

# Importance sampling type estimators based on approximate marginal MCMC

---

Matti Vihola

Department of Mathematics and Statistics, University of Jyväskylä

Bayesian Computation for High-Dimensional Statistical Models, Singapore, 31 Aug 2018

Joint work with: Jouni Helske (Linköping), Jordan Franks (Jyväskylä),  
Ajay Jasra (Singapore) and Kody Law (Manchester)



Bayesian latent variable models

Some popular inference algorithms

Importance sampling type estimator based on marginal MCMC

Examples

Discussion

# Bayesian latent variable models

---



The following abstraction applies to many relevant statistical models:

- $\Theta$  is a vector of (hyper)parameters
- $\mathbf{X}$  is a vector of latent variables
- $\mathbf{Y}$  is a vector of observations

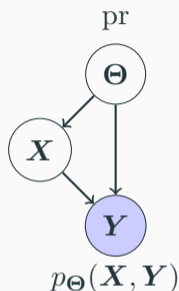
Only  $\mathbf{Y} = \mathbf{y}^*$  observed, both  $\Theta$  and  $\mathbf{X}$  are unknown.

The model is defined in terms of the following conditional laws:

- $\Theta \sim \text{pr}(\cdot)$ .
- $(\mathbf{X}, \mathbf{Y}) \mid \Theta \sim p_{\Theta}(\cdot)$   
(Often,  $p_{\Theta}(\mathbf{X}, \mathbf{Y}) = f_{\Theta}(\mathbf{X})g_{\Theta}(\mathbf{Y} \mid \mathbf{X})$ , but this is not relevant here.)

We are interested in the posterior of  $(\Theta, \mathbf{X})$  after observing  $\mathbf{Y} = \mathbf{y}^*$ :

$$\pi(\boldsymbol{\theta}, \mathbf{x}) = p(\mathbf{x}, \boldsymbol{\theta} \mid \mathbf{y}^*) \propto p(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y}^*) = \text{pr}(\boldsymbol{\theta})p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}^*).$$

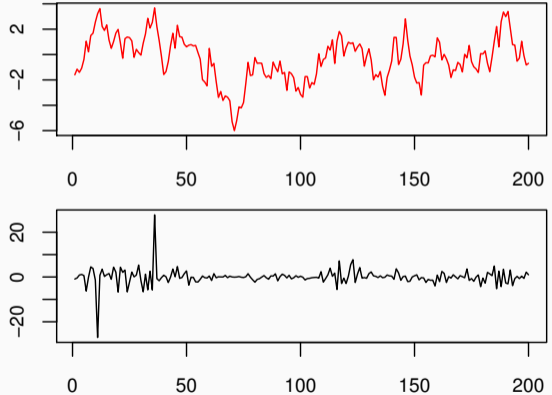




## Example of $p_{\Theta}(X, Y)$ : Stochastic volatility model

- $\Theta = (\phi, \sigma_x, \sigma_y)$
- $X = (X^{(1)}, \dots, X^{(T)})$   
stationary Gaussian AR(1) with  
parameters  $(\phi, \sigma_x^2)$ .
- The observations  
 $Y = (Y^{(1)}, \dots, Y^{(T)})$  are  
zero-mean Gaussian with  
 $\text{sd}(Y^{(t)}) = \sigma_y \exp(X^{(t)})$ ,

One realisation of  $X$  and  $Y$  with  
 $\theta = (0.9, 1, 2)$  and  $T = 200$ .



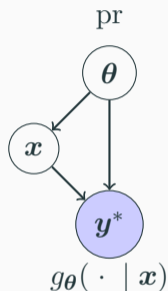


Typical scenario in a latent variable model:

- The hyperparameters  $\Theta$  are low-dimensional
  - $\dim(\Theta) = 3$  in the SV-example.
- The latent variables  $\mathbf{X}$  are high-dimensional
  - Often,  $\dim(\mathbf{X}) \propto \dim(\mathbf{Y})$ .
  - $\dim(\mathbf{X}) = 200$  in the SV-example.

