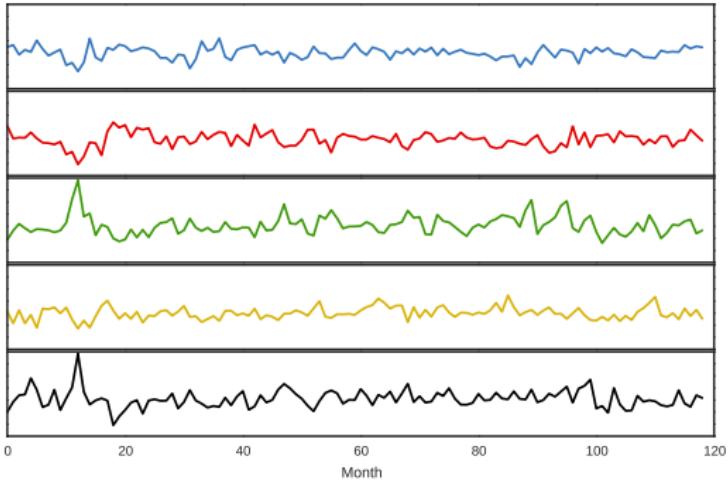


# Approximate Bayesian inference

## Variational and Monte Carlo methods

Christian A. Naesseth

# Exchange rate data



# Image data



- 1 Bayesian inference
- 2 Variational inference
- 3 Stochastic gradient variational inference
- 4 Rejection sampling variational inference
- 5 Variational sequential Monte Carlo

# Bayesian inference

- $\mathbf{y}$  - data
- $\mathbf{x}$  - latent (unknown) variables
- $p(\mathbf{x}, \mathbf{y})$  - probabilistic (generative) model

Posterior distribution:

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

Expectations:

$$\mathbb{E}_{p(\mathbf{x} | \mathbf{y})} [\varphi(\mathbf{x})] = \int \varphi(\mathbf{x}) p(\mathbf{x} | \mathbf{y}) d\mathbf{x}$$

# Examples of probabilistic generative models

- *Exchangeable data models*

$$p(\mathbf{x}, \mathbf{y}_{1:T}) = \mu(\mathbf{x}) \prod_{t=1}^T g(\mathbf{y}_t | \mathbf{x})$$

Applications: Regression, classification, ...

- *State space models*

$$p(\mathbf{x}_{1:T}, \mathbf{y}_{1:T}) \triangleq \mu(\mathbf{x}_1) \prod_{t=2}^T f(\mathbf{x}_t | \mathbf{x}_{t-1}) \prod_{t=1}^T g(\mathbf{y}_t | \mathbf{x}_t)$$

Applications: Filtering, smoothing, forecasting, ...

# Approximate Bayesian inference

Because  $p(\mathbf{x} | \mathbf{y})$  is typically intractable we approximate:

*Variational methods* - Turn integration into optimization

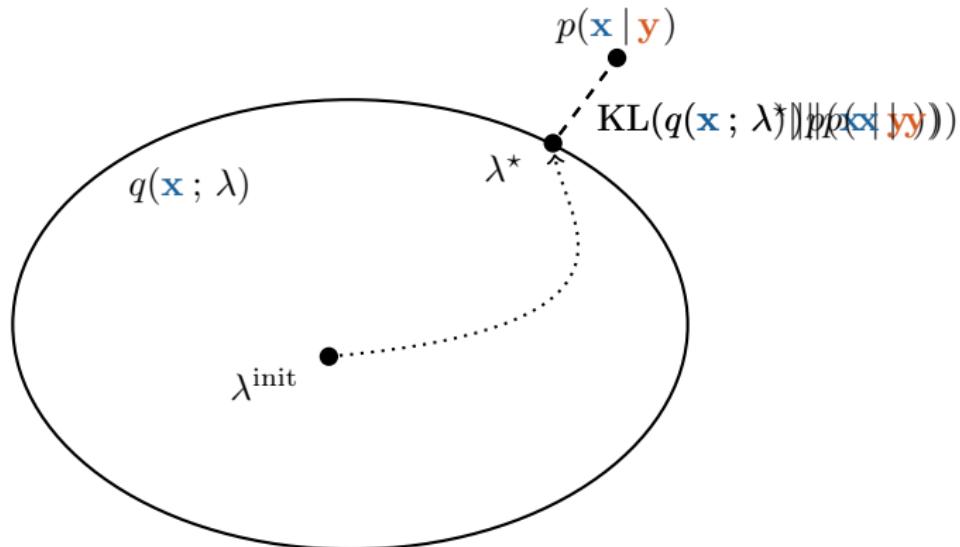
$$\mathbb{E}_{p(\mathbf{x} | \mathbf{y})} [\varphi(\mathbf{x})] \approx \mathbb{E}_{q(\mathbf{x}; \lambda^*)} [\varphi(\mathbf{x})]$$

