

Sequential Monte Carlo Methods

Seong-Hwan Jun

February 17, 2025

Bayesian Inference

Posterior distribution of $x \in \mathcal{X}$ given we observe $y \in \mathcal{Y}$:

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

- $p(y|x)$: Data likelihood.
- $p(x)$: Prior distribution.
- $p(y)$: Marginal likelihood.

Goal

Compute expectation of a function h w.r.t. posterior distribution:

$$I = \mathbb{E}_{X \sim p(x|y)}[h(X)] = \int h(x)p(x|y)dx,$$

where h is a known function that we can evaluate.

Example 1: $h(x) = x$,

$$\mathbb{E}_{X \sim p(x|y)}[X] = \int xp(x|y)dx.$$

Example 2: $h(x) = 1[x > a]$ for $x, a \in \mathbb{R}$,

$$\mathbb{E}_{X \sim p(x|y)}[1[X > a]] = \mathbb{P}(X > a|y) = \int 1[x > a]p(x|y)dx.$$

Solution I: Monte Carlo

If $p(x|y)$ is a known distribution such as Normal distribution, we can use Monte carlo approximation to compute the expectation no matter how complex h is:

$$\frac{1}{K} \sum_{k=1}^K h(x_k) \rightarrow \mathbb{E}_{X \sim p(x|y)}[h(X)] \text{ where } x_k \sim p(x|y) \text{ as } K \rightarrow \infty,$$

by the Law of Large Numbers (LLN).

Reality...

- In many cases, integral $p(y) = \int p(x, y)dx$ is not analytically available.
- Hence, $p(x|y)$ is not known in practice.

An important exception is if the likelihood and prior are *conjugate distributions*. In this case, $p(x|y)$ is known and can be sampled from.

Solution II: Importance Sampling

- Find a distribution q that is easy to sample from and $q(x) > 0$ whenever $p(x|y) > 0$.
- Propose $x_k \sim q(x)$ for $k = 1, \dots, K$.
- Let $w(x) = p(x|y)/q(x)$. Then,

$$\frac{1}{K} \sum_{k=1}^K w(x_k) h(x_k) \rightarrow I,$$

by LLN.

Derivation

$$\begin{aligned} I &= \int h(x)p(x|y)dx \\ &= \int h(x)\frac{p(x|y)}{q(x)}q(x)dx \\ &= \int h(x)w(x)q(x)dx \\ &= \mathbb{E}_{X \sim q}[h(X)w(X)]. \end{aligned}$$

Therefore, importance sampling is a Monte Carlo algorithm that approximates the expectation w.r.t to q with test function $h'(x) = h(x)w(x)$.

Self Normalization

Problem: We assumed that $p(x|y)$ can be evaluated. However, this requires computing $p(y)$,

$$p(y) = \int p(y|x)p(x)dx,$$

which is intractable in many settings. Therefore, $p(x|y) = p(x, y)/p(y)$ cannot be evaluated.

Self Normalization

- Let $\gamma(x) = p(x, y) = p(y|x)p(x)$ and $Z = p(y)$.
- Define weight function: $w(x) = \frac{\gamma(x)}{Zq(x)}$.
- Normalize: $\bar{w}_k = \frac{w(x_k)}{\sum_j w(x_j)}$.

Then,

$$\sum_{k=1}^K \bar{w}_k h(x_k) \rightarrow_p I.$$

Proof Sketch

$$\begin{aligned}\sum_{k=1}^K \bar{w}_k h(x_k) &= \sum_{k=1}^K \frac{w(x_k) h(x_k)}{\sum_j w(x_j)} \\ &= \sum_{k=1}^K h(x_k) \frac{\gamma(x_k)/Zq(x_k)}{\sum_j \gamma(x_j)/Zq(x_j)} \\ &= \frac{K^{-1} \sum_{k=1}^K h(x_k) \frac{\gamma(x_k)}{q(x_k)}}{K^{-1} \sum_{j=1}^K \frac{\gamma(x_j)}{q(x_j)}}.\end{aligned}$$

By LLN, the numerator converges to $\int h(x)p(x, y)dx$.

