

---

# Sequential Core-Set Monte Carlo (Supplementary Material)

---

Boyan Beronov<sup>1</sup>

Christian Weillbach<sup>1</sup>

Frank Wood<sup>1</sup>

Trevor Campbell<sup>2</sup>

<sup>1</sup>Dept. of Computer Science, University of British Columbia, Vancouver, Canada

<sup>2</sup>Dept. of Statistics, University of British Columbia, Vancouver, Canada

## 1 ALGORITHMS

### 1.1 SEQUENTIAL MONTE CARLO

---

**Algorithm 1** Resample-move SMC [Chopin, 2002, Gilks and Berzuini, 2001] for data tempering in exchangeable models

---

```
1: procedure SMC(population size  $K$ , proc. REJUV)
  ▷ initialize particles and memory
2:  $(\pi_k^{(0)}, \theta_k^{(0)}) \leftarrow$  Equation (2),  $M^{(0)} \leftarrow 0$ 
3: for  $t \leftarrow 1, \dots, T$  do
  ▷ update particle weights with new data batch  $(x_b^{(t)})_{b=1}^B$ 
4:  $(\bar{\pi}_k^{(t)})_{k=1}^K \leftarrow$  Equation (3)
  ▷ expand memory with new data batch  $(x_b^{(t)})_{b=1}^B$ 
5:  $M^{(t)} \leftarrow M^{(t-1)} + B$ ,  $(1, \bar{u}_j^{(t)})_{j=1}^{M^{(t)}} \leftarrow$  Equation (11)
  ▷ resample particles
6: if  $\text{ESS}(\bar{\pi}^{(t)}) < \text{threshold}$  then
7:  $(\hat{\pi}_k^{(t)}, \hat{\theta}_k^{(t)})_{k=1}^K \leftarrow$  Equation (5)
8: else
9:  $(\hat{\pi}_k^{(t)}, \hat{\theta}_k^{(t)})_{k=1}^K \leftarrow (\bar{\pi}_k^{(t)}, \theta_k^{(t-1)})_{k=1}^K$ 
10: end if
  ▷ rejuvenate particles using memory (Algorithm 2)
11:  $(\pi_k^{(t)}, \theta_k^{(t)})_{k=1}^K \leftarrow \text{REJUV}((\hat{\pi}_k^{(t)}, \hat{\theta}_k^{(t)})_{k=1}^K, (1, \bar{u}_j^{(t)})_{j=1}^{M^{(t)}})$ 
12: end for
13: return particles  $(\pi_k^{(T)}, \theta_k^{(T)})_{k=1}^K$ 
14: end procedure
```

---

For completeness, we provide in Algorithm 1 the original resample-move SMC method for sequential inference in static models, using the notation in this work. In comparison to SCMC (Algorithm 1), it has uniform data weights in line 5, lacks the core-set projection in lines 12-17, and does not produce a core-set memory as an output.

### 1.2 WEIGHTED METROPOLIS-HASTINGS FOR CORE-SET REJUVENATION

The simplest possible rejuvenation method, used inside Algorithm 1 for all experiments, is the Metropolis-Hastings kernel in Algorithm 2. The only difference to the standard procedure is the non-uniform weighting of data likelihoods in line 4, cf. Section 3.2. See Section 2 for details on the choice of proposal kernels  $q$  in individual experiments.

---

**Algorithm 2** Weighted rejuvenation kernel using Metropolis-Hastings

---

```
1: procedure REJUV(population  $(\widehat{\pi}_k^{(t)}, \widehat{\theta}_k^{(t)})_{k=1}^K$ , memory  $(\bar{w}_j^{(t)}, \bar{u}_j^{(t)})_{j=1}^{\bar{C}^{(t)}}$ , steps  $J$ , proposal  $q$ )
2:    $\forall k \in [K] : \theta_k^{(0)} \leftarrow \widehat{\theta}_k^{(t)}$ 
3:   for  $j \leftarrow 1, \dots, J$  do
4:      $\forall k \in [K] : \theta'_k \sim q(\cdot | \theta_k^{(j-1)})$ ,  $a_k^{(j)} \leftarrow \alpha(\theta'_k, \theta_k^{(j-1)})$ 
5:      $\forall k \in [K] : \theta_k^{(j)} \sim a_k^{(j)} \cdot \delta_{\theta'_k} + (1 - a_k^{(j)}) \cdot \delta_{\theta_k^{(j-1)}}$ 
6:   end for
7:   return  $(\widehat{\pi}_k^{(t)}, \theta_k^{(J)})_{k=1}^K$ 
8: end procedure
```