*Monte Carlo methods* - Random sampling

$$\mathbb{E}_{p(\mathbf{x} | \mathbf{y})} [\varphi(\mathbf{x})] \approx \frac{1}{N} \sum_{i=1}^N \varphi(\mathbf{x}^i), \quad \mathbf{x}^i \sim p(\mathbf{x} | \mathbf{y})$$

- 1 Bayesian inference
- 2 **Variational inference**
- 3 Stochastic gradient variational inference
- 4 Rejection sampling variational inference
- 5 Variational sequential Monte Carlo

# Variational inference



# Evidence lower bound

$$\begin{aligned} \text{KL}(q(\mathbf{x}; \lambda) \| p(\mathbf{x} | \mathbf{y})) &= \log p(\mathbf{y}) - \mathbb{E}_q [\log p(\mathbf{x}, \mathbf{y}) - \log q(\mathbf{x}; \lambda)] \\ \Rightarrow \log p(\mathbf{y}) &\geq \underbrace{\mathbb{E}_q [\log p(\mathbf{x}, \mathbf{y}) - \log q(\mathbf{x}; \lambda)]}_{\triangleq \mathcal{L}(\lambda)} \end{aligned}$$

- Maximizing the *ELBO*  $\mathcal{L}(\lambda)$  minimizes KL divergence between  $q(\mathbf{x}; \lambda)$  and  $p(\mathbf{x} | \mathbf{y})$
- *Coordinate ascent* when conditionally conjugate
- *Stochastic optimization* using noisy  $\nabla_\lambda \mathcal{L}(\lambda)$  estimated with Monte Carlo methods

# Bayesian logistic regression

- Data are pairs  $(\mathbf{y}_t, \mathbf{z}_t)$ 
  - $\mathbf{z}_t$  is a covariate
  - $\mathbf{y}_t$  is binary label
  - $\theta$  is regression coefficient
- The generative model is

$$p(\theta, \mathbf{y}_{1:T} | \mathbf{z}_{1:T}) = \mathcal{N}(\theta; 0, 1) \prod_{t=1}^T \text{Bernoulli}(\mathbf{y}_t; \sigma(\theta \mathbf{z}_t)),$$

where  $\sigma(\cdot)$  is the logistic function mapping reals to  $(0, 1)$ .

# VI for Bayesian logistic regression

- Let's consider a single data pair  $(\mathbf{y}, \mathbf{z})$
- Variational approximation  $q(\boldsymbol{\theta}; \lambda) = \mathcal{N}(\boldsymbol{\theta}; \mu, \sigma^2)$
- The ELBO is given by

$$\mathcal{L}(\mu, \sigma^2) = \mathbb{E}_q [\log p(\boldsymbol{\theta}) + \log p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{z}) - \log q(\boldsymbol{\theta}; \mu, \sigma^2)]$$

# VI for Bayesian logistic regression

$$\begin{aligned}
 \mathcal{L}(\mu, \sigma^2) &= \mathbb{E}_q [\log p(\boldsymbol{\theta}) + \log p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{z}) - \log q(\boldsymbol{\theta}; \mu, \sigma^2)] \\
 &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbb{E}_q [\log p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{z})] + \text{const.} \\
 &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbb{E}_q [\mathbf{y}\mathbf{z}\boldsymbol{\theta} - \log(1 + \exp(\mathbf{z}\boldsymbol{\theta}))] \\
 &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbf{y}\mathbf{z}\mu - \mathbb{E}_q [\log(1 + \exp(\mathbf{z}\boldsymbol{\theta}))]
 \end{aligned}$$

We are stuck because of the expectation. What can we do?