$$K^{-1} \sum_{k=1}^K h(x_k) \frac{\gamma(x_k)}{q(x_k)} \rightarrow \mathbb{E}_{X \sim q} \left[\frac{h(X)\gamma(X)}{q(X)} \right].$$

$$\begin{aligned} \text{RHS} &= \int h(x) \frac{\gamma(x)}{q(x)} q(x) dx \\ &= \int h(x) \gamma(x) dx \\ &= \int h(x) p(x, y) dx. \end{aligned}$$

Denominator converges to $Z = p(y)$.

$$K^{-1} \sum_{j=1}^K \frac{\gamma(x_j)}{q(x_j)} \rightarrow \mathbb{E}_{X \sim q} \left[\frac{\gamma(X)}{q(X)} \right].$$

$$\begin{aligned} \text{RHS} &= \int \frac{\gamma(x)}{q(x)} q(x) dx \\ &= \int p(x, y) dx \\ &= p(y). \end{aligned}$$

Therefore,

$$\sum_{k=1}^K \bar{w}_k h(x_k) \rightarrow \frac{\int h(x)p(x, y)dx}{p(y)} = \int h(x)p(x|y)dx.$$

To conclude the proof of consistency, we invoke continuous mapping theorem [Durrett, 2010, Thm 3.2.4]:

Let h be a measurable function and

$D_h = \{x : h \text{ is discontinuous at } x\}$. If $X_k \rightarrow_p X$ and $P(X \in D_h) = 0$, then $h(X_k) \rightarrow_p h(X)$.

Take $W_k = (X_k, Y_k)$, where $X_k = \frac{1}{K} \sum_{k=1}^K w_k \phi(x_{1:R}^k)$ and

$$Y_k = \left(\frac{1}{K} \sum_{k=1}^K w_k \right).$$

And, take $h(W_k) = X_k/Y_k$ to apply the continuous mapping theorem.

This proof sketch should also apply to almost sure convergence, if we invoke Strong Law of Large Numbers.

Example: Small Tail Probabilities

From [Robert and Casella, 2013, Example 3.11]. Importance sampling can be useful in many settings beyond Bayesian statistics.

Let $Z \sim N(0, 1)$. Estimate $\mathbb{P}(Z > 4.5)$.

- Solution 1: $Z_k \sim N(0, 1)$. Compute

$$\frac{1}{K} \sum_k 1[z_k > 4.5].$$

- Solution 2: $X_k \sim q = \text{Exponential}(0.1)$. Compute

$$\frac{1}{K} \sum_k 1[x_k > 4.5] \frac{\phi(x_k)}{q(x_k)}.$$

Note: We need to choose q to cover the region of interest (i.e., x for which $h(x)p(x) > 0$).

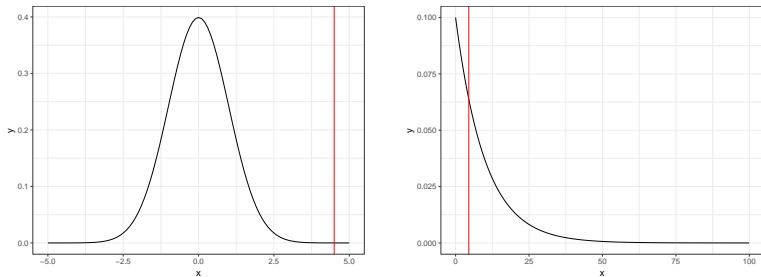


Figure: Left: Standard normal distribution. Right: $\text{Exp}(0.1)$. Red vertical line is the threshold, $a = 4.5$.

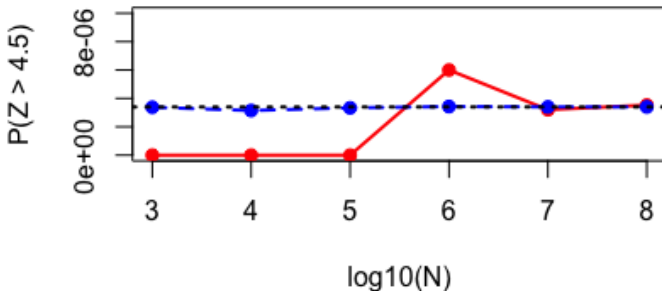


Figure: Monte Carlo estimate of $P(X > 4.5)$. Red: Simple Monte Carlo sampling. Blue dashed: Importance sampling. Dotted black: truth.

