemic,” *Mathematical biosciences*, vol. 255, pp. 21–32, 2014.
- [23] C. Rosato, J. Harris, J. Panovska-Griffiths, and S. Maskell, “Inference of stochastic disease transmission models using particle-MCMC and a gradient based proposal,” in *2022 25th International Conference on Information Fusion (FUSION)*, pp. 1–8, IEEE, 2022.
- [24] K. Kim, Z. Xu, J. R. Gardner, and T. Campbell, “Tuning sequential Monte Carlo samplers via greedy incremental divergence minimization,” *arXiv preprint arXiv:2503.15704*, 2025.
- [25] A. Millard, J. Murphy, D. Frisch, and S. Maskell, “Incorporating the chees criterion into sequential monte carlo samplers,” *arXiv preprint arXiv:2504.02627*, 2025.