Standard 'out-of-the-box' inference (e.g. using BUGS, Stan...):

- Simulate MCMC chain  $\mathbf{Z}_k = (\Theta_k, \mathbf{X}_k)$ , targeting  $\pi$ .
- High overall dimension & high correlations  
 $\implies$  inefficient



## Some popular inference algorithms

---



Consider the following factorisation of the posterior:

$$\pi(\boldsymbol{\theta}, \mathbf{x}) = \pi_m(\boldsymbol{\theta})r(\mathbf{x} | \boldsymbol{\theta}),$$

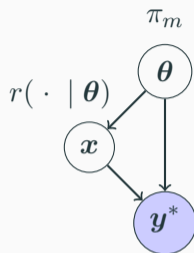
where the marginal posterior density and the corresponding conditional are given as follows:

$$\pi_m(\boldsymbol{\theta}) = \int \pi(\boldsymbol{\theta}, \mathbf{x})d\mathbf{x} \propto \text{pr}(\boldsymbol{\theta})L(\boldsymbol{\theta})$$

$$r(\mathbf{x} | \boldsymbol{\theta}) = \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}^*)}{L(\boldsymbol{\theta})}$$

with the marginal likelihood  $L$  taking the form:

$$L(\boldsymbol{\theta}) = \int p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}^*)d\mathbf{x}.$$







- $\theta$  low-dimensional, but  $\pi_m$  often non-standard  
     $\implies$  Non-parametric approximation, such as MCMC
- Problems (unless  $p_{\theta}(\mathbf{x}, \mathbf{y}^*)$  is of specific form such as Gaussian):
  - $L(\theta)$  is intractable.
  - $r(\mathbf{x} | \theta)$  is intractable.
- Two (successful branches of) solutions:
  - Approximate  $L(\theta)$  and  $r(\mathbf{x} | \theta)$  analytically.
  - Approximate  $L(\theta)$  and  $r(\mathbf{x} | \theta)$  using a specialised Monte Carlo algorithm.



(Rue, Martino & Chopin, 2009, *J. R. Stat. Soc. Ser. B. Stat. Methodol.*)

- Suppose  $p_{\theta}(\mathbf{x}, \mathbf{y}) = f_{\theta}(\mathbf{x})g_{\theta}(\mathbf{y})$  where  $f_{\theta}(\cdot)$  is Gaussian (&  $g_{\theta}$  of certain form)
- For any given  $\theta$ , use **Gaussian** (Laplace) **approximation**  $\hat{p}_{\theta}(\mathbf{x}, \mathbf{y}^*) \approx p_{\theta}(\mathbf{x}, \mathbf{y}^*)$   
 $\implies$  approximate likelihood  $L_a(\theta) = \int \hat{p}_{\theta}(\mathbf{x}, \mathbf{y}^*)d\mathbf{x}$
- Take a finite number of points  $(\theta_1, \dots, \theta_n)$  and approximate the full posterior as (something like)

$$\hat{\pi}(d\theta, d\mathbf{x}) = \frac{\sum_{k=1}^n w_k \text{Pr}(\theta_k) \hat{p}_{\theta_k}(\mathbf{x}, \mathbf{y}^*) \delta_{\theta_k}(d\theta) d\mathbf{x}}{\sum_{j=1}^n w_j \text{Pr}(\theta_j) L_a(\theta_j)}$$

where the weight  $w_k$  depends on the strategy how  $(\theta_i)$  chosen...

- (Further marginal corrections may be applied as well...)
- There is an **approximation error** (which does not vanish if  $n \rightarrow \infty$ )



(Andrieu, Doucet & Holenstein, 2010, *J. R. Stat. Soc. Ser. B. Stat. Methodol.*)

- For any given  $\theta$ , it is straightforward to generate random variables  $(V^{(i)}, \mathbf{X}^{(i)})$ , with  $V^{(i)} \geq 0$ , using particle filter (PF), which satisfy

$$\mathbb{E} \left[ \sum_{i=1}^m V^{(i)} \right] = L(\theta), \quad \text{and} \quad \mathbb{E} \left[ \sum_{i=1}^m V^{(i)} f(\mathbf{X}^{(i)}) \right] = \int p_{\theta}(\mathbf{x}, \mathbf{y}^*) f(\mathbf{x}) d\mathbf{x}.$$

- The algorithm:
  - Implement Metropolis-Hastings  $(\Theta_k)_{k \geq 1}$  targetting  $\pi_m$ , using  $\sum_{i=1}^m V^{(i)}$  in place of  $L(\theta)$ .
  - Construct an approximation of full  $\pi(\theta, \mathbf{x})$  using  $(V^{(i)}, \mathbf{X}^{(i)})$  above.