Brief Summary

- Monte Carlo sampling can be used to approximate complex integral numerically.
- IS can improve efficiency in terms of number of samples required.
- IS can be useful with just a simple q even when direct sampling is not possible.
- IS yields estimate of the marginal likelihood: $Z = p(y)$.
- Remember to choose q such that $q(x) > 0$ whenever $p(x) > 0$ (or $h(x)p(x) > 0$).
- Generally, we want to choose $q(x)$ to be similar to $p(x)$.

Sequential Importance Sampling

Now, let $\mathbf{x} = (x_1, \dots, x_d)$ be a d -dimensional vector.

Goal: Compute

$$I = \int h(x_1, \dots, x_d) p(x_1, \dots, x_d | \mathbf{y}) dx_1 \dots dx_d.$$

- Solution I: Importance sampling. We need to find a multivariate proposal distribution that is easy to sample from.

Sequential Importance Sampling

Now, let $\mathbf{x} = (x_1, \dots, x_d)$ be a d -dimensional vector.

Goal: Compute

$$I = \int h(x_1, \dots, x_d) p(x_1, \dots, x_d | \mathbf{y}) dx_1 \dots dx_d.$$

- Solution I: Importance sampling. We need to find a multivariate proposal distribution that is easy to sample from.
- If $x_i \in \mathbb{R}$, we may be able to use *multivariate Normal* distribution. But for general setting, finding q may be difficult.

Sequential Importance Sampling

Now, let $\mathbf{x} = (x_1, \dots, x_d)$ be a d -dimensional vector.

Goal: Compute

$$I = \int h(x_1, \dots, x_d) p(x_1, \dots, x_d | \mathbf{y}) dx_1 \dots dx_d.$$

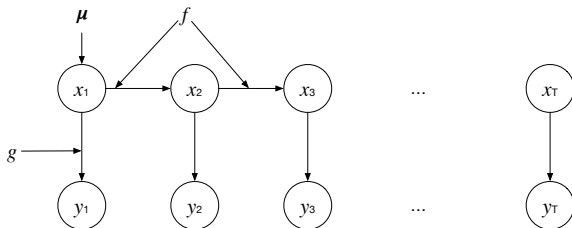
- Solution I: Importance sampling. We need to find a multivariate proposal distribution that is easy to sample from.
- If $x_i \in \mathbb{R}$, we may be able to use *multivariate Normal* distribution. But for general setting, finding q may be difficult.
- Also, curse of dimensionality: number of samples needed to sufficiently approximate the integral grows exponentially with dimension.

Idea: Propose one dimension at a time from $x_i \sim q_i$ for $i = 1, \dots, d$.

$$\begin{array}{ll} \text{Proposal:} & x_i^k \sim q_i \\ \text{Sample extension:} & \mathbf{x}_{1:i}^k = (x_1^k, \dots, x_{i-1}^k, x_i^k) \\ \text{Weight computation:} & w(x_i^k) = \frac{\gamma_i(x_{1:i}^k)}{\nu_i(x_{1:i}^k)}, \end{array}$$

where $\nu_i(x_{1:i}) = \prod_{j=1}^i q_j(x_j | x_{1:j-1})$ and $x_{i:j} = (x_i, \dots, x_j)$ for $0 < i < j$.

Application: Hidden Markov Model



$$x_1 \sim \mu(x_1)$$

$$x_t | x_{t-1} \sim f(x_t | x_{t-1}) \text{ for } t = 2, \dots, T$$

$$y_t | x_t \sim g(y_t | x_t) \text{ for } t = 1, \dots, T.$$

- $p(\mathbf{x}) = \mu(x_1) \prod_{t=2}^T f(x_t|x_{t-1})$
- $p(\mathbf{y}|\mathbf{x}) = \prod_{t=1}^T g(y_t|x_t)$
- $\gamma(\mathbf{x}) = p(\mathbf{x}, \mathbf{y}) = \mu(x_1)g(y_1|x_1) \prod_{t=2}^T f(x_t|x_{t-1})g(y_t|x_t)$.
- $Z = p(\mathbf{y}) = \int p(\mathbf{x}, \mathbf{y})d\mathbf{x}$

Recursive weight Update

$$\begin{aligned}w(x_{1:t}) &= \frac{\gamma_t(x_{1:t})}{\nu_t(x_{1:t})} \\ &= \frac{\gamma_{t-1}(x_{1:t-1})}{\nu_{t-1}(x_{1:t-1})} \frac{f(x_t|x_{t-1})g(y_t|x_t)}{q_t(x_t|x_{1:t-1})} \\ &= w(x_{1:t-1})\alpha(x_{1:t-1}, x_t).\end{aligned}$$

Therefore, we should store the weight from previous iteration and compute only the weight update function $\alpha(x_{1:t-1}, x_t)$ at current iteration.