---

### 1.3 GREEDY ITERATIVE GEODESIC ASCENT

Since Hilbert core-set projection can be reduced to the standard sparse nonnegative least-squares (SNNLS) problem in Equation (16), SMC (Algorithm 1) is parametrized by the choice of an SNNLS solver. In this work, we choose GIGA [Campbell and Broderick, 2018] for this purpose, because it was designed specifically for this application, and because it provides the currently best accuracy-cost trade-off. Bayesian core-set methods such as [Campbell and Beronov, 2019] can reach higher levels of precision, but require their own internal sequential inference procedures with significantly higher runtimes.

We note that prior work on model-agnostic Bayesian core-set construction was limited to the batch setting [Huggins et al., 2016, Campbell and Broderick, 2018, 2019, Campbell and Beronov, 2019], in which the input data have uniform weights and the target vector in Equation (15) is given by  $b = A \cdot \vec{1}$ . The iterative core-set construction task performed by CPF instead requires the solution of an SNNLS problem with general weights as in Equation (16), and we provide the natural generalization of GIGA to this scenario in Algorithm 3.

## 2 EXPERIMENT DETAILS

In this section, we provide analytical background on the models and the parameters for each experiment. The hyperparameters for rejuvenation kernels are listed in Table 1, and all experiments are conducted with 3 rejuvenation steps. Following the concentration rate of a conjugate Gaussian posterior, we reduce the variance of proposal kernels linearly in the number of data points, where we denote the total number of observations at SMC step  $t$  by  $\sigma_t$ .

### 2.1 AUTO-REGRESSIVE PROCESS

#### 2.1.1 Rejuvenation Proposal

We use the same Gaussian kernel as in Equation (6), with parameters listed in Table 1 and  $\Sigma = 3$ . We have done no tuning and, given the simplicity of the experiment, assume that many other parameters would work as well.

### 2.2 NORMAL-INVERSE-WISHART

When the generative model constitutes an exponential family, the space of core-set posteriors generally lies within the same family. Our first experiment makes use of this circumstance in order to directly evaluate the approximation error of CPF in Figure 3, without introducing an additional inference problem.

---

**Algorithm 3** Greedy iterative geodesic ascent—generalized from  $b = A \cdot \bar{1}$  in [Campbell and Broderick, 2018]

