

Variational Hamiltonian Monte Carlo via Score Matching

Cheng Zhang (Joint work with Prof. Shahbaba and Prof. Zhao)

Department of Mathematics
University of California, Irvine

Jan 6, 2017

Outline

- 1 Background
 - Bayesian Inference
- 2 Markov chain Monte Carlo
 - Metropolis-Hastings Algorithm
 - Hamiltonian Monte Carlo
 - Scalable MCMC
- 3 Fixed-Form Variational Bayes
 - Lower Bounds and Free Energy
 - Variational Bayes as Linear Regression
- 4 Variational Hamiltonian Monte Carlo
 - Approximation with Random Bases
 - Variational HMC
 - Experiments
- 5 Conclusion

Bayesian Inference

- Bayesian inference model
 - $\mathcal{D} = \{y^1, \dots, y^N\}$: observed data
 - $\theta \in \mathbb{R}^d$: model parameter
 - $p(\mathcal{D}|\theta)$: model density
 - $p(\theta)$: prior

Bayesian Inference

- Bayesian inference model
 - $\mathcal{D} = \{y^1, \dots, y^N\}$: observed data
 - $\theta \in \mathbb{R}^d$: model parameter
 - $p(\mathcal{D}|\theta)$: model density
 - $p(\theta)$: prior
- Goal : learning parameter θ from data

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta) \cdot p(\theta)}{p(\mathcal{D})} \propto p(\mathcal{D}|\theta) \cdot p(\theta)$$

Bayesian Inference

- Bayesian inference model
 - $\mathcal{D} = \{y^1, \dots, y^N\}$: observed data
 - $\theta \in \mathbb{R}^d$: model parameter
 - $p(\mathcal{D}|\theta)$: model density
 - $p(\theta)$: prior
- Goal : learning parameter θ from data

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta) \cdot p(\theta)}{p(\mathcal{D})} \propto p(\mathcal{D}|\theta) \cdot p(\theta)$$

- Difficulty : $p(\mathcal{D})$ unknown \Rightarrow **intractable posterior distribution** $p(\theta|\mathcal{D})$
e.g., probabilistic graphical models, Bayesian hierarchical models

Bayesian Inference

■ Bayesian inference model

- $\mathcal{D} = \{y^1, \dots, y^N\}$: observed data
- $\theta \in \mathbb{R}^d$: model parameter
- $p(\mathcal{D}|\theta)$: model density
- $p(\theta)$: prior

■ Goal : learning parameter θ from data

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta) \cdot p(\theta)}{p(\mathcal{D})} \propto p(\mathcal{D}|\theta) \cdot p(\theta)$$

■ Difficulty : $p(\mathcal{D})$ unknown \Rightarrow **intractable posterior distribution** $p(\theta|\mathcal{D})$ e.g., probabilistic graphical models, Bayesian hierarchical models

■ Two popular approximations

- **Markov chain Monte Carlo**. Sample by running a Markov chain : asymptotically unbiased but computationally slow
- **Variational Bayes**. Approximate via tractable distributions : computationally fast but may result in poor approximation.

Markov chain Monte Carlo

- Intuitive idea : evolve a Markov chain to sample from a target distribution $\pi(\theta)$ (METROPOLIS et al. 1953).



Markov chain Monte Carlo

- Intuitive idea : evolve a Markov chain to sample from a target distribution $\pi(\theta)$ (METROPOLIS et al. 1953).
- Conditions for transition kernel $T(\cdot|\cdot)$
 - *Irreducibility* : any state has positive probability of visiting any other state.
 - *Aperiodicity* : The chain should not get trapped in cycles.
 - *Detailed Balance condition (sufficient)* :

$$\pi(\theta)T(\theta'|\theta) = \pi(\theta')T(\theta|\theta')$$



Markov chain Monte Carlo

- Intuitive idea : evolve a Markov chain to sample from a target distribution $\pi(\theta)$ (METROPOLIS et al. 1953).
- Conditions for transition kernel $T(\cdot|\cdot)$
 - *Irreducibility* : any state has positive probability of visiting any other state.
 - *Aperiodicity* : The chain should not get trapped in cycles.
 - *Detailed Balance condition (sufficient)* :

$$\pi(\theta)T(\theta'|\theta) = \pi(\theta')T(\theta|\theta')$$

- Metropolis-Hastings algorithms (one iteration)

- 1 sample $\theta' \sim q(\theta'|\theta)$
- 2 update the current state to θ' with probability $\alpha(\theta, \theta') = \min[1, \frac{\pi(\theta')q(\theta|\theta')}{\pi(\theta)q(\theta'|\theta)}]$

Markov chain Monte Carlo

- Intuitive idea : evolve a Markov chain to sample from a target distribution $\pi(\theta)$ (METROPOLIS et al. 1953).
- Conditions for transition kernel $T(\cdot|\cdot)$
 - *Irreducibility* : any state has positive probability of visiting any other state.
 - *Aperiodicity* : The chain should not get trapped in cycles.
 - *Detailed Balance condition (sufficient)* :

$$\pi(\theta)T(\theta'|\theta) = \pi(\theta')T(\theta|\theta')$$

- Metropolis-Hastings algorithms (one iteration)
 - 1 sample $\theta' \sim q(\theta'|\theta)$
 - 2 update the current state to θ' with probability $\alpha(\theta, \theta') = \min[1, \frac{\pi(\theta')q(\theta|\theta')}{\pi(\theta)q(\theta'|\theta)}]$
- Pros & Cons for simple MCMCs (e.g., RWM and Gibbs sampling)
 - Pro : easy to implement and computationally cheap
 - Con : slow mixing due to random walk behaviors, especially in complicate, high-dimensional models.