Particle marginal Metropolis-Hastings (PMMH) algorithm:

- Draw a new proposal  $\tilde{\Theta}_k \sim q(\Theta_{k-1}, \cdot)$
- Run PF with  $\theta = \tilde{\Theta}_k \rightarrow (\tilde{V}_k^{(i)}, \tilde{\mathbf{X}}_k^{(i)})$ .
- Accept and set  $(\Theta_k, V_k^{(i)}, \mathbf{X}_k^{(i)}) \leftarrow (\tilde{\Theta}_k, \tilde{V}_k^{(i)}, \tilde{\mathbf{X}}_k^{(i)})$  With probability

$$\min \left\{ 1, \frac{\text{pr}(\tilde{\Theta}_k) \left( \sum_{i=1}^m \tilde{V}_k^{(i)} \right) q(\tilde{\Theta}_k, \Theta_{k-1})}{\text{pr}(\Theta_{k-1}) \left( \sum_{i=1}^m V_{k-1}^{(i)} \right) q(\Theta_{k-1}, \tilde{\Theta}_k)} \right\};$$

otherwise reject and set  $(\Theta_k, V_k^{(i)}, \mathbf{X}_k^{(i)}) \leftarrow (\Theta_{k-1}, V_{k-1}^{(i)}, \mathbf{X}_{k-1}^{(i)})$



What is nice about this is that:

- This is valid MCMC, in the sense that

$$\frac{1}{n} \sum_{k=1}^n \frac{\sum_{i=1}^m V_k^{(i)} f(\Theta_k, \mathbf{X}_k)}{\sum_{i=1}^m V_k^{(i)}} \xrightarrow{n \rightarrow \infty} \int f(\boldsymbol{\theta}, \mathbf{x}) \pi(\boldsymbol{\theta}, \mathbf{x}) d\boldsymbol{\theta} d\mathbf{x}, \quad (\text{a.s.})$$

(under a minimal Harris recurrence assumption)

- $\implies$  PMCMC provides (asymptotically) **exact inference** (as  $n \rightarrow \infty$ ).

(NB: There is no asymptotic in ‘number of particles’  $m$ ...

... but  $m$  must be ‘large enough’ to make the MCMC mix sufficiently well...)



- (When applicable) INLA<sup>1</sup> is fast, and **often very accurate**
  - How accurate? How can you tell?
- PMCMC computationally demanding, but **exact** (asymptotically)
  - Does not require 'nearly Gaussian' structure of  $p_{\theta}(\mathbf{x}, \mathbf{y})$ .
  - Might need large  $m$  to work well  $\implies$  slow
  - Might still be 'sticky' (slower than geometric if  $\sum_{i=1}^m V^{(i)}$  are unbounded. . .)
- How about combining ideas both from INLA and PMCMC:
  - $\rightarrow$  Monte Carlo **correction/diagnosis** for INLA output, or
  - $\rightarrow$  Laplace approximations to **speed up** PMCMC. . .

---

<sup>1</sup>The same arguments hold with any other approximate scheme in place of INLA!



(Christen & Fox, 2005, *J. Comput. Graph. Statist.*)

Trick to make MCMC faster by using an approximation-based ‘screening’.

- Draw a new proposal  $\tilde{\Theta}_k \sim q(\Theta_{k-1}, \cdot)$
- With probability

$$\min \left\{ 1, \frac{\text{pr}(\tilde{\Theta}_k) L_a(\tilde{\Theta}_k) q(\tilde{\Theta}_k, \Theta_{k-1})}{\text{pr}(\Theta_{k-1}) L_a(\Theta_{k-1}) q(\Theta_{k-1}, \tilde{\Theta}_k)} \right\}$$

continue to the next step, otherwise reject.

- Run PF with  $\theta = \tilde{\Theta}_k \rightarrow (\tilde{V}_k^{(i)}, \tilde{X}_k^{(i)})$
- With probability

$$\min \left\{ 1, \frac{(\sum_{i=1}^m \tilde{V}_k^{(i)}) / L_a(\tilde{\Theta}_k)}{(\sum_{i=1}^m V_{k-1}^{(i)}) / L_a(\Theta_{k-1})} \right\}$$

accept, otherwise reject.