Proposal

- Prior: $q_t = f(x_t|x_{t-1})$.
 - Weight function: $\alpha(x_{1:t-1}, x_t) = g(y_t|x_t)$.
 - Pro: Simplicity.
 - Con: May require large number of samples if $f(x_t|x_{t-1})$ differs significantly from $p(x_t|x_{1:t-1}, y_t)$.

- Adapted: $q_t = p(x_t|x_{1:t-1}, y_t)$

$$p(x_t|x_{1:t-1}, y_t) = \frac{p(x_t, y_t|x_{1:t-1})}{p(y_t|x_{1:t-1})} = \frac{g(y_t|x_t)f(x_t|x_{t-1})}{\int g(y_t|x_t)f(x_t|x_{t-1})dx_t}.$$

- Weight update function: $p(y_t|x_{1:t-1})^{-1}$.
- Pro: Makes use of the latest observation to build a smart proposal. Generally requires less number of samples compared to prior (for example, to attain similar accuracy of approximation).
- Con: Need to analytically compute $p(y_t|x_{1:t-1})$.

Example: Stochastic Volatility Model

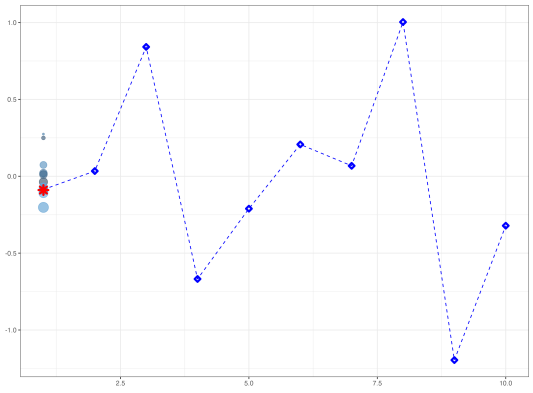
$$X_1 \sim \mathcal{N}(x_1|0, \sigma^2)$$

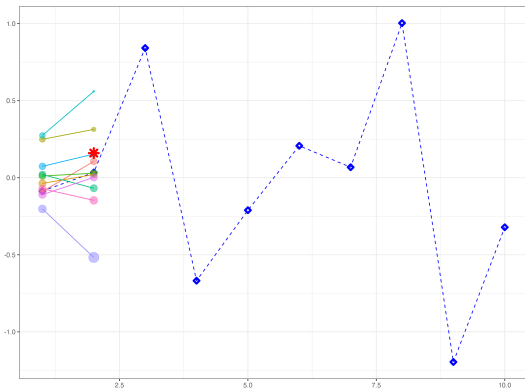
$$X_t|(X_{t-1} = x_{t-1}) \sim \mathcal{N}(x_t|\phi x_{t-1}, \sigma^2), \quad t = 2, \dots, T,$$

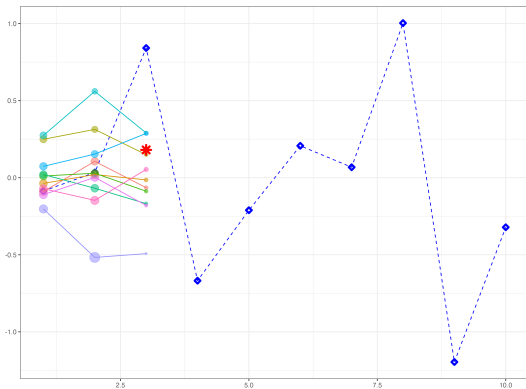
$$Y_t|(X_t = x_t) \sim \mathcal{N}(y_t|0, \beta^2 \exp(x_t)), \quad t = 2, \dots, T.$$

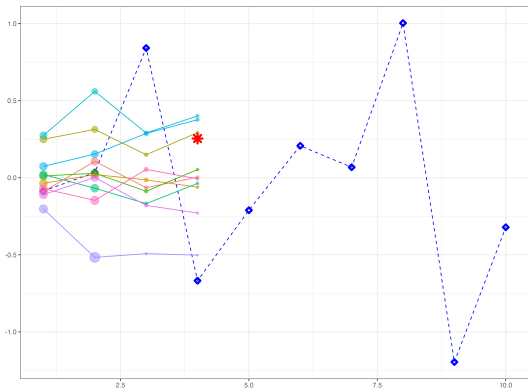
- X_t : Unobserved volatility of an asset (e.g., stock price).
- Y_t : Observed change in the price of the asset.