Hamiltonian Monte Carlo

- Intuition : Leveraging a Hamiltonian dynamical system to generate trial moves in MCMC samplers. (DUANE et al. 1987, NEAL 2011)



Hamiltonian Monte Carlo

- Intuition : Leveraging a Hamiltonian dynamical system to generate trial moves in MCMC samplers. (DUANE et al. 1987, NEAL 2011)
- Model based energy function : the Hamiltonian

$$H(\boldsymbol{\theta}, \mathbf{r}) = U(\boldsymbol{\theta}) + K(\mathbf{r})$$

- Potential : $U(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta}, \mathcal{D}) = -[\log p(\boldsymbol{\theta}) + \log p(\mathcal{D}|\boldsymbol{\theta})]$
- Kinetic : $K(\mathbf{r}) = \frac{1}{2}\mathbf{r}^\top M^{-1}\mathbf{r} \Rightarrow \pi(\mathbf{r}) \sim \mathcal{N}(\mathbf{0}, M)$

The joint density of $\mathbf{z} = (\boldsymbol{\theta}, \mathbf{r})$ is

$$\pi(\mathbf{z}) \propto \exp(-U(\boldsymbol{\theta}) - K(\mathbf{r})) \propto p(\boldsymbol{\theta}|\mathcal{D}) \cdot \pi(\mathbf{r})$$

Hamiltonian Monte Carlo

- Intuition : Leveraging a Hamiltonian dynamical system to generate trial moves in MCMC samplers. (DUANE et al. 1987, NEAL 2011)
- Model based energy function : the Hamiltonian

$$H(\boldsymbol{\theta}, \mathbf{r}) = U(\boldsymbol{\theta}) + K(\mathbf{r})$$

- Potential : $U(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta}, \mathcal{D}) = -[\log p(\boldsymbol{\theta}) + \log p(\mathcal{D}|\boldsymbol{\theta})]$
- Kinetic : $K(\mathbf{r}) = \frac{1}{2}\mathbf{r}^\top M^{-1}\mathbf{r} \Rightarrow \pi(\mathbf{r}) \sim \mathcal{N}(\mathbf{0}, M)$

The joint density of $\mathbf{z} = (\boldsymbol{\theta}, \mathbf{r})$ is

$$\pi(\mathbf{z}) \propto \exp(-U(\boldsymbol{\theta}) - K(\mathbf{r})) \propto p(\boldsymbol{\theta}|\mathcal{D}) \cdot \pi(\mathbf{r})$$

- Hamilton's equations :

$$\frac{d\boldsymbol{\theta}}{dt} = \nabla_{\mathbf{r}} H = \nabla_{\mathbf{r}} K(\mathbf{r}), \quad \frac{d\mathbf{r}}{dt} = -\nabla_{\boldsymbol{\theta}} H = -\nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta})$$

Hamiltonian flow $\phi_S^H : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$, $\mathbf{z}(0) = \mathbf{z} \mapsto \mathbf{z}^* = \mathbf{z}(s)$

Hamiltonian Monte Carlo

- Intuition : Leveraging a Hamiltonian dynamical system to generate trial moves in MCMC samplers. (DUANE et al. 1987, NEAL 2011)
- Model based energy function : the Hamiltonian

$$H(\boldsymbol{\theta}, \mathbf{r}) = U(\boldsymbol{\theta}) + K(\mathbf{r})$$

- Potential : $U(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta}, \mathcal{D}) = -[\log p(\boldsymbol{\theta}) + \log p(\mathcal{D}|\boldsymbol{\theta})]$
- Kinetic : $K(\mathbf{r}) = \frac{1}{2} \mathbf{r}^\top M^{-1} \mathbf{r} \Rightarrow \pi(\mathbf{r}) \sim \mathcal{N}(\mathbf{0}, M)$

The joint density of $\mathbf{z} = (\boldsymbol{\theta}, \mathbf{r})$ is

$$\pi(\mathbf{z}) \propto \exp(-U(\boldsymbol{\theta}) - K(\mathbf{r})) \propto p(\boldsymbol{\theta}|\mathcal{D}) \cdot \pi(\mathbf{r})$$

- Hamilton's equations :

$$\frac{d\boldsymbol{\theta}}{dt} = \nabla_{\mathbf{r}} H = \nabla_{\mathbf{r}} K(\mathbf{r}), \quad \frac{d\mathbf{r}}{dt} = -\nabla_{\boldsymbol{\theta}} H = -\nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta})$$