# Importance sampling type estimator based on marginal MCMC

---





(Review, consistency and CLTs: V, Helske, Franks, arXiv:1609.02541)

**Phase 1:** MCMC which targets the approximate marginal  $\pi_a(\boldsymbol{\theta}) \propto \text{pr}(\boldsymbol{\theta}L_a(\boldsymbol{\theta}))$

- Draw a new proposal  $\tilde{\boldsymbol{\Theta}}_k \sim q(\boldsymbol{\Theta}_{k-1}, \cdot)$
- With probability

$$\min \left\{ 1, \frac{\text{pr}(\tilde{\boldsymbol{\Theta}}_k)L_a(\tilde{\boldsymbol{\Theta}}_k)q(\tilde{\boldsymbol{\Theta}}_k, \boldsymbol{\Theta}_{k-1})}{\text{pr}(\boldsymbol{\Theta}_{k-1})L_a(\boldsymbol{\Theta}_{k-1})q(\boldsymbol{\Theta}_{k-1}, \tilde{\boldsymbol{\Theta}}_k)} \right\}$$

accept  $\boldsymbol{\Theta}_k = \tilde{\boldsymbol{\Theta}}_k$ ; otherwise reject  $\boldsymbol{\Theta}_k = \boldsymbol{\Theta}_{k-1}$ .

**Phase 2:** For  $k = 1, \dots, n$ , run PF with  $\boldsymbol{\theta} = \boldsymbol{\Theta}_k$   $\longrightarrow (V_k^{(i)}, X_k^{(i)})$  and calculate

$$E_n = \frac{\sum_{k=1}^n \sum_{i=1}^m W_k^{(i)} f(\boldsymbol{\Theta}_k, \mathbf{X}_k^{(i)})}{\sum_{j=1}^n \sum_{\ell=1}^m W_j^{(\ell)}} \quad \text{where} \quad W_k^{(i)} = \frac{V_k^{(i)}}{L_a(\boldsymbol{\Theta}_k)}.$$



- Phase 2 corrections entirely independent ('post-processing')  
⇒ parallelisable ⇒ scalable.
- Allows for calculating the correction only for accepted states ('jump chain')  
⇒ less expensive than DA
- Allows for (further) *thinning* before (expensive) correction  
⇒ further savings
- The approximate marginal MCMC ( $\Theta_k$ ) need not rely on estimators  
⇒ safer & easier to implement efficiently (e.g. adaptive MCMC...)
- The MCMC ( $\Theta_k$ ) need not be reversible  
⇒ new exciting non-reversible samplers readily applicable!
- Non-negativity of the estimator  $W_k$  not required  
⇒ allows for direct 'debiasing' tricks (or 'randomised multi-level Monte Carlo')  
(cf. Rhee & Glynn, *Oper. Res.* 2015; V, *Oper. Res.*, 2018)



## General setup & assumptions

General setup for IS *type* estimators based on approximate marginal MCMC:

- $\pi(\boldsymbol{\theta}, \mathbf{x}) = \pi_m(\boldsymbol{\theta})r(\mathbf{x} | \boldsymbol{\theta})$ .
- $\pi_m \lll \pi_a$
- $(\Theta_k)_{k \geq 1}$  MCMC **Harris ergodic** wrt  $\pi_a$
- $(\xi_k)_{k \geq 1}$  conditionally independent finite random signed measures given  $(\Theta_k)_{k \geq 1}$ , which form “**proper weighting**”:

General	The LVM example
$\pi_m(\boldsymbol{\theta})$	$\propto \text{pr}(\boldsymbol{\theta})L(\boldsymbol{\theta})$
$\pi_a(\boldsymbol{\theta})$	$\propto \text{pr}(\boldsymbol{\theta})L_a(\boldsymbol{\theta})$
$\xi_k(f)$	$\sum_{i=1}^m W_k^{(i)} f(\Theta_k, \mathbf{X}_k^{(i)})$

$$\mathbb{E}[\xi_k(\mathbf{1}) | \Theta_k = \boldsymbol{\theta}] = w_u(\boldsymbol{\theta}), \quad \text{where} \quad w_u(\boldsymbol{\theta}) = c_w \frac{\pi_m(\boldsymbol{\theta})}{\pi_a(\boldsymbol{\theta})}, \quad c_w > 0$$

$$\mathbb{E}[\xi_k(f) | \Theta_k = \boldsymbol{\theta}] = w_u(\boldsymbol{\theta}) \int r(\mathbf{x} | \boldsymbol{\theta}) f(\boldsymbol{\theta}, \mathbf{x}) d\mathbf{x}$$