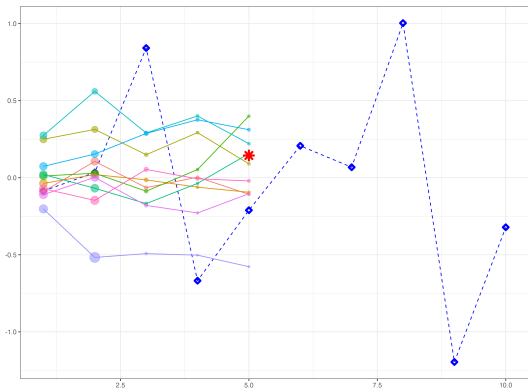
Illustration of SIS

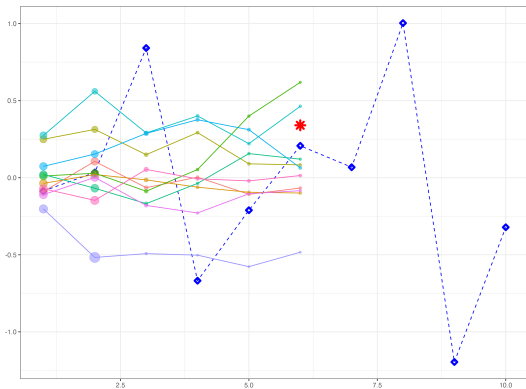


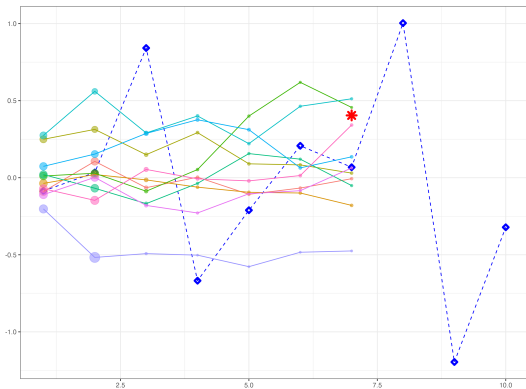


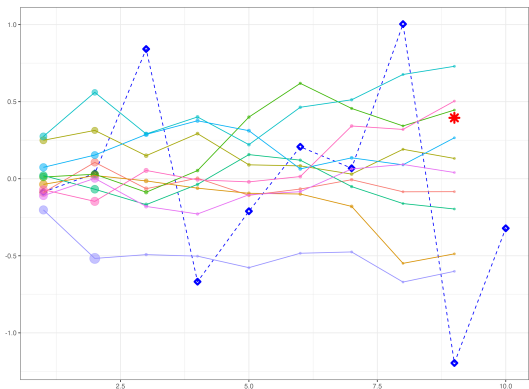




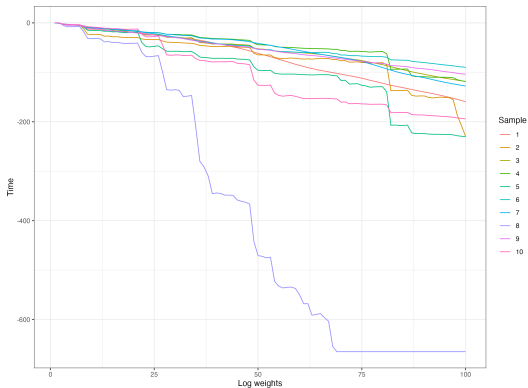








Weight degeneracy



Brief Summary

- SIS was originally designed for settings where we need to approximate high dimensional integral or perform imputation [Kong, Liu, and Wong. JASA, (1994)].
- Particularly useful if the model exhibits a temporal structure.
- Only need to find local (low-dimensional) proposal distributions.
- Weights decay with T . For large T , SIS usually does not work well (contradictory to the first point).
- Only a handful of samples become relevant as T increases, leading to waste of computational resources.

