

Sequential Monte Carlo Methods for Bayesian Computation

A. Doucet

Feb. 2015

- SMC allows us to sample approximately from $\{p(x_{1:t} | y_{1:t})\}_{t \geq 1}$ and compute $\{p(y_{1:t})\}_{t \geq 1}$ for state-space models.
- SMC are useful far beyond this class of models
 - More general time series models
 - General “static” Bayesian computation
 - Some control problems
 - Rare event simulation
- We show here how to extend trivially SMC techniques previously discussed to a general class of problems.

Generic Problem

- Consider a sequence of probability distributions $\{\pi_t\}_{t \geq 1}$ defined on a sequence of spaces $\{\mathcal{X}_t\}_{t \geq 1}$ where $\mathcal{X}_t = \mathcal{X}^t$.
- Each distribution $\pi_t(x_{1:t})$ is assumed known *up to a normalizing constant*, i.e.

$$\pi_t(x_{1:t}) = \frac{\gamma_t(x_{1:t})}{Z_t}$$

where $\gamma_t : \mathcal{X}_t \rightarrow \mathbb{R}^+$ can be computed pointwise but Z_t cannot.

- SMC methods can be used to sample sequentially from $\{\pi_t\}_{t \geq 1}$ and compute Z_t .

Examples

- **State-space models**

$$\gamma_t(x_{1:t}) = p(x_{1:t}, y_{1:t}) = \mu(x_1) \prod_{k=2}^t f(x_k | x_{k-1}) \prod_{k=1}^t g(y_k | x_k),$$

$$Z_t = p(y_{1:t}) = \int \cdots \int \mu(x_1) \prod_{k=2}^t f(x_k | x_{k-1}) \prod_{k=1}^t g(y_k | x_k) dx_{1:t}$$

- **General time series models**

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

$$= \mu(x_1) \prod_{k=2}^t f(x_k | y_{1:k-1}, x_{1:k-1}) \prod_{k=1}^t g(y_k | y_{1:k-1}, x_{1:k})$$

$$Z_t = p(y_{1:t})$$

- **Nonparametric Bayes:** Wood & Griffiths, 2006; Caron & D., 2007, Saeedi & Bouchard-Côté, 2011; Ahmed & Smola, 2012 etc.
- **Graphical models, coalescent models, phylogenetic trees:** Gorur & Teh, 2008, Ihler, Frank, Smyth, 2009 etc.

- Risk sensitive control (Kantas, 2009):

$$\gamma_t(x_{1:t}) = \mu(x_1) \prod_{k=2}^t f(x_k | x_{1:k-1}) \prod_{k=1}^t \exp(C(x_k))$$

$$Z_t = \mathbb{E}_\mu \left[\exp \left(\sum_{k=1}^t C(x_k) \right) \right]$$

- Eigenmeasure/Eigenvalue (Heterington, 1984): $K(x, y) \geq 0$ with $\int K(x, y) dy \neq 1$, $\nu(y) = \lambda \int \nu(x) K(x, y) dx$

$$\gamma_t(x_{1:t}) = \mu(x_1) \prod_{k=2}^t K(x_{t-1}, x_t),$$

$$\pi_t(x_t) \rightarrow \nu(x_t) \text{ and } \frac{Z_{t+1}}{Z_t} \rightarrow \lambda$$

SMC for General Target Distributions

- Assume at time $t - 1$, we have $\{W_{t-1}^{(i)}, X_{1:t-1}^{(i)}\}$ approximating $\pi_{t-1}(x_{1:t-1})$.
- We introduce $q_t(x_t | x_{1:t-1})$ and use

$$\pi_t(x_{1:t}) = \frac{\gamma_t(x_{1:t})}{Z_t} = \frac{w_t(x_{1:t}) \pi_{t-1}(x_{1:t-1}) q_t(x_t | x_{1:t-1})}{Z_t}$$

where the incremental weight is

$$w_t(x_{1:t}) = \frac{\gamma_t(x_{1:t})}{\gamma_{t-1}(x_{1:t-1}) q_t(x_t | x_{1:t-1})}$$