---

1: **procedure** GIGA(dictionary  $A \in \mathbb{R}^{K \times N}$ , target  $b \in \mathbb{R}^K$ , core-set size  $M$ , tolerance  $\varepsilon$ )  
      $\triangleright$  *normalize vectors and initialize weights to 0*  
 2:  $\beta \leftarrow \frac{b}{\|b\|_2}$ ,  $\forall n \in [N] : \alpha_n \leftarrow \frac{A \cdot e_n}{\|A \cdot e_n\|_2}$   
 3:  $w_0 \leftarrow 0$ ,  $\beta_0 \leftarrow 0$ ,  $\epsilon_0 \leftarrow \|b\|_2$   
 4: **for**  $t \in \{0, \dots, M-1\}$  **do**  
      $\triangleright$  *compute the geodesic direction for each data point*  
 5:  $d_t \leftarrow \frac{\beta - \langle \beta, \beta_t \rangle \beta_t}{\|\beta - \langle \beta, \beta_t \rangle \beta_t\|_2}$ ,  $\forall n \in [N] : d_{t,n} \leftarrow \frac{\alpha_n - \langle \alpha_n, \beta_t \rangle \beta_t}{\|\alpha_n - \langle \alpha_n, \beta_t \rangle \beta_t\|_2}$   
      $\triangleright$  *choose the best geodesic*  
 6:  $n_t \leftarrow \operatorname{argmax}_{n \in [N]} \langle d_t, d_{t,n} \rangle$   
 7:  $\zeta_0 \leftarrow \langle \beta, a_{n_t} \rangle$ ,  $\zeta_1 \leftarrow \langle \beta, \beta_t \rangle$ ,  $\zeta_2 \leftarrow \langle a_{n_t}, \beta_t \rangle$   
      $\triangleright$  *compute the step size*  
 8:  $\gamma_t \leftarrow \frac{\zeta_0 - \zeta_1 \zeta_2}{(\zeta_0 - \zeta_1 \zeta_2) + (\zeta_1 - \zeta_0 \zeta_2)}$   
      $\triangleright$  *update the core-set*  
 9:  $w_{t+1} \leftarrow \frac{(1-\gamma_t)w_t + \gamma_t 1_{n_t}}{\|(1-\gamma_t)\beta_t + \gamma_t \alpha_{n_t}\|_2}$ ,  $\beta_{t+1} \leftarrow \frac{(1-\gamma_t)\beta_t + \gamma_t a_{n_t}}{\|(1-\gamma_t)\beta_t + \gamma_t \alpha_{n_t}\|_2}$   
 10:  $\epsilon_{t+1} \leftarrow \|A \cdot w_{t+1} - b\|_2$   
      $\triangleright$  *check improvement in numerical precision*  
 11: **if**  $\epsilon_{t+1} > \epsilon_t \cdot (1 + \varepsilon)$  **then**  
 12:      $w_M \leftarrow w_t$ ,  $\beta_M \leftarrow \beta_t$   
 13:     **break**  
 14: **end if**  
 15: **end for**  
      $\triangleright$  *scale the weights optimally*  
 16:  $\forall n \in [N] : (w_M)_n \leftarrow (w_M)_n \cdot \frac{\|b\|_2}{\|A \cdot e_n\|_2} \langle \beta, \beta_M \rangle$   
 17: **return**  $w_M$   
 18: **end procedure**

---

## 2.2.1 Generative Model

The prior and likelihood of the NIW model read

$$p_0(m, \Sigma) = p(m, \Sigma \mid \mu_0, \sigma_0, \Psi_0, \nu_0) = \mathcal{N}\left(m \mid \mu_0, \frac{\Sigma}{\sigma_0}\right) \mathcal{W}^{-1}(\Sigma \mid \Psi_0, \nu_0) \quad (1)$$

$$p(x \mid m, \Sigma) = \mathcal{N}(x \mid m, \Sigma) . \quad (2)$$

## 2.2.2 Core-Set Posterior

The conjugacy of the standard NIW model generalizes to the case of weighted observations: Denoting by  $|w| := \sum_n w_n$  and  $\bar{x} := \frac{1}{|w|} \sum_n w_n \cdot x_n$  the total weight and the weighted mean of observations, the sufficient statistics of the exact posterior  $p_1$  after observing the weighted data batch  $(w_n, x_n)_n$  are

$$\sigma_1 = \sigma_0 + |w| \quad \nu_1 = \nu_0 + |w| \quad \mu_1 = \frac{\sigma_0 \cdot \mu_0 + |w| \cdot \bar{x}}{\sigma_0 + |w|} \quad (3)$$

$$\Psi_1 = \Psi_0 + \frac{\sigma_0 \cdot |w|}{\sigma_0 + |w|} (\mu_0 - \bar{x})(\mu_0 - \bar{x})^T + \sum_n w_n (x_n - \bar{x})(x_n - \bar{x})^T . \quad (4)$$