- If  $\pi_a(m^{(1)}) < \infty$  where  $m^{(1)}(\boldsymbol{\theta}) = \mathbb{E}[|\xi_k(1)| + |\xi_k(f)| \mid \Theta_k = \boldsymbol{\theta}]$ , then

$$E_n = \frac{\sum_{k=1}^n \xi_k(f)}{\sum_{j=1}^n \xi_j(\mathbf{1})} \xrightarrow[\text{a.s.}]{n \rightarrow \infty} \pi(f) = \int f(\boldsymbol{\theta}, \mathbf{x}) \pi(\boldsymbol{\theta}, \mathbf{x}) d\boldsymbol{\theta} d\mathbf{x}.$$

- Suppose further that (for instance):

- $\pi_a(m^{(2)}) < \infty$  with  $m^{(2)} = \mathbb{E}[\xi_k(\bar{f})^2 \mid \Theta_k = \boldsymbol{\theta}]$  where  $\bar{f}(\boldsymbol{\theta}, \mathbf{x}) = f(\boldsymbol{\theta}, \mathbf{x}) - \pi(f)$ ,
- $(\Theta_k)_{k \geq 1}$  follows  $P$  which is aperiodic and reversible, with asymptotic variance  $\text{Var}(w_u \bar{f}^*, P) < \infty$ , where  $\bar{f}^*(\boldsymbol{\theta}, \mathbf{x}) = \int \bar{f}(\boldsymbol{\theta}, \mathbf{x}') r(\mathbf{x}' \mid \boldsymbol{\theta}) d\mathbf{x}'$ ,

Then,

$$\sqrt{n}[E_n - \pi(f)] \xrightarrow[\text{d}]{n \rightarrow \infty} N\left(0, \overbrace{\frac{\text{Var}(w_u \bar{f}^*, P)}{c_w^2}}^{\text{MCMC}} + \overbrace{\frac{\pi_a(v)}{c_w^2}}^{\text{IS corr}}\right),$$

where  $v(\boldsymbol{\theta}) = \text{Var}(\xi_k(\bar{f}) \mid \Theta_k = \boldsymbol{\theta})$ .



- The IS-type correction may be applied also when  $(\Theta_k, U_k)_{k \geq 1}$  is a pseudo-marginal chain arising from estimators  $\tilde{U}_\theta$  satisfying  $E[\tilde{U}_\theta] = L_a(\theta)$ .

In the pseudo-marginal case, consistency is more delicate:

- If  $\tilde{U}_\theta > 0$  a.s., then then we may always use  $W_k^{(i)} = V_k^{(i)} / U_k$ .
- When  $\mathbb{P}(\tilde{U}_\theta = 0)$  depends on  $\theta$ , this must be accounted for.
- For instance,  $V_k^{(i)}$  is constructed independent of  $U_k$ , then we must compensate for an extra factor  $p(\theta) = \mathbb{P}(U_\theta > 0) \dots \rightsquigarrow$  lazy ABC

## Examples

---



- State dynamics linear-Gaussian.
- Family of non-linear/non-Gaussian observation models.
- Approximate inference based on Laplace approximation (Durbin & Koopman, *Biometrika*, 1997).
- Proper weighting based on:
  - (i) Bootstrap particle filter (BSF).
  - (ii) Simple importance sampling & antithetic variables (SPDK, Shephard & Pitt, *Biometrika*, 1997; Durbin & Koopman, *Biometrika*, 1997)
  - (iii)  $\psi$ -auxiliary particle filter ( $\psi$ -APF: bootstrap PF for 'Laplace twisted model;' see Guarniero, Johansen & Lee, *JASA*, 2017)
- Simple IS-corrected estimator (IS1) or estimator based on jump chain (IS2).
- Compare against direct pseudo-marginal (PM) and delayed acceptance (DA).



(Stochastic volatility model with  $T = 5473$  observations, S&P index data. The numbers are 'inverse relative efficiencies' (avg. time (h)  $\times$  MSE) — lower is better.)