- By sampling $X_t^{(i)} \sim q_t(x_t | x_{1:t-1})$, we obtain

$$\hat{\pi}_t(x_{1:t}) = \sum_{i=1}^N W_{1:t}^{(i)} \delta_{X_{1:t}^{(i)}}(x_{1:t}), \quad W_t^{(i)} \propto W_{t-1}^{(i)} w_t(X_{1:t}^{(i)})$$

- At time 1, you are interested in

$$\pi_1(x_1) = \frac{\gamma_1(x_1)}{Z_1}$$

- Introduce a distribution $q_1(x_1)$ easy to sample from and use

$$\pi_1(x_1) = \frac{w_1(x_1) q_1(x_1)}{Z_1}, \quad w_1(x_1) = \frac{\gamma_1(x_1)}{q(x_1)}$$

- Sample $X_1^{(i)} \sim q_1(x_1)$ thus $\hat{q}_1(x_1) = \frac{1}{N} \sum_{i=1}^N \delta_{X_1^{(i)}}(x_1)$ and

$$\hat{\pi}_1(x_1) = \sum_{i=1}^N W_1^{(i)} \delta_{X_1^{(i)}}(x_1), \quad W_1^{(i)} \propto w_1(X_1^{(i)})$$

A General SMC Algorithm

Assume we have N weighted particles $\{W_{t-1}^{(i)}, X_{1:t-1}^{(i)}\}$ approximating $\pi_{t-1}(x_{1:t-1})$ then at time t ,

- Sample $\tilde{X}_t^{(i)} \sim q_t(x_t | X_{1:t-1}^{(i)})$, set $\tilde{X}_{1:t} = (X_{1:t-1}^{(i)}, \tilde{X}_t^{(i)})$ and

$$\tilde{\pi}_t(x_{1:t}) = \sum_{i=1}^N W_{1:t}^{(i)} \delta_{X_{1:t}^{(i)}}(x_{1:t}),$$

$$W_t^{(i)} \propto W_{t-1}^{(i)} w_t(X_{1:t-1}^{(i)}, \tilde{X}_t^{(i)}),$$

$$\widehat{Z_t / Z_{t-1}} = \sum_{i=1}^N W_{t-1}^{(i)} w_t(X_{1:t-1}^{(i)}, \tilde{X}_t^{(i)})$$

- If $ESS < N/2$ resample $X_{1:t}^{(i)} \sim \tilde{\pi}_t(x_{1:t})$ and set $W_t^{(i)} \leftarrow \frac{1}{N}$ to obtain $\hat{\pi}_t(x_{1:t}) = \frac{1}{N} \sum_{i=1}^N \delta_{X_{1:t}^{(i)}}(x_{1:t})$.

- The proposal minimizing the variance of incremental weight

$$w_t(x_{1:t}) = \frac{\gamma_t(x_{1:t})}{\gamma_{t-1}(x_{1:t-1}) q_t(x_t | x_{1:t-1})} \propto \frac{\pi_t(x_{1:t})}{\pi_{t-1}(x_{1:t-1}) q_t(x_t | x_{1:t-1})}$$

is

$$q_t^{\text{opt}}(x_t | x_{1:t-1}) = \pi_t(x_t | x_{1:t-1}).$$

- In this case, we have

$$\begin{aligned} w_t(x_{1:t}) &= w_{t-1}(x_{1:t-1}) \frac{\gamma_t(x_{1:t})}{\gamma_{t-1}(x_{1:t-1}) q_t^{\text{opt}}(x_t | x_{1:t-1})} \\ &= w_{t-1}(x_{1:t-1}) \frac{\gamma_t(x_{1:t-1})}{\gamma_{t-1}(x_{1:t-1})}. \end{aligned}$$

Another Generic Sampling Problem

- Let $\{\pi_t\}_{t \geq 1}$ be a *sequence of probability distributions* defined on \mathcal{X} such that $\pi_t(x) = \pi_t(x)$ and each $\pi_t(x)$ is known up to a *normalizing constant*, i.e.

$$\pi_t(x) = \underbrace{Z_t^{-1}}_{\text{unknown}} \cdot \underbrace{\gamma_t(x)}_{\text{known}}.$$

- Standard SMC do not apply as all the distributions are defined on the same space.
- It is possible to adapt SMC to this problem.

- “A Good Monte Carlo is a Dead Monte Carlo”: Whenever you can integrate out analytically variables, do it.