Hamiltonian flow $\phi_s^H : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$, $\mathbf{z}(0) = \mathbf{z} \mapsto \mathbf{z}^* = \mathbf{z}(s)$

- Properties : **reversibility**, **volume preservation** and **constant Hamiltonian** over time t



Hamiltonian Monte Carlo

- Numerical integrators : the leap-frog scheme

$$\begin{aligned} \mathbf{r}^{(t+1/2)} &= \mathbf{r}^{(t)} - \frac{\varepsilon}{2} \nabla_{\theta} U(\theta^{(t)}) \\ \theta^{(t+1)} &= \theta^{(t)} + \varepsilon M^{-1} \mathbf{r}^{(t+1/2)} \\ \mathbf{r}^{(t+1)} &= \mathbf{r}^{(t+1/2)} - \frac{\varepsilon}{2} \nabla_{\theta} U(\theta^{(t+1)}) \end{aligned} \tag{1}$$

The leap-frog scheme (1) remains **time reversible** and **volume preserving**. 😊
Constant Hamiltonian does not hold. 😞



Hamiltonian Monte Carlo

- Numerical integrators : the leap-frog scheme

$$\begin{aligned}
 \mathbf{r}^{(t+1/2)} &= \mathbf{r}^{(t)} - \frac{\varepsilon}{2} \nabla_{\theta} U(\theta^{(t)}) \\
 \theta^{(t+1)} &= \theta^{(t)} + \varepsilon M^{-1} \mathbf{r}^{(t+1/2)} \\
 \mathbf{r}^{(t+1)} &= \mathbf{r}^{(t+1/2)} - \frac{\varepsilon}{2} \nabla_{\theta} U(\theta^{(t+1)})
 \end{aligned} \tag{1}$$

The leap-frog scheme (1) remains **time reversible** and **volume preserving**. 😊
Constant Hamiltonian does not hold. 😞

- Metropolis correction

$$\alpha_{hmc}(\mathbf{z}, \mathbf{z}^*) = \min\{1, \exp(-H(\mathbf{z}^*) + H(\mathbf{z}))\} \tag{2}$$

Hamiltonian Monte Carlo

- Numerical integrators : the leap-frog scheme

$$\begin{aligned}
 \mathbf{r}^{(t+1/2)} &= \mathbf{r}^{(t)} - \frac{\varepsilon}{2} \nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta}^{(t)}) \\
 \boldsymbol{\theta}^{(t+1)} &= \boldsymbol{\theta}^{(t)} + \varepsilon M^{-1} \mathbf{r}^{(t+1/2)} \\
 \mathbf{r}^{(t+1)} &= \mathbf{r}^{(t+1/2)} - \frac{\varepsilon}{2} \nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta}^{(t+1)})
 \end{aligned} \tag{1}$$

The leap-frog scheme (1) remains **time reversible** and **volume preserving**. 😊
Constant Hamiltonian does not hold. 😞

- Metropolis correction

$$\alpha_{hmc}(\mathbf{z}, \mathbf{z}^*) = \min\{1, \exp(-H(\mathbf{z}^*) + H(\mathbf{z}))\} \tag{2}$$

- Some variants.

- Automatically tuning of the hyper-parameters (e.g., step size ε , number of leap-frog steps L) (HOFFMAN et GELMAN 2011 ; WANG, MOHAMED et NANDO 2013)
- Riemannian Manifold HMC (GIROLAMI et CALDERHEAD 2011)

Stochastic Gradient MCMC

- Challenge from massive data.

$$\nabla_{\theta} U(\theta) = - \sum_{y \in \mathcal{D}} \nabla_{\theta} \log p(y|\theta) - \nabla_{\theta} \log p(\theta) \sim \mathcal{O}(N)$$

Stochastic Gradient MCMC

- Challenge from massive data.

$$\nabla_{\theta} U(\theta) = - \sum_{y \in \mathcal{D}} \nabla_{\theta} \log p(y|\theta) - \nabla_{\theta} \log p(\theta) \sim \mathcal{O}(N)$$

- Stochastic gradient MCMC : use stochastic gradients

$$\nabla_{\theta} \tilde{U}(\theta) = - \frac{|\mathcal{D}|}{|\tilde{\mathcal{D}}|} \sum_{y \in \tilde{\mathcal{D}}} \nabla_{\theta} \log p(y|\theta) - \nabla_{\theta} \log p(\theta), \quad \tilde{\mathcal{D}} \subset \mathcal{D}$$

e.g., SGLD (WELLING et TEH 2011), SGHMC (CHEN, E. B. FOX et GUESTRIN 2014), SGNHT (DING et al. 2014)

Stochastic Gradient MCMC

- Challenge from massive data.

$$\nabla_{\theta} U(\theta) = - \sum_{y \in \mathcal{D}} \nabla_{\theta} \log p(y|\theta) - \nabla_{\theta} \log p(\theta) \sim \mathcal{O}(N)$$