- 1 Bayesian inference
- 2 Variational inference
- 3 Stochastic gradient variational inference**
- 4 Rejection sampling variational inference
- 5 Variational sequential Monte Carlo

# Stochastic Optimization

*Gradient-based optimization:*

$$\nabla \mathcal{L}(\lambda) = 0$$

*Stochastic approximation:* [Robbins and Monro, 1951]

$$\lambda^{k+1} = \lambda^k + \alpha^k \widehat{\nabla} \mathcal{L}(\lambda^k),$$

with  $\mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda^k)] = \nabla \mathcal{L}(\lambda^k)$ .

## VI for Bayesian logistic regression

$$\mathcal{L}(\mu, \sigma^2) = -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbf{y} \mathbf{z} \mu - \mathbb{E}_q [\log (1 + \exp(\mathbf{z} \boldsymbol{\theta}))]$$

$$\begin{aligned}\nabla \mathbb{E}_q [\log (1 + \exp(\mathbf{z} \boldsymbol{\theta}))] &= \int \nabla q(\boldsymbol{\theta}; \mu, \sigma^2) \log (1 + \exp(\mathbf{z} \boldsymbol{\theta})) d\boldsymbol{\theta} \\ &= \mathbb{E}_q [\nabla \log q(\boldsymbol{\theta}; \mu, \sigma^2) \log (1 + \exp(\mathbf{z} \boldsymbol{\theta}))] \\ &\approx \frac{1}{N} \sum_{i=1}^N \nabla \log q(\boldsymbol{\theta}^i; \mu, \sigma^2) \log (1 + \exp(\mathbf{z} \boldsymbol{\theta}^i)) \triangleq \widehat{\nabla} \mathcal{L}(\mu, \sigma^2),\end{aligned}$$

with  $\boldsymbol{\theta}^i \sim q(\boldsymbol{\theta}; \mu, \sigma^2)$ . Unbiased MC estimator of the gradient!

# Stochastic gradient variational inference

Monte Carlo estimates of gradients to optimize the ELBO.

*Score function:* [Paisley et al., 2012, Ranganath et al., 2014]

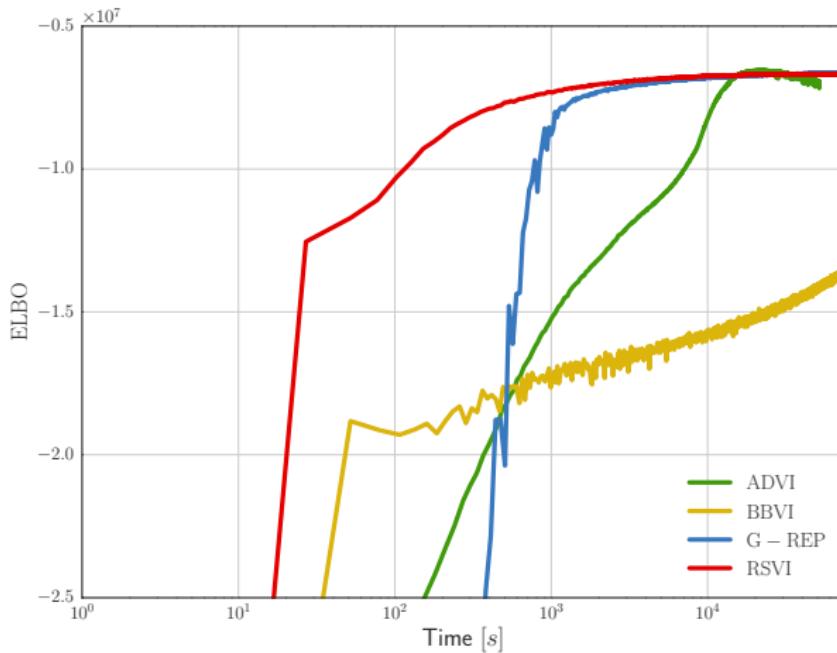
$$\begin{aligned}\nabla \mathcal{L}(\lambda) &= \nabla \mathbb{E}_q [\log p(\mathbf{x}, \mathbf{y}) - \log q(\mathbf{x}; \lambda)] \\ &= \mathbb{E}_q [\nabla \log q(\mathbf{x}; \lambda) \cdot (\log p(\mathbf{x}, \mathbf{y}) - \log q(\mathbf{x}; \lambda))]\end{aligned}$$