			BSF		SPDK				$\psi$ -APF			
	AI	AI <sup>G</sup>	IS2	IS2 <sup>8</sup>	PM	DA	IS1	IS2	PM	DA	IS1	IS2
Time	1.3	0.2	25.2	4.6	4.4	1.9	2.8	1.5	2.4	1.4	1.5	1.3
$\phi$	0.083	0.062	0.304	0.050	1.015	0.696	0.684	0.483	0.021	0.024	0.009	0.017
$\sigma_\eta$	0.726	0.298	0.483	0.096	3.090	3.307	0.603	0.710	0.044	0.055	0.016	0.028
$\nu$	0.008	0.747	0.287	0.042	1.208	2.544	0.228	0.404	0.026	0.027	0.010	0.020
$X_1$	0.133	0.035	0.321	0.071	3.054	1.883	0.346	0.373	0.029	0.026	0.007	0.018
$X_{5473}$	1.887	0.417	0.540	0.112	6.574	1.871	0.444	0.810	0.057	0.064	0.012	0.039





- 'Ideal' state dynamics follows a stochastic differential equation (SDE).
  - Cannot simulate exactly from the ideal transition.
  - Easy to simulate from time-discretised model (Euler, Milstein, ...).
  - The denser discretisation, the more simulation costs.
- Conditionally independent observations at discrete times.
- Approximate inference: particle marginal Metropolis-Hastings (PMMH) with 'coarse' (and cheap) time-discretisation.
- Correction with particle filter using 'fine' time-discretisation.

# Discretely observed (time-discretised) diffusion: empirical results



(Geometric Brownian motion observed at integer times, linear-Gaussian observations of log-state, Milstein discretisation. Parallel implementation with 48 cores, time mins.)

Init.	Mean						IRE				
	GT	Prior mean			Prior sample		Prior mean			Prior sample	
		DA	IS2	IS2 <sup>t</sup>	DA	IS2	DA	IS2	IS2 <sup>t</sup>	DA	IS2
Time	—	12.3	3.4	1.9	14.0	3.3	12.3	3.4	1.9	14.0	3.3
$\nu$	0.053	0.061	0.053	0.053	0.064	0.053	0.069	0.004	0.002	0.135	0.004
$\sigma_x$	0.253	0.278	0.253	0.253	0.251	0.252	0.576	0.029	0.019	0.336	0.022
$\sigma_y$	1.058	1.054	1.058	1.058	1.083	1.058	0.088	0.020	0.014	1.010	0.022
$X_1$	1.254	1.273	1.254	1.246	1.243	1.252	0.670	0.109	0.119	0.805	0.103
$X_{50}$	2.960	2.953	2.966	2.935	20.773	2.971	12.605	1.880	2.074	$4 \times 10^6$	2.308

# Discretely observed (time-discretised) diffusion: empirical results



(Geometric Brownian motion observed at integer times, linear-Gaussian observations of log-state, Milstein discretisation. Parallel implementation with 48 cores, time mins.)

Init.	Mean						IRE				
	GT	Prior mean			Prior sample		Prior mean			Prior sample	
		DA	IS2	IS2 <sup>t</sup>	DA	IS2	DA	IS2	IS2 <sup>t</sup>	DA	IS2
Time	—	12.3	3.4	1.9	14.0	3.3	12.3	3.4	1.9	14.0	3.3
$\nu$	0.053	0.061	0.053	0.053	0.064	0.053	0.069	0.004	0.002	0.135	0.004
$\sigma_x$	0.253	0.278	0.253	0.253	0.251	0.252	0.576	0.029	0.019	0.336	0.022
$\sigma_y$	1.058	1.054	1.058	1.058	1.083	1.058	0.088	0.020	0.014	1.010	0.022
$X_1$	1.254	1.273	1.254	1.246	1.243	1.252	0.670	0.109	0.119	0.805	0.103
$X_{50}$	2.960	2.953	2.966	2.935	20.773	2.971	12.605	1.880	2.074	$4 \times 10^6$	2.308

- Instead of correcting with ‘fine’ dynamics, it is possible to do IS correction for **ideal SDE dynamics** (Franks, Jasra, Law & V, arXiv:1807.10259, 2018).
- The correction is based on
  - Debiasing trick/randomised MLMC (Rhee & Glynn, *Oper. Res.*, 2015) with
  - ‘ $\Delta$ -PF’ (Jasra, Kamatani, Law & Zhou. *SIAM J. Sci. Comp.*, 2018).
- Detailed presentation in the closing workshop!

## Discussion

---

Is IS correction always better than DA (with same  $n$ )?



Is IS correction always better than DA (with same  $n$ )?



No.