This can be shown as follows: Using  $\langle x, y \rangle_\Sigma := x^T \Sigma^{-1} y$  and  $\|x\|_\Sigma^2 := \langle x, x \rangle_\Sigma$  for the inner product and norm induced by  $\Sigma$ ,

$$\begin{aligned} & \left( \frac{p_0(\cdot)}{p_1(\cdot)} \cdot \prod_n p(x_n \mid \cdot)^{w_n} \right) (m, \Sigma) \\ & \stackrel{(1),(2)}{=} \frac{\mathcal{N}\left(m \mid \mu_0, \frac{\Sigma}{\sigma_0}\right)}{\mathcal{N}\left(m \mid \mu_1, \frac{\Sigma}{\sigma_1}\right)} \cdot \frac{\mathcal{W}^{-1}(\Sigma \mid \Psi_0, \nu_0)}{\mathcal{W}^{-1}(\Sigma \mid \Psi_1, \nu_1)} \cdot \exp \left\{ \sum_n w_n \cdot \log \mathcal{N}(x_n \mid m, \Sigma) \right\} \\ & \stackrel{\mathcal{N}, \mathcal{W}}{\propto} \frac{\exp \left\{ -\frac{1}{2} \left[ \ln \left| \frac{\Sigma}{\sigma_0} \right| + \|m - \mu_0\|_{\Sigma/\sigma_0}^2 \right] \right\}}{\exp \left\{ -\frac{1}{2} \left[ \ln \left| \frac{\Sigma}{\sigma_1} \right| + \|m - \mu_1\|_{\Sigma/\sigma_1}^2 \right] \right\}} \cdot \frac{\exp \left\{ -\frac{1}{2} \left[ (\nu_0 + D + 1) \ln |\Sigma| + \text{Tr}(\Psi_0 \Sigma^{-1}) \right] \right\}}{\exp \left\{ -\frac{1}{2} \left[ (\nu_1 + D + 1) \ln |\Sigma| + \text{Tr}(\Psi_1 \Sigma^{-1}) \right] \right\}} \\ & \quad \cdot \exp \left\{ -\frac{1}{2} \left[ |w| \cdot \ln |\Sigma| + \sum_n w_n \|x_n - m\|_\Sigma^2 \right] \right\} \\ & \stackrel{(3)}{\propto} \exp \left\{ \text{Tr}([\Psi_1 - \Psi_0] \Sigma^{-1}) - \sum_n w_n \|x_n - m\|_\Sigma^2 - \sigma_0 \|m - \mu_0\|_\Sigma^2 + \sigma_1 \|m - \mu_1\|_\Sigma^2 \right\}^{\frac{1}{2}} \\ & \stackrel{(3)}{=} \exp \left\{ \text{Tr}([\Psi_1 - \Psi_0] \Sigma^{-1}) - |w| \|m\|_\Sigma^2 + 2 \langle m, |w| \bar{x} \rangle_\Sigma - \sum_n w_n \|x_n\|_\Sigma^2 - \sigma_0 \|m - \mu_0\|_\Sigma^2 \right. \\ & \quad \left. + (\sigma_0 + |w|) \left\| m - \frac{\sigma_0 \cdot \mu_0 + |w| \cdot \bar{x}}{\sigma_0 + |w|} \right\|_\Sigma^2 \right\}^{\frac{1}{2}} \\ & \stackrel{(5)}{=} \exp \left\{ \sigma_0 \left( \text{Tr}(\mu_0 \mu_0^T \Sigma^{-1}) - \|\mu_0\|_\Sigma^2 \right) - \sigma_1 \left( \text{Tr}(\mu_1 \mu_1^T \Sigma^{-1}) - \|\mu_1\|_\Sigma^2 \right) + \sum_n w_n \left( \text{Tr}(x_n x_n^T \Sigma^{-1}) - \|x_n\|_\Sigma^2 \right) \right\}^{\frac{1}{2}} \\ & \stackrel{\text{Tr}}{=} 1 , \end{aligned}$$

Experiment	$\alpha_\sigma$	$\alpha_\nu$	$\beta$
NIW	2	2	1.015
BLR	0.002	-	1
AR(1)	1.5	-	1

Table 1: Parameters for the rejuvenation kernels.