*Reparameterization:* [Kingma & Welling 2014, Rezende et al, 2014]

$$\begin{aligned}\nabla \mathcal{L}(\lambda) &= \nabla \mathbb{E}_q [\log p(\mathbf{x}, \mathbf{y}) - \log q(\mathbf{x}; \lambda)] \\ &= \nabla \mathbb{E}_{s(\varepsilon)} [\log p(h(\varepsilon, \lambda), \mathbf{y}) - \log q(h(\varepsilon, \lambda); \lambda)] \\ &= \mathbb{E}_{s(\varepsilon)} [\nabla_{\mathbf{x}} (\log p(\mathbf{x}, \mathbf{y}) - \log q(\mathbf{x}; \lambda)) \nabla_{\lambda} h(\varepsilon, \lambda)]\end{aligned}$$

e.g.  $\mathbf{x} = h(\varepsilon, \mu, \sigma) = \mu + \sigma \varepsilon$ ,  $s(\varepsilon) = \mathcal{N}(0, 1)$ .

## Reparameterization vs. score function



# Why doesn't reparameterization always work?

*All random variables we simulate on our computers are ultimately transformations of uniforms.*

Why can't we do reparameterization gradients for everything?

Often transformations  $h(\cdot)$  used in practice are *not differentiable*.

- 1 Bayesian inference
- 2 Variational inference
- 3 Stochastic gradient variational inference
- 4 Rejection sampling variational inference**
- 5 Variational sequential Monte Carlo

# Reparameterizing gamma

- **Gamma( $\lambda, 1$ ):** Common distribution in VI
- Used to sample beta, Dirichlet, Student's t, chi-squared, ...
- *Simulation:*
  - Propose

$$h(\varepsilon, \lambda) = \left( \lambda - \frac{1}{3} \right) \left( 1 + \frac{\varepsilon}{\sqrt{9\lambda - 3}} \right)^3, \quad \varepsilon \sim \mathcal{N}(0, 1)$$

- Accept–reject

$$u \leq a(\varepsilon, \lambda)? \quad u \sim \mathcal{U}[0, 1]$$

- Tricky to reparameterize because of the accept–reject step!

A simple method for generating gamma variables, Marsaglia & Tsang, 2000

## Our strategy: partial reparameterization

- *Idea:* Use  $\varepsilon \sim \pi(\varepsilon; \lambda)$  (only *weakly dependent* on  $\lambda$ ) s.t.  
 $\mathbf{x} = h(\varepsilon, \lambda) \sim q(\mathbf{x}; \lambda)$
- *Then:*

$$\begin{aligned}\nabla \mathcal{L}(\lambda) &= \nabla \mathbb{E}_{\pi(\varepsilon; \lambda)} [\log p(h(\varepsilon, \lambda), \mathbf{y}) - \log q(h(\varepsilon, \lambda); \lambda)] \\ &= \underbrace{\mathbb{E}_{\pi(\varepsilon; \lambda)} [\nabla_{\mathbf{x}} (\log p(\mathbf{x}, \mathbf{y}) - \log q(\mathbf{x}; \lambda)) \nabla_{\lambda} h(\varepsilon, \lambda)]}_{\text{reparameterization}} \\ &\quad + \underbrace{\mathbb{E}_{\pi(\varepsilon; \lambda)} [(\log p(h(\varepsilon, \lambda), \mathbf{y}) - \log q(h(\varepsilon, \lambda); \lambda)) \nabla_{\lambda} \log \pi(\varepsilon; \lambda)]}_{\text{correction}}\end{aligned}$$

- *Intuition:* Interpolating between *score function gradients* and *reparameterization gradients*.

# Rejection sampling and reparameterization

- *Idea:*

- Use  $h(\varepsilon, \lambda)$  from the *proposal*
- Let  $\pi(\varepsilon; \lambda)$  be the distribution of the *accepted*  $\varepsilon$

- *Rejection sampling:*

- Propose  $\mathbf{x} = h(\varepsilon, \lambda)$ ,  $\varepsilon \sim s(\varepsilon)$
- Accept if  $u \leq a(\varepsilon, \lambda) = \frac{q(h(\varepsilon, \lambda); \lambda) \left| \frac{dh}{d\varepsilon}(\varepsilon, \lambda) \right|}{M_\lambda s(\varepsilon)}$ ,  $u \sim \mathcal{U}[0, 1]$