- Example.** Consider $X_t = (U_t, Z_t)$ where

$$\begin{aligned} U_t | U_{t-1} &\sim f(\cdot | U_{t-1}), \\ Z_t &= A_{U_t} Z_{t-1} + B_{U_t} V_t, \quad V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_{t_v}), \\ Y_t &= C_{U_t} Z_t + D_{U_t} W_t, \quad W_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_{t_w}). \end{aligned}$$

- SIS can be used to approximate $p(x_{1:t} | y_{1:t})$ but this is inefficient as

$$p(x_{1:t} | y_{1:t}) = p(u_{1:t} | y_{1:t}) \underbrace{p(z_{1:t} | y_{1:t}, u_{1:t})}_{\text{Gaussian}}.$$

- Rao-Blackwellized SIS methods target

$p(u_{1:t} | y_{1:t}) \propto p(y_{1:t} | u_{1:t}) p(u_{1:t})$ instead of $p(x_{1:t} | y_{1:t})$ (D, Godsill & Andrieu, 2000; Liu & Chen, 2000). Also apply to dynamic Bayesian nets (Murphy, D., De Freitas, Russell, 2000).

Applications

- Sequential Bayesian Inference:* $\pi_t(x) = p(x | y_{1:t})$.
- Global optimization:* $\pi_t(x) \propto [\pi(x)]^{\eta_t}$ with $\{\eta_t\}$ increasing sequence such that $\eta_t \rightarrow \infty$.
- Sampling from a fixed target* π : $\pi_t(x) \propto \mu(x) [g(y|x)]^{\phi_t}$ where μ easy to sample and $\phi_t = 0$, $\phi_t > \phi_{t-1}$ and $\phi_P = 1$.
- Rare event simulation* $\pi(A) \ll 1$: $\gamma_t(x) = \pi(x) 1_{\mathcal{X}_t}(x)$ with Z_1 known, $\mathcal{X}_1 = \mathcal{X}$, $\mathcal{X}_t \subset \mathcal{X}_{t-1}$ and $\mathcal{X}_P = A$ then $Z_T = \pi(A)$.

- Run a MCMC (e.g. Metropolis-Hastings) algorithm to sample from each target distribution π_t ; i.e. build a Markov kernel $K_t(x'|x)$ such that

$$\pi_t(x') = \int_{\mathcal{X}} \pi_t(x) K_t(x'|x) dx$$

and simulate a Markov chain $\{X_t^{(i)}\}$: $X_t^{(1)} \sim \mu_t$ and $X_t^{(i)} \sim K_t(x'|X_t^{(i-1)})$.

- Under weak assumptions, we have $\lim_{i \rightarrow \infty} \mathcal{L}(X_t^{(i)}) = \pi_t$ i.e. $X_t^{(i)}$ is asymptotically distributed according to π_t and

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \varphi(X_t^{(i)}) = \int \varphi(x) \pi_t(x) dx.$$

Importance Sampling

- Let the *target distribution* be $\pi_t(x) = Z_t^{-1} \gamma_t(x)$ and μ_t be an *importance distribution* then

$$\pi_t(x) = \frac{w_t(x) q_t(x)}{\int w_t(x) q_t(x) dx} \text{ where } w_t(x) = \frac{\gamma_t(x)}{q_t(x)},$$

$$Z_t = \int w_t(x) q_t(x) dx$$

- By sampling N i.i.d. particles $X_t^{(i)} \sim q_t$ then $\hat{q}_t(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X_t^{(i)}}(x)$ and

$$\hat{\pi}_t(x) = \sum_{i=1}^N W_t^{(i)} \delta_{X_t^{(i)}}(x) \text{ where } W_t^{(i)} \propto w_t(X_t^{(i)}),$$

$$\hat{Z}_t = \frac{1}{N} \sum_{i=1}^N w_t(X_t^{(i)}).$$

- Convergence to π_t can be slow and is difficult to diagnose.
- It does not give an estimate of Z_t with 'good' properties.
- If π_{t-1} and π_t are 'close', then it should be possible to devise a cleverer strategy.
- We can use instead importance Sampling.