where the penultimate step uses the following identity:

$$\begin{aligned}
\Psi_1 - \Psi_0 &\stackrel{(4)}{=} \frac{\sigma_0 \cdot |w|}{\sigma_0 + |w|} (\mu_0 - \bar{x}) (\mu_0 - \bar{x})^T - |w| \cdot \bar{x} \cdot \bar{x}^T + \sum_n w_n \{x_n x_n^T + 2 \cdot \bar{x} \cdot \bar{x}^T - (x_n \bar{x}^T + \bar{x} x_n^T)\} \\
&= \frac{\sigma_0 \cdot |w|}{\sigma_0 + |w|} \left\{ \mu_0 \mu_0^T + \left(1 - \left(\frac{|w|}{\sigma_0} + 1\right)\right) \bar{x} \cdot \bar{x}^T - (\mu_0 \bar{x}^T + \bar{x} \mu_0^T) \right\} + \sum_n w_n \cdot x_n x_n^T \\
&= \sigma_0 \cdot \mu_0 \mu_0^T - \frac{1}{\sigma_0 + |w|} \left\{ \sigma_0^2 \cdot \mu_0 \mu_0^T + |w|^2 \cdot \bar{x} \cdot \bar{x}^T + \sigma_0 |w| (\mu_0 \bar{x}^T + \bar{x} \mu_0^T) \right\} + \sum_n w_n \cdot x_n x_n^T \\
&\stackrel{(3)}{=} \sigma_0 \cdot \mu_0 \mu_0^T - \sigma_1 \cdot \mu_1 \mu_1^T + \sum_n w_n \cdot x_n x_n^T
\end{aligned} \tag{5}$$

□

### 2.2.3 Convergence Metric

Posterior convergence in this experiment was measured using the maximum mean discrepancy,

$$\text{MMD}_k^2(\pi, \hat{\pi}) := \mathbb{E}_{x, x' \sim \pi} [k(x, x')] + \mathbb{E}_{y, y' \sim \hat{\pi}} [k(y, y')] - 2 \mathbb{E}_{x \sim \pi, y \sim \hat{\pi}} [k(x, y)],$$

a reproducing kernel Hilbert space metric which is well-suited for comparing distributions in different representations [Gretton et al., 2012]. The kernel  $k : \Theta \times \Theta \rightarrow [0, 1]$  in the space of Gaussian sufficient statistics was chosen as a Gaussian RBF, with a radius provided by the analytically computable Jeffrey divergence between model parameters,

$$k(m, \Sigma | \hat{m}, \hat{\Sigma}) \propto \exp \left\{ -\frac{1}{\alpha} \mathcal{D}_J \left( \mathcal{N}(m, \Sigma) \| \mathcal{N}(\hat{m}, \hat{\Sigma}) \right)^2 \right\},$$

and with a scale adjusted as  $\alpha = \frac{1}{2} \mathcal{D}_J(p_0 \| p_T)$  in terms of the prior  $p_0$  and the exact posterior  $p_T$ .

### 2.2.4 Rejuvenation Proposal

Given  $\sigma_t, \nu_t$  according to (3) for the SMC step  $t$ , the Metropolis-Hastings proposal kernel used for rejuvenation in the NIW experiments is constructed as follows:

$$q_t(\hat{m}, \hat{\Sigma} | m, \Sigma; \alpha_\sigma, \alpha_\nu, \beta) = \mathcal{N}\left(\hat{m} | m, \frac{\hat{\Sigma}}{\alpha_\sigma \sigma_t^\beta}\right) \mathcal{W}^{-1}\left(\hat{\Sigma} | (\alpha_\nu \nu_t^\beta - n - 1)\Sigma, \alpha_\nu \nu_t^\beta\right).$$

## 2.3 LOGISTIC REGRESSION

Unless stated otherwise in Section 5.3, we use the same setup for our Bayesian logistic regression experiment as in [Huggins et al., 2016]. The rejuvenation proposal is a Gaussian kernel centred at each particle,

$$q_t(\cdot | \Sigma; \alpha_\sigma, \beta) = \mathcal{N}\left(\cdot | 0, \frac{\Sigma}{\alpha_\sigma \sigma_t^\beta}\right), \tag{6}$$

where the tuning parameters are listed in Table 1, and clipping is introduced to ensure that  $\alpha_\sigma \sigma_t^\beta > 1.0$  for all  $t$ . To calibrate all SMC variants equally well,  $\Sigma$  is chosen by estimation from 10,000 posterior samples of a STAN [Carpenter et al., 2017] oracle. We set only two control parameters for STAN explicitly: `adapt_delta = 0.9` and `max_treedepth = 15`.