- *Then:*

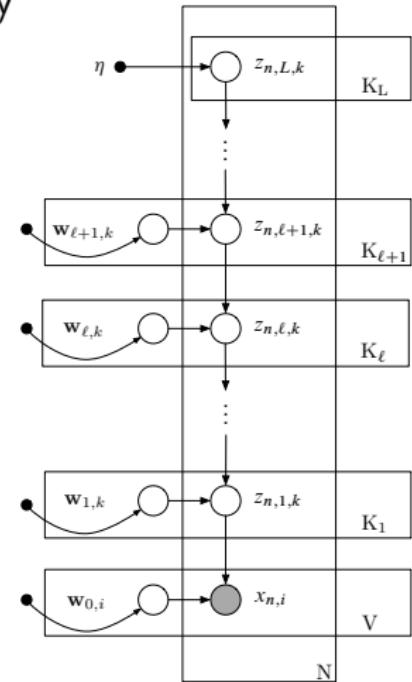
$$\begin{aligned}\pi(\varepsilon; \lambda) &= \int_0^1 \pi(\varepsilon, u; \lambda) du = M_\lambda s(\varepsilon) \int_0^1 \mathbf{1}[u \leq a(\varepsilon, \lambda)] du \\ &= q(h(\varepsilon, \lambda); \lambda) \left| \frac{dh}{d\varepsilon}(\varepsilon, \lambda) \right|\end{aligned}$$

# Sparse gamma deep exponential family

$$\mathbf{w}_{k,k'} \sim \text{Gamma}(0.1, 0.1)$$

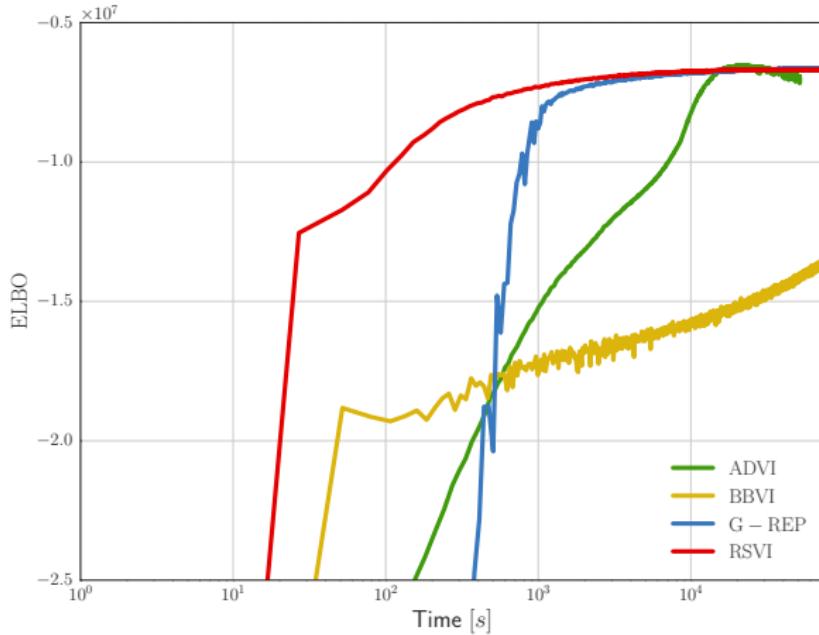
$$\mathbf{x}_{n,k}^\ell \sim \text{Gamma}\left(0.1, \frac{0.1}{\sum_{k'} \mathbf{w}_{k,k'}^\ell \mathbf{x}_{n,k'}^{\ell+1}}\right)$$

$$\mathbf{y}_{n,d} \sim \text{Poisson}\left(\sum_k \mathbf{w}_{k,d}^0 \mathbf{x}_{n,k}^1\right)$$



Deep exponential families, Ranganath et al., 2015

## Sparse gamma DEF – Olivetti faces



[Kucukelbir et al., 2015, Ranganath et al., 2014, Ruiz et al., 2016]

# Variational Monte Carlo methods

- Monte Carlo methods can enable variational inference for more complex (non-conjugate) models  $p(\mathbf{x}, \mathbf{y})$
- But what if we can use Monte Carlo methods to also define more flexible  $q(\mathbf{x} ; \lambda)$ ?
- I will illustrate how we can do this using sequential Monte Carlo for the state space model:

$$p(\mathbf{x}_{1:T}, \mathbf{y}_{1:T}) \triangleq \mu(\mathbf{x}_1) \prod_{t=2}^T f(\mathbf{x}_t | \mathbf{x}_{t-1}) \prod_{t=1}^T g(\mathbf{y}_t | \mathbf{x}_t)$$