- Stochastic gradient MCMC : use stochastic gradients

$$\nabla_{\theta} \tilde{U}(\theta) = - \frac{|\mathcal{D}|}{|\tilde{\mathcal{D}}|} \sum_{y \in \tilde{\mathcal{D}}} \nabla_{\theta} \log p(y|\theta) - \nabla_{\theta} \log p(\theta), \quad \tilde{\mathcal{D}} \subset \mathcal{D}$$

e.g., SGLD (WELLING et TEH 2011), SGHMC (CHEN, E. B. FOX et GUESTRIN 2014), SGNHT (DING et al. 2014)

- Properties of stochastic gradient MCMC

- Convergence relies on stochastic differential equation of the form

$$d\mathbf{z} = f(\mathbf{z})dt + \sqrt{2D(\mathbf{z})}dW(t)$$

with appropriate $f(\mathbf{z})$ and $D(\mathbf{z})$. See a complete recipe (MA, CHEN et E. FOX 2015)

- Need to anneal (or use small) step size, sacrificing the exploration efficiency to compromise scalability (BETANCOURT 2015).

Surrogate Method

Find cheap surrogate function $U^*(\theta)$ and transition kernel $T_s(\cdot|\cdot)$ that leaves $q^*(\theta) \propto \exp(-U^*(\theta))$ invariant.

$$q^*(\theta)T_s(\theta'|\theta) = q^*(\theta')T_s(\theta|\theta')$$

Use T_s to generate proposals

Proposition (Liu 2001)

The target distribution $p(\theta|\mathcal{D}) \propto \exp(-U(\theta))$ is the stationary distribution for a Markov chain simulated according to the following procedure : given the current state θ , let $\vartheta_0 = \theta$ and recursively sample $\vartheta_i \sim T_s(\cdot|\vartheta_{i-1})$ for $i = 1, \dots, k$. Then, accept the proposal $\theta^ = \vartheta_k$ with the following probability*

$$\alpha_s(\theta, \theta^*) = \min \left\{ 1, \frac{p(\theta^*|\mathcal{D})q^*(\theta)}{p(\theta|\mathcal{D})q^*(\theta^*)} \right\}$$

Some of the existing surrogate methods : Gaussian process surrogate (RASMUSSEN 2003; LAN et al. 2015), reproducing kernel Hilbert space surrogate (STRATHMANN et al. 2015) and random network surrogate (ZHANG, SHAHBABA et ZHAO 2015).

Lower Bounds and Free Energy

Fixed-Form Variational Bayes (HONKELA et al. 2010 ; SAUL et JORDAN 1996) uses a parametrized distribution

$$q_{\eta}(\theta) = \exp[T(\theta)\eta - A(\eta)] \quad (3)$$

to approximate the target posterior $p(\theta|\mathcal{D})$

Lower Bounds and Free Energy

Fixed-Form Variational Bayes (HONKELA et al. 2010 ; SAUL et JORDAN 1996) uses a parametrized distribution

$$q_{\eta}(\theta) = \exp[T(\theta)\eta - A(\eta)] \quad (3)$$

to approximate the target posterior $p(\theta|\mathcal{D})$

- Distance between distributions

$$\begin{aligned} D_{KL}(q_{\eta}(\theta) \| p(\theta|\mathcal{D})) &= \int q_{\eta}(\theta) \log \left(\frac{q_{\eta}(\theta)}{p(\theta|\mathcal{D})} \right) d\theta \\ &= \log p(\mathcal{D}) - \underbrace{\int q_{\eta}(\theta) \log \left(\frac{p(\theta, \mathcal{D})}{q_{\eta}(\theta)} \right) d\theta}_{\text{Free Energy (Lower bound)}} \end{aligned}$$

Lower Bounds and Free Energy

Fixed-Form Variational Bayes (HONKELA et al. 2010 ; SAUL et JORDAN 1996) uses a parametrized distribution

$$q_{\eta}(\theta) = \exp[T(\theta)\eta - A(\eta)] \quad (3)$$

to approximate the target posterior $p(\theta|\mathcal{D})$

- Distance between distributions

$$\begin{aligned} D_{KL}(q_{\eta}(\theta) \| p(\theta|\mathcal{D})) &= \int q_{\eta}(\theta) \log \left(\frac{q_{\eta}(\theta)}{p(\theta|\mathcal{D})} \right) d\theta \\ &= \log p(\mathcal{D}) - \underbrace{\int q_{\eta}(\theta) \log \left(\frac{p(\theta, \mathcal{D})}{q_{\eta}(\theta)} \right) d\theta}_{\text{Free Energy (Lower bound)}} \end{aligned}$$

- An optimization problem

$$\hat{\eta} = \arg \max_{\eta} \mathbb{E}_{q_{\eta}} [\log p(\theta, \mathcal{D}) - \log q_{\eta}(\theta)] \quad (4)$$

more accurate than using mean-field assumptions ; requires analytically evaluation of $\mathbb{E}_{q_{\eta}} \log q_{\eta}(\theta)$, $\mathbb{E}_{q_{\eta}} \log p(\theta, \mathcal{D})$ and the derivatives.