Limitations of Importance Sampling

- Importance Sampling (IS) is a straightforward method to use if q_t is easy to sample.
- For the estimates to have reasonable variances, we need to select very carefully the importance distribution.
- Naive strategies provide typically estimates with exponential variance in the dimension.
- For state-space models discussed previously, $\dim(\mathcal{X}) < 10$ in most cases. For static problems, we often have $\dim(\mathcal{X}) > 1000$.

- “Philosophy”: Start by doing simple things before trying to do complex things.
- Develop a sequential/iterative IS strategy where we start by approximating a simple target distribution π_1 . Then targets evolve over time and we *build the importance distribution sequentially*; i.e. at time t , we use q_{t-1} to build q_t .
- This approach only makes sense if the sequence $\{\pi_t\}$ is not arbitrary; i.e. π_{t-1} somewhat close to π_t .

- At time 1, sample N ($N \gg 1$) particles $X_1^{(i)} \sim q_1$ to obtain the following IS estimates

$$\hat{\pi}_1(x) = \sum_{i=1}^N W_1^{(i)} \delta_{X_1^{(i)}}(x) \text{ where } W_1^{(i)} \propto w_1(X_1^{(i)}),$$

$$\hat{Z}_1 = \frac{1}{N} \sum_{i=1}^N w_1(X_1^{(i)})$$

- *Remark*: Estimates have reasonable variance only if discrepancy between π_1 and μ_1 small; hence the need to start with easy to sample or approximate π_1 .

- At time $t - 1$, one has N particles $\{X_{t-1}^{(i)}, W_{t-1}^{(i)}\}$

$$X_{t-1}^{(i)} \sim q_{t-1}(x) \text{ and } W_{t-1}^{(i)} \propto \frac{\pi_{t-1}(X_{t-1}^{(i)})}{q_{t-1}(X_{t-1}^{(i)})}$$

- Move the particles according to transition kernel

$$X_t^{(i)} \sim K_t(x | X_{t-1}^{(i)}) \Rightarrow q_t(x') = \int q_{t-1}(x) K_t(x' | x) dx$$

- Optimal transition kernel $K_t(x' | x) = \pi_t(x')$ cannot be used so we need alternatives.

- $K_t(x' | x) = K_t(x')$ with
 - simple parametric form (e.g. Gaussian, multinomial etc.) (e.g. Cappé et al, 2005);
 - semi-parametric based on $\hat{\mu}_{t-1}(x)$ (e.g. West, 1993; Titterton, 2001)
- $K_t(x' | x)$ MCMC kernel of invariant distribution π_t .
 - burn-in correction by importance sampling (Gilks & Berzuini, 2001; Neal, 2001; Crisan & D., 2000).
- $K_t(x' | x)$ approximation of a Gibbs sampler of invariant distribution π_t .

- At time $t = 1$, sample $X_1^{(i)} \sim q_1(x)$ and set $w_1(X_1^{(i)}) = \frac{\gamma_1(X_1^{(i)})}{q_1(X_1^{(i)})}$.
- At time $t \geq 2$, sample $X_t^{(i)} \sim K_t(x|X_{t-1}^{(i)})$ and set $w_t(X_t^{(i)}) = \frac{\gamma_t(X_t^{(i)})}{q_t(X_t^{(i)})}$ where

$$q_t(x') = \int q_{t-1}(x) K_t(x'|x) dx.$$

- We have

$$\hat{\pi}_t(x) = \sum_{i=1}^N w_t^{(i)} \delta_{X_t^{(i)}}(x), \quad \hat{Z}_t = \frac{1}{N} \sum_{i=1}^N w_t(X_t^{(i)}).$$

An Artificial Target Distribution

- Problem summary:** It is impossible to compute pointwise $q_t(x_t)$ hence $\gamma_t(x_t) / q_t(x_t)$ except when $t = 1$.
- Solution:** Perform importance sampling on extended space.
- At time 2,

$$\frac{\pi_2(x_2)}{q_2(x_2)} = \frac{\pi_2(x_2)}{\int q_1(x_1) K_2(x_2|x_1) dx_1}$$
 cannot be evaluated

but alternative weights can be defined

$$\frac{\text{new joint target distribution}}{\text{joint importance distribution}} = \frac{\pi_2(x_2) L_1(x_1|x_2)}{q_1(x_1) K_2(x_2|x_1)}$$

where $L_1(x_1|x_2)$ is an *arbitrary* (backward) Markov kernel.