# Is IS correction always better than DA (with same $n$ )?



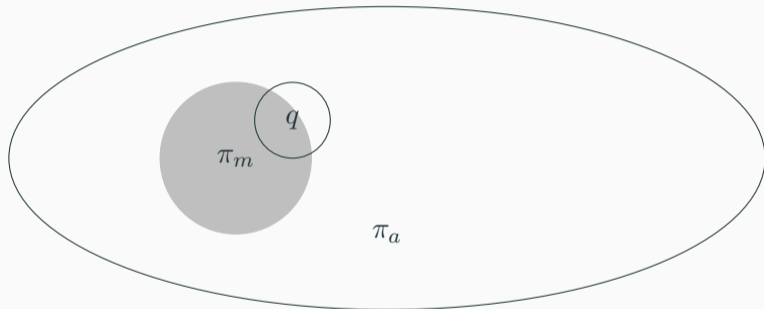
No.

More details Wed 12 Sep at 4pm...





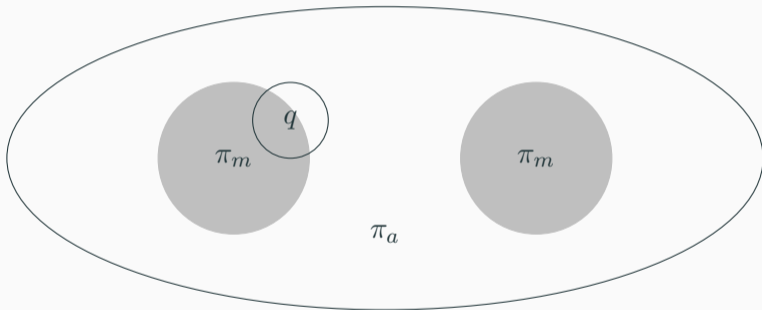
(Franks, V: arXiv:1706.09873)



DA better than IS:  $\pi_m$  and  $\pi_a$  are uniform,  $q$  uniform random walk. Approximate chain spends a lot of time outside the support of  $\pi_m$ .



(Franks, V: arXiv:1706.09873)



IS better than DA chain, which is reducible (cannot switch mode of  $\pi_m$ ).



- In practice, we have  $\pi_a \approx \pi_m$ , which is clearly not the case in the examples above.
- Empirical results suggest that IS often improves on DA slightly

(Franks, V: arXiv:1706.09873):

- If  $c_w^{-1}W_k \leq C$  a.s., then

$$\text{Var}(\text{IS}) \leq C\text{Var}(\text{DA}) + \bar{\pi}(\xi^2[C - c_w^{-1}W])$$

where  $\bar{\pi}$  corresponds to the stationary distribution of the DA chain.

$\rightsquigarrow$  With parallelisation, IS might be a better choice. . .

- NB: In the LVM setting we may modify the likelihood approximation:
  - $L_a(\boldsymbol{\theta}) \rightarrow L_a(\boldsymbol{\theta}) + \epsilon$

This leads to bounded weights if the likelihood estimators are bounded.



- If there is an approximation available, use it!
- IS type correction is a natural way to use the approximation
  - May be a useful alternative to DA pseudo-marginal algorithm (because of the several possible advantages)...
  - ... but not guaranteed to be uniformly better
- Our contributions:
  - arXiv:1609.02541: Review, consistency/CLT results; application in the state-space context, using Laplace approximation and coarse discretisation of diffusion model
  - arXiv:1706.09873: Theoretical bounds relating the efficiencies of IS/DA
  - arXiv:1807.10259: Full inference of SDE driven HMM based on randomised MLMC.
- Ongoing work:
  - Application beyond the state-space context.
  - Some insights for ABC-MCMC...

## References

- **C. Andrieu, A. Doucet and R. Holenstein.**  
Particle Markov chain Monte Carlo methods.  
*J. R. Stat. Soc. Ser. B. Stat. Methodol.* 72(3), 269–342, 2010.
- **J. A. Christen and C. Fox.**  
Markov chain Monte Carlo using an approximation.  
*J. Comput. Graph. Statist.* 14(4), 795-810, 2005.
- **J. Franks, A. Jasra, K. Law and M. Vihola.**  
Unbiased inference for discretely observed hidden Markov model diffusions  
*arXiv:1807.10259*, 2018.
- **J. Franks and M. Vihola.**  
Importance sampling correction versus standard averages of reversible MCMCs in terms of the asymptotic variance  
*arXiv:1706.09873*, 2017.
- **M. Vihola, J. Helske and J. Franks.**  
Importance sampling type estimators based on approximate marginal MCMC.  
*arXiv:1609.02541*, 2016.

Thank you!