Stochastic Linear Regression

Optimization in (4) can be solved using stochastic linear regression (SALIMANS et KNOWLES 2013). Rewrite (3) in the unnormalized form

$$\tilde{q}_{\tilde{\eta}}(\boldsymbol{\theta}) = \exp[\tilde{T}(\boldsymbol{\theta})\tilde{\eta}], \quad \tilde{T}(\boldsymbol{\theta}) = (1, T(\boldsymbol{\theta})), \quad \tilde{\eta} = (\eta_0, \boldsymbol{\eta}^\top)^\top$$

Unnormalized KL divergence :

$$\begin{aligned} \tilde{D}_{KL}(\tilde{q}_{\tilde{\eta}} \| p(\boldsymbol{\theta}, \mathcal{D})) &= \int \tilde{q}_{\tilde{\eta}}(\boldsymbol{\theta}) \log \frac{\tilde{q}_{\tilde{\eta}}(\boldsymbol{\theta})}{p(\boldsymbol{\theta}, \mathcal{D})} d\boldsymbol{\theta} - \int \tilde{q}_{\tilde{\eta}}(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \int \exp[\tilde{T}(\boldsymbol{\theta})\tilde{\eta}] [\tilde{T}(\boldsymbol{\theta})\tilde{\eta} - \log p(\boldsymbol{\theta}, \mathcal{D})] d\boldsymbol{\theta} - \int \exp[\tilde{T}(\boldsymbol{\theta})\tilde{\eta}] d\boldsymbol{\theta} \end{aligned} \quad (5)$$

Find the minimum by differentiation

$$\nabla_{\tilde{\eta}} \tilde{D}_{KL}(\tilde{q}_{\tilde{\eta}} \| p(\boldsymbol{\theta}, \mathcal{D})) = \int \tilde{q}_{\tilde{\eta}}(\boldsymbol{\theta}) [\tilde{T}(\boldsymbol{\theta})^\top \tilde{T}(\boldsymbol{\theta})\tilde{\eta} - \tilde{T}(\boldsymbol{\theta})^\top \log p(\boldsymbol{\theta}, \mathcal{D})] d\boldsymbol{\theta} = \mathbf{0}$$

The minimum :

$$\tilde{\eta} = \left(\mathbb{E}_q[\tilde{T}(\boldsymbol{\theta})^\top \tilde{T}(\boldsymbol{\theta})] \right)^{-1} \mathbb{E}_q[\tilde{T}(\boldsymbol{\theta})^\top \log p(\boldsymbol{\theta}, \mathcal{D})] \quad (6)$$

Monte Carlo Estimation

(6) is not a solution yet, but can be used to derive a fix point iteration. Let

$$C = \mathbb{E}_q[\tilde{T}(\theta)^\top \tilde{T}(\theta)], \quad g = \mathbb{E}_q[\tilde{T}(\theta)^\top \log p(\theta, \mathcal{D})]$$

then $\tilde{\eta} = C^{-1}g$.

Stochastic Optimization for Fixed-Form VB

- 1: Initialize $\tilde{\eta}_1, C_1, g_1 = C_1 \tilde{\eta}_1$, set step-size $w \in [0, 1]$
- 2: **for** $t = 1$ to T **do**
- 3: Drawing a single sample θ_t^* from the current approximation $q_{\eta_t}(\theta)$
- 4: Update C_{t+1}, g_{t+1} as follows :
- 5: $g_{t+1} = (1 - w)g_t + w\hat{g}_t, \quad \hat{g}_t = \tilde{T}(\theta_t^*)^\top \log p(\theta_t^*, \mathcal{D})$
- 6: $C_{t+1} = (1 - w)C_t + w\hat{C}_t, \quad \hat{C}_t = \tilde{T}(\theta_t^*)^\top \tilde{T}(\theta_t^*)$
- 7: Update the parameters :
- 8: $\tilde{\eta}_{t+1} = C_{t+1}^{-1}g_{t+1}$
- 9: **end for**

See SALIMANS et KNOWLES 2013 for more variants.

Combine VB and HMC via Surrogate

- Surrogate method revisit

Target : $p(\theta|\mathcal{D}) \propto \exp(-U(\theta))$, Surrogate : $q^*(\theta) \propto \exp(-U^*(\theta))$

Acceptance probability

$$\alpha_s(\theta, \theta^*) = \min \left\{ 1, \frac{p(\theta^*|\mathcal{D})q^*(\theta)}{p(\theta|\mathcal{D})q^*(\theta^*)} \right\}$$

High approximation quality \Rightarrow high acceptance rate.

Combine VB and HMC via Surrogate

- Surrogate method revisit

Target : $p(\theta|\mathcal{D}) \propto \exp(-U(\theta))$, Surrogate : $q^*(\theta) \propto \exp(-U^*(\theta))$

Acceptance probability

$$\alpha_s(\theta, \theta^*) = \min \left\{ 1, \frac{p(\theta^*|\mathcal{D})q^*(\theta)}{p(\theta|\mathcal{D})q^*(\theta^*)} \right\}$$

High approximation quality \Rightarrow high acceptance rate.