- "Proof" of validity:

$$\int \pi_2(x_2) L_1(x_1|x_2) dx_1 = \pi_2(x_2) \underbrace{\int L_1(x_1|x_2) dx_1}_{=1 \text{ whatever being } L_1} = \pi_2(x_2)$$

- In most cases, we *cannot* compute the marginal importance distribution

$$q_t(x_t) = \int q_{t-1}(x_{t-1}) K_t(x_t|x_{t-1}) dx_{t-1} \\ = \int q_1(x_1) \prod_{k=2}^t K_k(x_k|x_{k-1}) dx_{1:t-1}.$$

- Monte Carlo approximation is possible

$$\tilde{q}_t(x_t) = \int \hat{q}_{t-1}(x_{t-1}) K_t(x_t|x_{t-1}) dx_{t-1} = \frac{1}{N} \sum_{i=1}^N K_t(x|X_{t-1}^{(i)})$$

but is computationally intensive $\mathcal{O}(N^2)$.

An Extended Target Distribution Trick... Again

- Similarly at time t , $q_t(x_t)$ cannot be computed so perform importance sampling on an extended space between

$$\tilde{q}_t(x_{1:t}) = q_1(x_1) \prod_{k=2}^t K_k(x_k|x_{k-1})$$

and an extended artificial joint target distribution

$$\tilde{\pi}_t(x_{1:t}) = \pi_t(x_t) \prod_{k=1}^{t-1} L_k(x_k|x_{k+1})$$

where $\{L_k\}$ is an arbitrary sequence of "backward" Markov kernels.

- "Proof" of validity

$$\int \tilde{\pi}_t(x_{1:t}) dx_{1:t-1} = \pi_t(x_t) \underbrace{\int \prod_{k=1}^{t-1} L_k(x_k|x_{k+1}) dx_{1:t-1}}_{=1 \text{ whatever being } \{L_k\}} = \pi_t(x_t).$$

- By extending the integration space, the variance of the importance weights can only increase.
- The optimal kernels $\{L_k^{\text{opt}}\}$ are the ones bringing us back to the case where there is no space extension; i.e.

$$L_{t-1}^{\text{opt}}(x_{t-1}|x_t) = \frac{q_{t-1}(x_{t-1}) K_t(x_t|x_{t-1})}{q_t(x_t)}$$

- The result follows straightforwardly from the forward-backward formula for Markov processes

$$\tilde{q}_t(x_{1:t}) = q_1(x_1) \prod_{k=2}^t K_k(x_k|x_{k-1}) = q_t(x_t) \prod_{k=2}^t L_{t-1}^{\text{opt}}(x_{t-1}|x_t)$$

- L_{t-1}^{opt} cannot typically be computed (though there are important exceptions) but can be properly approximated in numerous cases.
- *Even if an approximation is used, the estimates are still asymptotically consistent.*

SMC Sampler

- At time $t = 1$, sample $X_1^{(i)} \sim \mu_1(\cdot)$ and set $w_1(X_1^{(i)}) = \frac{\gamma_1(X_1^{(i)})}{q_1(X_1^{(i)})}$.
- Resample $\{X_1^{(i)}, W_1^{(i)}\}$ to obtain new particles also denoted $\{X_1^{(i)}\}$
- At time $t \geq 2$
 - sample $X_t^{(i)} \sim K_t(X_{t-1}^{(i)}, x_t)$.
 - compute $w_t(X_{t-1}^{(i)}, X_t^{(i)}) = \frac{\gamma_t(X_t^{(i)}) L_{t-1}(X_{t-1}^{(i)}|X_t^{(i)})}{\gamma_{t-1}(X_{t-1}^{(i)}) K_t(X_t^{(i)}|X_{t-1}^{(i)})}$.
- Resample $\{X_t^{(i)}, W_t^{(i)}\}$ to obtain new particles also denoted $\{X_t^{(i)}\}$.