SIS with Resampling

- Idea: Interleave resampling step to choose promising particles.
- Use the weights to prune the particles.
- Sequential Monte Carlo methods refer to a class of algorithms that involve sequential proposal, weight computation, followed by (optional) resampling.
- Best tutorial to get started in SMC (in my opinion): [Doucet and Johansen, 2009].

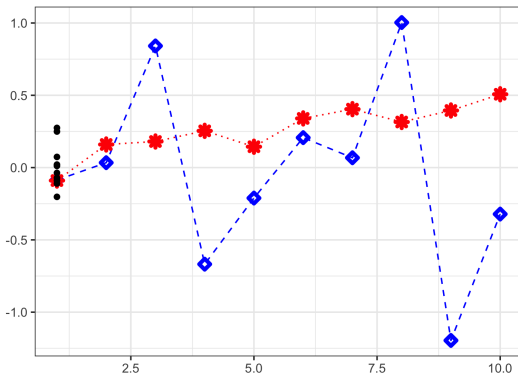
$t = 1$:

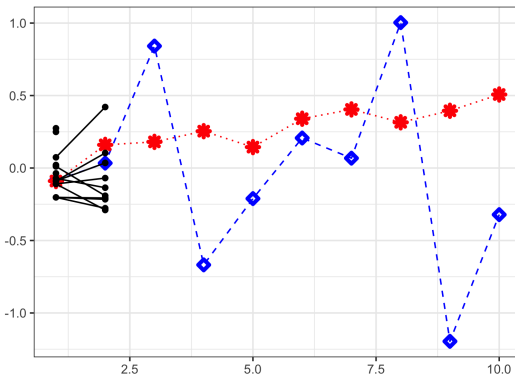
- Proposal: $x_1^k \sim q_1(x_1)$.
- Weight computation: $w(x_1^k) = \alpha(x_1^k)$.
- Weight normalization: $\bar{w}_1^k = w(x_1^k) / \sum_{k'} w(x_1^{k'})$.

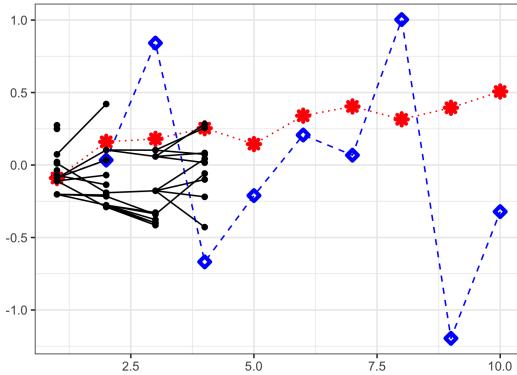
$t \geq 2$:

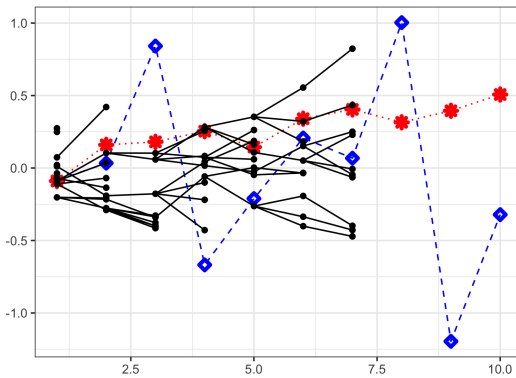
- Resampling: $j \sim \text{Multinomial}(\bar{w}_{t-1}^1, \dots, \bar{w}_{t-1}^K)$.
- Proposal: $x_t^k \sim q_t(x_t | x_{1:t-1}^j)$.
- Extension: $\mathbf{x}^k = (x_{1:t-1}^j, x_t^k)$.
- Weight computation: $w(x_{1:t}^k) = \alpha(x_{1:t-1}^j, x_t^k)$.
- Normalize the weights: $\bar{w}_t^k = w(x_{1:t}^k) / \sum_{k'} w(x_{1:t}^{k'})$.