- 1 Bayesian inference
- 2 Variational inference
- 3 Stochastic gradient variational inference
- 4 Rejection sampling variational inference
- 5 Variational sequential Monte Carlo

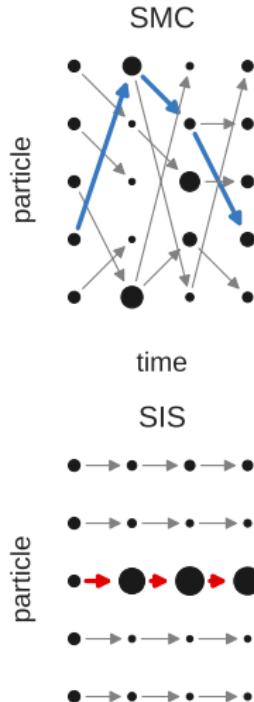
# Sequential Monte Carlo

Approximation:

$$\hat{p}(\mathbf{x}_{1:t-1} \mid \mathbf{y}_{1:t-1}) = \sum_{i=1}^N \frac{w_{t-1}^i}{\sum_\ell w_{t-1}^\ell} \delta_{\mathbf{x}_{1:t-1}^i}$$

Update:

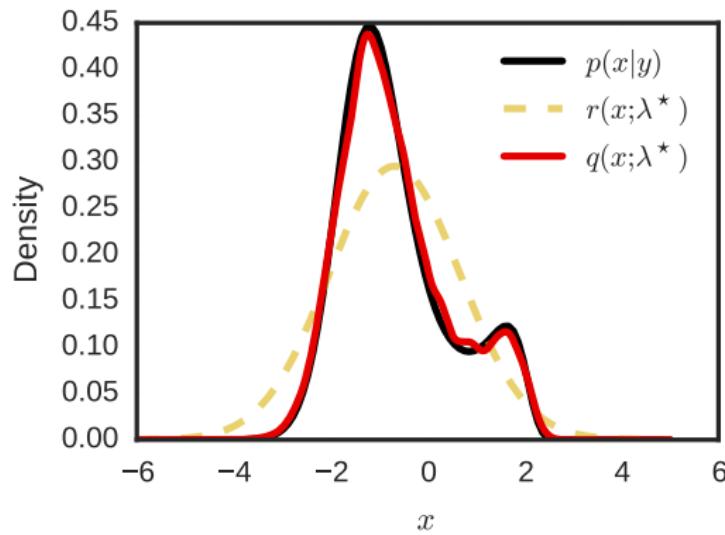
- $a_{t-1}^i \sim \text{Discrete}\left(w_{t-1}^j / \sum_\ell w_{t-1}^\ell\right)$
- $\mathbf{x}_t^i \sim r(\mathbf{x}_t \mid \mathbf{x}_{t-1}^{a_{t-1}^i}; \lambda)$
- $w_t^i = \frac{f(\mathbf{x}_t^i \mid \mathbf{x}_{t-1}^{a_{t-1}^i})g(\mathbf{y}_t \mid \mathbf{x}_t^i)}{r(\mathbf{x}_t^i \mid \mathbf{x}_{t-1}^{a_{t-1}^i}; \lambda)}$



# Variational Sequential Monte Carlo

*Key idea:*

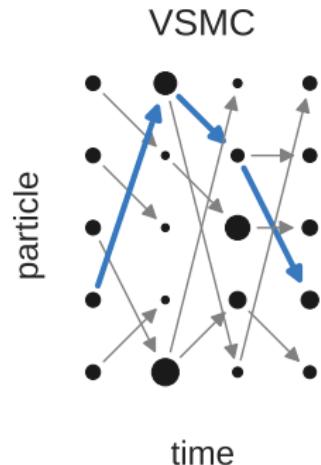
$q(\mathbf{x}_{1:T}; \lambda)$  defined by running SMC and sample  $\hat{p}(\mathbf{x}_{1:T} | \mathbf{y}_{1:T})$



# Variational Sequential Monte Carlo Generative Process

1. Run SMC to get  $\{(\mathbf{x}_{1:T}^i, w_T^i)\}_{i=1}^N$
2.  $b_T \sim \text{Discrete}(w_T^i / \sum_\ell w_T^\ell)$
3. **return**  $\mathbf{x}_{1:T} \triangleq \mathbf{x}_{1:T}^{b_T}$

$$\mathbf{x}_{1:T} \sim \hat{p}(\mathbf{x}_{1:T} \mid \mathbf{y}_{1:T}) = \sum_{i=1}^N \frac{w_T^i}{\sum_\ell w_T^\ell} \delta_{\mathbf{x}_{1:T}^i}$$