- Flexible and efficient surrogate based on random network

$$U^*(\theta) = z(\theta) = \sum_{i=1}^s v_i a(\theta; \gamma_i) \sim \mathcal{O}(s) \quad (7)$$

where $\{\gamma_i\}_{i=1}^n$ are random samples from some distribution.

Combine VB and HMC via Surrogate

- Surrogate method revisit

Target : $p(\theta|\mathcal{D}) \propto \exp(-U(\theta))$, Surrogate : $q^*(\theta) \propto \exp(-U^*(\theta))$

Acceptance probability

$$\alpha_s(\theta, \theta^*) = \min \left\{ 1, \frac{p(\theta^*|\mathcal{D})q^*(\theta)}{p(\theta|\mathcal{D})q^*(\theta^*)} \right\}$$

High approximation quality \Rightarrow high acceptance rate.

- Flexible and efficient surrogate based on random network

$$U^*(\theta) = z(\theta) = \sum_{i=1}^s v_i a(\theta; \gamma_i) \quad \sim \quad \mathcal{O}(s) \quad (7)$$

where $\{\gamma_i\}_{i=1}^n$ are random samples from some distribution.

- Use **variational Bayes** to improve approximation.

$$\hat{q}(\theta) = \arg \min_{q^*} D(q^*(\theta), p(\theta, \mathcal{D}))$$

A Rich Class of Functions : \mathcal{F}_p

Bases : $\{a(\theta; \gamma) : \gamma \in \Gamma\}$, $\Theta \subset \mathbb{R}^d$, a distribution on $\Gamma : p(\gamma)$, consider functions

$$f(\theta) = \int_{\Gamma} \alpha(\gamma) a(\theta; \gamma) d\gamma \quad (8)$$

where $|\alpha(\gamma)| \leq C|p(\gamma)|$. Define a norm $\|f\|_p = \sup_{\gamma} \left| \frac{\alpha(\gamma)}{p(\gamma)} \right|$ and the set

$$\mathcal{F}_p \equiv \left\{ f(\theta) = \int_{\Gamma} \alpha(\gamma) a(\theta; \gamma) d\gamma \mid \|f\|_p < \infty \right\}$$

Theorem (Rahimi 2008)

Let μ be any probability measure on Θ , $\|f\|_{\mu}^2 = \int_{\Theta} f^2(\theta) \mu(d\theta)$. Suppose $\sup_{\theta, \gamma} |a(\theta; \gamma)| \leq 1$. Fix $f \in \mathcal{F}_p$. $\forall \delta > 0$, with probability at least $1 - \delta$ over $\gamma_i \stackrel{iid}{\sim} p(\gamma)$, there exist v_1, \dots, v_s such that $z(\theta) = \sum_{i=1}^s v_i a(\theta; \gamma_i)$ satisfies

$$\|z - f\|_{\mu} < \frac{\|f\|_p}{\sqrt{s}} \left(1 + \sqrt{2 \log \frac{1}{\delta}} \right)$$

\mathcal{F}_p is dense in \mathcal{H}

A reproducing kernel Hilbert space (RKHS) \mathcal{H} with kernel k on $\Theta \times \Theta$

$$k(\theta, \theta') = \int_{\Gamma} p(\gamma) a(\theta; \gamma) a(\theta'; \gamma) d\gamma$$

\mathcal{H} can be constructed based on functions of the form (8).

Proposition (Rahimi 2008)

Let $\hat{\mathcal{H}}$ be the completion of the set of functions of the form (8) such that

$$\int_{\Gamma} \frac{\alpha(\gamma)^2}{p(\gamma)} d\gamma < \infty \quad (9)$$

with the inner product

$$\langle f, g \rangle = \int_{\Gamma} \frac{\alpha(\gamma)\beta(\gamma)}{p(\gamma)} d\gamma$$

where $g(\theta) = \int_{\Gamma} \beta(\gamma) a(\theta; \gamma) d\gamma$. Then $\hat{\mathcal{H}} = \mathcal{H}$

Note that $\forall f \in \mathcal{F}_p$, (9) is satisfied. Therefore, \mathcal{F}_p is a subset of \mathcal{H} . In fact, \mathcal{F}_p is dense in \mathcal{H} . See RAHIMI et RECHT 2008 for detailed proofs.