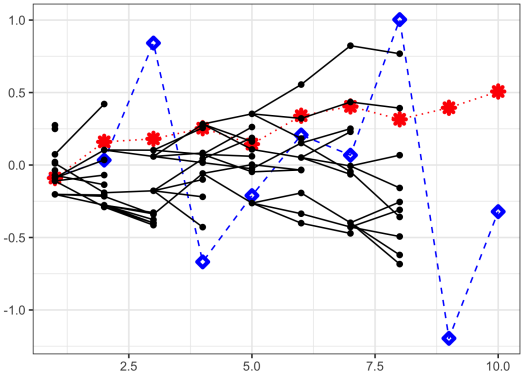
Illustration of SMC on SV Model











Filtering

- Samples and the weights can be used to approximate the *filtering* distribution:

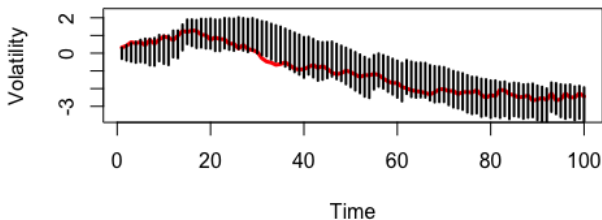
$$\hat{p}(x_t|y_{1:t}) = \sum_{k=1}^K \bar{w}_t^k \delta_{x_t^k}(x_t) \text{ for } t = 1, \dots, T.$$

or after resampling:

$$\hat{p}(x_t|y_{1:t}) = \frac{1}{K} \sum_{k=1}^K \delta_{x_t^k}(x_t) \text{ for } t = 1, \dots, T.$$

Effectiveness of SMC on SV Model

Ran with 10,000 particles. Computed empirical 95% confidence interval. Contains the true x_t about 93% of the time.



Predictive Distribution

- The generated samples can be used to build a predictive distribution:

$$p(x_{t+1}|y_{1:t}) = \int p(x_{t+1}|x_t)p(x_t|y_{1:t})dx_t.$$

Therefore, take the test function $h(x_{t+1}) = p(x_{t+1}|x_t)$ (e.g., $p(x_{t+1}|x_t) = f(x_{t+1}|x_t)$ in HMM application) and,

$$\hat{p}(x_{t+1}|y_{1:t}) = \sum_{k=1}^K p(x_{t+1}|x_t^k)\bar{w}_t^k\delta_{x_t^k}(x_t) \text{ for } t = 1, \dots, T.$$

Applications

- Online estimation: as the observation arrives, infer the latent state.
 - E.g., fraud detection, missile tracking, robot localization, etc.
- An extension of SMC [Del Moral et al., 2006], can be used in problems that do not exhibit temporal structure.
 - Phylogenetic inference [Bouchard-Côté et al., 2012].
 - Graph matching [Jun et al., 2017].
- Inference over graphical models [Naesseth et al., 2014].
- Probabilistic programming [Murray et al., 2017].

Resampling Algorithms

Can reduce variance of the estimator by using better resampling algorithms [Douc and Cappé, 2005]:

- Stratified Resampling.
- Residual Resampling.
- Systematic Resampling.
- Adaptive Resampling.

- A. Bouchard-Côté, S. Sankararaman, and M. I. Jordan. Phylogenetic inference via sequential monte carlo. *Systematic biology*, 61(4):579–593, 2012.
- P. Del Moral, A. Doucet, and A. Jasra. Sequential monte carlo samplers. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(3):411–436, 2006.
- R. Douc and O. Cappé. Comparison of resampling schemes for particle filtering. In *Image and Signal Processing and Analysis, 2005.*, pages 64–69. IEEE, 2005.
- A. Doucet and A. M. Johansen. A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of nonlinear filtering*, 12(656-704):3, 2009.
- R. Durrett. *Probability: theory and examples*. Cambridge university press, 2010.

- S.-H. Jun, S. W. Wong, J. Zidek, and A. Bouchard-Côté. Sequential graph matching with sequential monte carlo. In *Artificial Intelligence and Statistics*, pages 1075–1084, 2017.
- L. M. Murray, D. Lundén, J. Kudlicka, D. Broman, and T. B. Schön. Delayed sampling and automatic rao-blackwellization of probabilistic programs. *arXiv preprint arXiv:1708.07787*, 2017.
- C. A. Naesseth, F. Lindsten, and T. B. Schön. Sequential monte carlo for graphical models. In *Advances in Neural Information Processing Systems*, pages 1862–1870, 2014.
- C. Robert and G. Casella. *Monte Carlo statistical methods*. Springer Science & Business Media, 2013.