- We need to sample from a sequence of (artificial) target distributions $\{\tilde{\pi}_t\}$ of increasing dimension.
- **Conceptual difference:** Given $\{K_t\}$, $\{\tilde{\pi}_t\}$ has been constructed in a "clever" way such that

$$\int \tilde{\pi}_t(x_{1:t}) dx_{1:t-1} = \pi_t(x_t)$$

whereas usually the sequence of targets $\{\tilde{\pi}_t\}$ is fixed and $\{K_t\}$ is designed accordingly.

- Because we typically cannot use $\{L_k^{\text{opt}}\}$, the variance of the weights typically increases over time and it is necessary to resample.

Monte Carlo Estimates

- We obtain

$$\hat{\pi}_t(x) = \sum_{i=1}^N W_t^{(i)} \delta_{X_t^{(i)}}(x).$$

- Ratio of Normalizing Constants

$$\frac{Z_t}{Z_{t-1}} = \int \frac{\gamma_t(x_t) L_{t-1}(x_{t-1}|x_t)}{\gamma_{t-1}(x_{t-1}) K_t(x_t|x_{t-1})} \pi_{t-1}(x_{t-1}) K_t(x_t|x_{t-1}) dx_{t-1:t}$$

so

$$\widehat{\frac{Z_t}{Z_{t-1}}} = \sum_{i=1}^N W_{t-1}^{(i)} \frac{\gamma_t(X_t^{(i)}) L_{t-1}(X_{t-1}^{(i)}|X_t^{(i)})}{\gamma_{t-1}(X_{t-1}^{(i)}) K_t(X_t^{(i)}|X_{t-1}^{(i)})}.$$

- This is a generalization of the celebrated Jarzynski-Crooks identity (1997); see Neal (2001).

- **First step:** Build a sequence of distributions $\{\pi_t\}$ going from π_1 easy to sample/approximate to $\pi_P = \pi$; e.g. $\pi_t(x) \propto \mu(x) [g(y|x)]^{\phi_t}$ where μ easy to sample and $\phi_t = 0, \phi_t > \phi_{t-1}$ and $\phi_P = 1$.
- **Second step:** Introduce a sequence of transition kernels $\{K_t\}$; e.g. K_t MCMC sampler of invariant distribution π_t .
- **Third step:** Introduce a sequence of backward kernels $\{L_t\}$ equal/approximating L_t^{opt} ; e.g.

$$L_{t-1}(x_{t-1}|x_t) = \frac{\pi_{t-1}(x_{t-1})K_t(x_t|x_{t-1})}{\int \pi_{t-1}(x)K_t(x_t|x)dx}$$

$$\Rightarrow \frac{\gamma_t(x_t)L_{t-1}(x_{t-1}|x_t)}{\gamma_{t-1}(x_{t-1})K_t(x_t|x_{t-1})} = \frac{\gamma_t(x_t)}{\int \gamma_{t-1}(x)K_t(x_t|x)dx}$$

$$L_{t-1}(x_{t-1}|x_t) = \frac{\pi_t(x_{t-1})K_t(x_t|x_{t-1})}{\pi_t(x_t)}$$

$$\Rightarrow \frac{\gamma_t(x_t)L_{t-1}(x_{t-1}|x_t)}{\gamma_{t-1}(x_{t-1})K_t(x_t|x_{t-1})} = \frac{\gamma_t(x_{t-1})}{\gamma_{t-1}(x_{t-1})}$$

Experimental Setups

- We build the sequence of P distributions

$$\pi_t(\mu_{1:k}) \propto p(\mu_{1:k}) [p(y_{1:m}|\mu_{1:k}, \sigma_{1:k}, w_{1:k})]^{\phi_t}$$

where $\phi_1 = 0 < \phi_2 < \dots < \phi_P = 1$.