Free-Form Variational Bayes

- Random bases surrogate induced distribution

$$q_{\mathbf{v}}(\boldsymbol{\theta}) \propto \exp(-z(\boldsymbol{\theta})) = \exp\left[-\sum_{i=1}^S v_i a(\boldsymbol{\theta}; \gamma_i) - \Phi(\mathbf{v})\right]$$

where γ_i are drawn iid from some distribution. The natural parameters \mathbf{v} of interest belong to the set

$$\Omega := \{\mathbf{v} \in \mathbb{R}^S \mid \Phi(\mathbf{v}) < +\infty\}$$

- A distance measure $D : (q, p) \rightarrow [0, +\infty)$, q, p are unnormalized densities.
- **Free-Form** variational inference

$$\hat{\mathbf{v}} = \arg \min_{\mathbf{v} \in \Omega} D(q_{\mathbf{v}}(\boldsymbol{\theta}) \| p(\boldsymbol{\theta}, \mathcal{D})) \quad (10)$$

- $q_{\mathbf{v}}(\boldsymbol{\theta})$ does not have to be tractable
- **Free style** construction of the random network surrogate.

Distance Between Unnormalized Densities

- Potential matching distance

$$\begin{aligned} D_{PM}(q_{\mathbf{v}}(\boldsymbol{\theta})\|p(\boldsymbol{\theta}, \mathcal{D})) &= \frac{1}{2} \min_{\mathbf{b}} \int q_{\mathbf{v}}(\boldsymbol{\theta}) \|z(\boldsymbol{\theta}) - U(\boldsymbol{\theta}) - \mathbf{b}\|^2 d\boldsymbol{\theta} \\ &= \frac{1}{2} \text{Var}_{q_{\mathbf{v}}}(z(\boldsymbol{\theta}) - U(\boldsymbol{\theta})) \end{aligned} \quad (11)$$

- “Score matching” distance

$$\tilde{D}_{SM}(q_{\mathbf{v}}(\boldsymbol{\theta})\|p(\boldsymbol{\theta}, \mathcal{D})) = \frac{1}{2} \int q_{\mathbf{v}}(\boldsymbol{\theta}) \|\nabla_{\boldsymbol{\theta}} z(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta})\|^2 d\boldsymbol{\theta} \quad (12)$$

The above distances are well-defined

$$\begin{aligned} D_{PM}(q_{\mathbf{v}}(\boldsymbol{\theta})\|p(\boldsymbol{\theta}, \mathcal{D})) = 0 \text{ or } \tilde{D}_{SM}(q_{\mathbf{v}}(\boldsymbol{\theta})\|p(\boldsymbol{\theta}, \mathcal{D})) = 0 \\ \Rightarrow z(\boldsymbol{\theta}) = U(\boldsymbol{\theta}) + \text{Constant} \Rightarrow q_{\mathbf{v}}(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\mathcal{D}) \end{aligned}$$

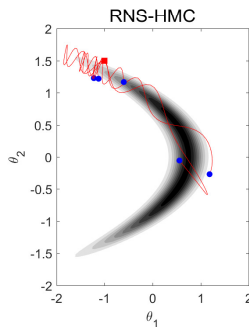
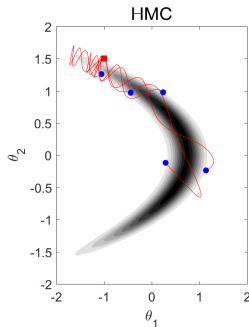
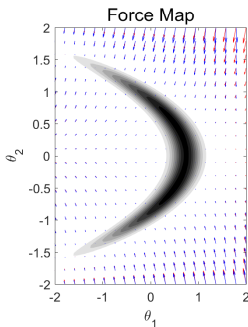
In practice, (12) is usually intractable, use **empirical version** instead.

Surrogate Induced Hamiltonian Flow

Define $\tilde{H}(\theta, \mathbf{r}) = z(\theta) + K(\mathbf{r})$, simulating the surrogate induced Hamiltonian Dynamics

$$\frac{d\theta}{dt} = M^{-1}\mathbf{r}, \quad \frac{d\mathbf{r}}{dt} = -\nabla_{\theta}z(\theta) \quad (13)$$

defines a mapping $\phi_S^{\tilde{H}} : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$, $(\theta, \mathbf{r}) \mapsto (\theta^*, \mathbf{r}^*)$



Variational Hamiltonian Monte Carlo

- 1 Simulate surrogate induced Hamiltonian flow to generate (θ^*, \mathbf{r}^*) and accept with probability

$$\alpha_{vhmc} = \min\{1, \exp(\tilde{H}(\theta, \mathbf{r}) - \tilde{H}(\theta^*, \mathbf{r}^*))\}$$

- 2 Add to the training data set.

$$\mathcal{T}_s^{(t)} := \mathcal{T}_s^{(t-1)} \cup \{(\theta_t, \nabla_{\theta} U(\theta_t))\}$$

- 3 Update surrogate by minimizing the empirical squared distance plus regularization.

$$\hat{\mathbf{v}}_t = \arg \min_{\mathbf{v}} \frac{1}{2} \sum_{n=1}^t \|\nabla_{\theta} z(\theta_n) - \nabla_{\theta} U(\theta_n)\|^2 + \frac{\lambda}{2} \|\mathbf{v}\|^2 \quad (14)$$

Regularized surrogate approximation to simulate Hamiltonian flow

$$V_t(\theta) = \mu_t z_t(\theta) + \frac{1}{2} (1 - \mu_t) (\theta - \theta^L)^\top \nabla_{\theta}^2 U(\theta^L) (\theta - \theta^L), \quad \mu_t : 0 \uparrow 1$$