# Distribution of VSMC

The distribution  $q$  is the marginal of  $\mathbf{x}_{1:T}^{b_T}$  over all random variables,  $\mathbf{x}_{1:T}^{1:N}, a_{1:T-1}^{1:N}, b_T$ , generated by VSMC

$$q(\mathbf{x}_{1:T} | \mathbf{y}_{1:T}; \lambda) = p(\mathbf{x}_{1:T}, \mathbf{y}_{1:T}) \mathbb{E} \left[ \widehat{p}(\mathbf{y}_{1:T})^{-1} \middle| \mathbf{x}_{1:T} \right],$$

where  $\widehat{p}(\mathbf{y}_{1:T}) \triangleq \prod_{t=1}^T \frac{1}{N} \sum_{i=1}^N w_t^i$ .

# Surrogate ELBO

We can not directly use the standard ELBO,  $\mathcal{L}(\lambda)$ , because evaluating  $q$  pointwise is intractable.

Consider instead:

$$\tilde{\mathcal{L}}(\lambda) \triangleq \mathbb{E} \left[ \sum_{t=1}^T \log \left( \frac{1}{N} \sum_{i=1}^N w_t^i \right) \right] = \mathbb{E} [\log \hat{p}(\mathbf{y}_{1:T})]$$

**Theorem 1** (Surrogate ELBO).

$$\log p(\mathbf{y}_{1:T}) \geq \mathcal{L}(\lambda) \geq \tilde{\mathcal{L}}(\lambda).$$

*Proof.* (*Idea*) Use the definition of  $\mathcal{L}(\lambda)$  for  $q(\mathbf{x}_{1:T} | \mathbf{y}_{1:T}; \lambda)$  and then Jensen's inequality.  $\square$

# Stochastic Optimization

We optimize the surrogate ELBO  $\tilde{\mathcal{L}}(\lambda)$  using stochastic gradient ascent,

$$\nabla \tilde{\mathcal{L}}(\lambda) = \nabla \mathbb{E} [\log \hat{p}(\mathbf{y}_{1:T})] = \mathbb{E} \left[ \nabla \log \hat{p}(\mathbf{y}_{1:T}) + \log \hat{p}(\mathbf{y}_{1:T}) \nabla \log \tilde{\phi} \right].$$

Variance reduction using the reparameterization trick, Rao-Blackwellization, and control variates.

# Theoretical Results

- For some  $c(\lambda) < \infty$ ,

$$\text{KL} \left( q(\mathbf{x}_{1:T}; \lambda) \middle\| p(\mathbf{x}_{1:T} | \mathbf{y}_{1:T}) \right) \leq \frac{c(\lambda)}{N}.$$

- With  $N = bT$ ,

$$\begin{aligned} \text{KL} \left( q(\mathbf{x}_{1:T}; \lambda) \middle\| p(\mathbf{x}_{1:T} | \mathbf{y}_{1:T}) \right) &\leq -\mathbb{E} \left[ \log \frac{\widehat{p}(\mathbf{y}_{1:T})}{p(\mathbf{y}_{1:T})} \right] \\ &\xrightarrow[T \rightarrow \infty]{} \frac{\sigma^2(\lambda)}{2b}, \end{aligned}$$

for  $\sigma^2(\lambda) < \infty$ .

# Experiments - Stochastic volatility

*Data:* 10 years of monthly returns for the exchange rate of 22 currencies wrt to USD