- MCMC sampler to sample from π_t : update $\mu_{1:4}$ via a MH kernel with additive normal random walk.
- We use

$$L_{t-1}(x_{t-1}|x_t) = \frac{\pi_t(x_{t-1})K_t(x_t|x_{t-1})}{\pi_t(x_t)}$$

- We have observations $y_{1:m}$ with

$$p(y|\mu_{1:k}, \sigma_{1:k}, w_{1:k}) = \sum_{i=1}^k w_i \mathcal{N}(y; \mu_i, \sigma_i)$$

- Assume that $k = 4, w_i = 1/k$ and $\sigma_i = 0.55$ are known, $\mu_{1:k}$ is uniform on the k -dimensional hypercube $[-10, 10]^k$.
- We simulate $m = 100$ observations for $\mu = \mu_{1:4} = (-3, 0, 3, 6)$ and want to sample from

$$p(\mu_{1:k}|y_{1:m}) \propto p(\mu_{1:k}) p(y_{1:m}|\mu_{1:k}, \sigma_{1:k}, w_{1:k}).$$

- Invariance of the posterior to permutation of the labels of the parameters gives it $k! = 24$ symmetric modes
- Basic random-walk MCMC and importance sampling methods fail.

Running Times SMC Samplers: CPU vs GPU

Table: Running times for the Sequential Monte Carlo Sampler for various values of N .

N	CPU (mins)	8800GT (secs)	Speedup	GTX280 (secs)	Speedup
8192	4.44	1.192	223.5	0.597	446
16384	8.82	2.127	249	1.114	475
32768	17.7	3.995	266	2.114	502
65536	35.3	7.889	268	4.270	496
131072	70.6	15.671	270	8.075	525
262144	141	31.218	271	16.219	522

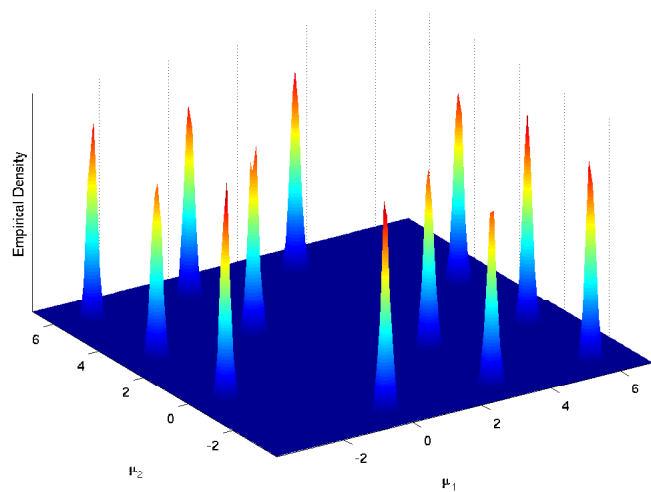


Figure: Estimated marginal posterior density $p(\mu_{1:2} | y_{1:m})$ from SMC samples, $N = 8192$.

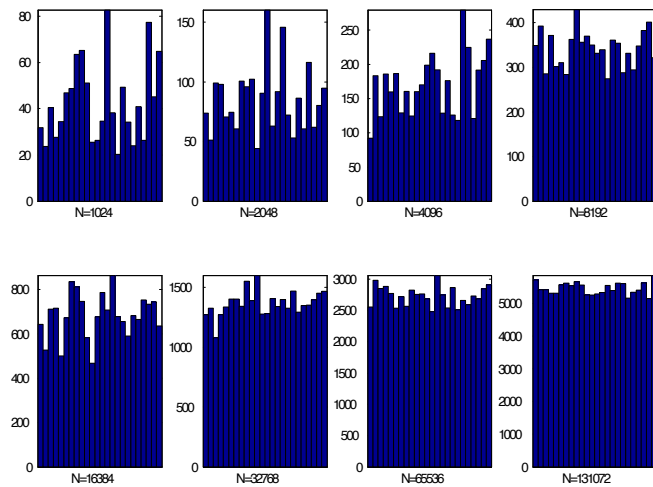


Figure: Effective number of SMC samples from each mode

- A powerful advantage of SMC samplers over MCMC is that they offer more flexibility in terms of adaptation compared to adaptive MCMC as no reliance on ergodicity.
- Adaptive schedule $\{\phi_t\}_{t \geq 1}$ can be built to ensure $ESS_t = \alpha ESS_{t-1}$.
- Adaptive proposals can be build; e.g. scaling of random walk can depend on $\{X_t^{(i)}\}$.
- Adaptive numbers of particles can be used.

- SMC methods can be used to sample from non-standard high-dimensional distributions.
- This is a powerful alternative/complement to MCMC useful in complex scenarios.
- Very easily parallelizable and GPU implementations already available.
- Adaptive strategies can easily be implemented without affecting convergence.