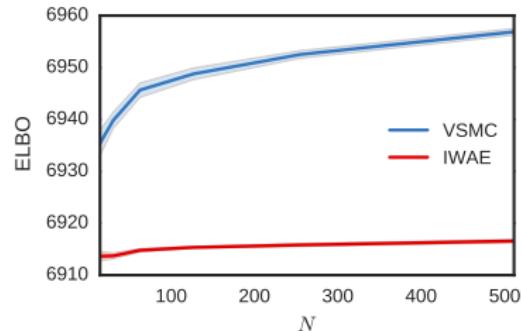
$$\mathbf{x}_t = \mu + \phi(\mathbf{x}_{t-1} - \mu) + v_t,$$

$$v_t \sim \mathcal{N}(0, Q)$$

$$\mathbf{y}_t = \beta \exp\left(\frac{\mathbf{x}_t}{2}\right) e_t,$$

$$e_t \sim \mathcal{N}(0, I)$$

	Method	ELBO
$N = 4$	Structured VI	6905.1
	IWAE	6911.2
$N = 8$	<b>VSMC</b>	<b>6921.6</b>
	IWAE	6912.4
$N = 16$	<b>VSMC</b>	<b>6935.8</b>
	IWAE	6913.3
	<b>VSMC</b>	<b>6936.6</b>



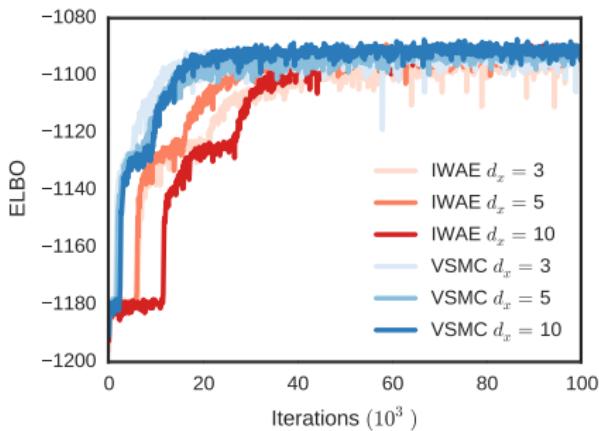
$$r(\mathbf{x}_t | \mathbf{x}_{t-1}; \lambda, \theta) \propto f(\mathbf{x}_t | \mathbf{x}_{t-1}; \theta) \mathcal{N}(\mathbf{x}_t; \mu_t, \Sigma_t)$$

# Experiments - Stochastic recurrent neural networks

*Data:* 105 motor cortex neurons simultaneously recorded in a macaque monkey

$$\mathbf{x}_t = \mu_\theta(\mathbf{x}_{t-1}) + \exp(\sigma_\theta(\mathbf{x}_{t-1})/2) v_t,$$

$$\mathbf{y}_t \sim \text{Poisson}(\exp(\eta_\theta(\mathbf{x}_t))),$$



Proposal:

$$r(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y}_t) \propto f(\mathbf{x}_t | \mathbf{x}_{t-1}; \theta)$$

$$\mathcal{N}\left(\mathbf{x}_t; \mu_\lambda(\mathbf{y}_t), e^{2\sigma_\lambda(\mathbf{y}_t)}\right),$$

$\mu_\lambda(\cdot), \sigma_\lambda(\cdot)$  are *inference networks.*

# Collaborators

Scott  
Linderman



Francisco  
Ruiz



Rajesh  
Ranganath



David Blei



Work done while visiting Columbia University



Reparameterization Gradients through Acceptance-Rejection

Sampling Algorithms, Naesseth et al., 2017

Variational Sequential Monte Carlo, Naesseth et al., 2017

# Take-home messages

- Stochastic optimization opens up for interesting new approximate Bayesian inference methods
- Monte Carlo methods  $\Rightarrow$  more flexible variational approximations with guarantees
- Variational methods  $\Rightarrow$  adaptation of Monte Carlo methods based on a global objective

Contact: [christian.a.naesseth@liu.se](mailto:christian.a.naesseth@liu.se)

Website: [users.isy.liu.se/en/rt/chran60/index.html](http://users.isy.liu.se/en/rt/chran60/index.html)



# Subsampling and Variational EM

- Data subsampling  $\mathbf{y}_{1:T} = \{\mathbf{y}_{1:T}^j\}_{j=1}^n$

$$\tilde{\mathcal{L}}(\lambda) = \sum_{j=1}^n \mathbb{E} \left[ \log \hat{p}(\mathbf{y}_{1:T}^j) \right]$$

At each iteration sample a mini-batch of data.

- Variational EM

$$\log p_\theta(\mathbf{y}_{1:T}) \geq \tilde{\mathcal{L}}(\lambda, \theta) = \mathbb{E} [\log \hat{p}_\theta(\mathbf{y}_{1:T})]$$

Maximize  $\tilde{\mathcal{L}}(\lambda, \theta)$  jointly with respect to  $\lambda, \theta$