Online Updating

Let $A(\theta) = (A_1(\theta), A_2(\theta), \dots, A_s(\theta))$, where

$$A_i(\theta) = \nabla_{\theta} a(\theta; \gamma_i), \quad i = 1, \dots, s$$

Variational Hamiltonian Monte Carlo

- 1: Set λ, μ_t, s and HMC parameters ε, L . Initialize $\theta^{(0)}$ to a first guess, $\mathbf{v}^{(0)} = \mathbf{0}$, $C^{(0)} = \frac{1}{\lambda} I_s$. Find θ^L and compute $\nabla_{\theta}^2 U(\theta)^L$.
- 2: **for** $t = 1$ to T **do**
- 3: Perform one HMC iteration for the regularized surrogate induced distribution $q_{\mathbf{v}}^{(t)}(\theta) \propto \exp(-V_t(\theta))$ to draw $(\theta^{(t+1)}, \mathbf{r}^{(t+1)})$
- 4: Acquire $\nabla_{\theta} U(\theta^{(t+1)})$ and $A_{t+1} = A(\theta^{(t+1)})$
- 5: Compute $W^{(t+1)} = C^{(t)} A_{t+1}^T [I_d + A_{t+1} C^{(t)} A_{t+1}^T]^{-1}$
- 6: Update $\mathbf{v}^{(t+1)}, C^{(t+1)}$ as follows
- 7: $\mathbf{v}^{(t+1)} = \mathbf{v}^{(t)} + W^{(t+1)} (\nabla_{\theta} U(\theta^{(t+1)}) - A_{t+1} \mathbf{v}^{(t)})$
- 8: $C^{(t+1)} = C^{(t)} - W^{(t+1)} A_{t+1} C^{(t)}$
- 9: **end for**

Stopping time and Connections to Existing Work

- Inefficient updating when surrogate is well trained \Rightarrow set up a stopping time t_0
 - $t < t_0$. Perform **Free-Form** variational Bayes to improve approximate quality of the surrogate
 - $t > t_0$. Perform standard **HMC** that samples from the surrogate induced distribution.
- Connections to some existing work
 - Stochastic linear regression for **Fix-Form** variational Bayes (SALIMANS et KNOWLES 2013) : allows **Free-Form** intractable approximate distributions, use HMC to draw samples.
 - Random bases surrogate HMC (ZHANG, SHAHBABA et ZHAO 2015) : further reduce the computation by allowing to use the fast surrogate in the Metropolis correction step, trade-off between approximation accuracy and computational efficiency.

A Beta-binomial Model for Overdispersion

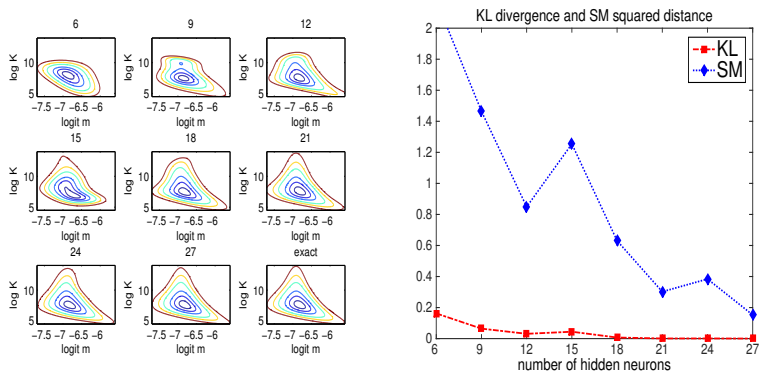


FIGURE – Left : Approximate posteriors for a varying number of hidden neurons. Exact posterior at bottom right. **Right :** KL-divergence and score matching squared distance between the surrogate approximation and the exact posterior density using an increasing number of hidden neurons.

Bayesian Probit Regression

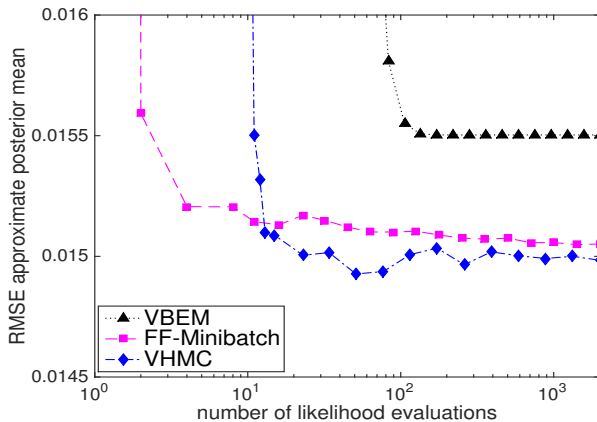


FIGURE – RMSE of the approximate posterior mean as a function of the number of likelihood evaluations for different variational Bayesian approaches and VHM algorithm.

Independent Component Analysis

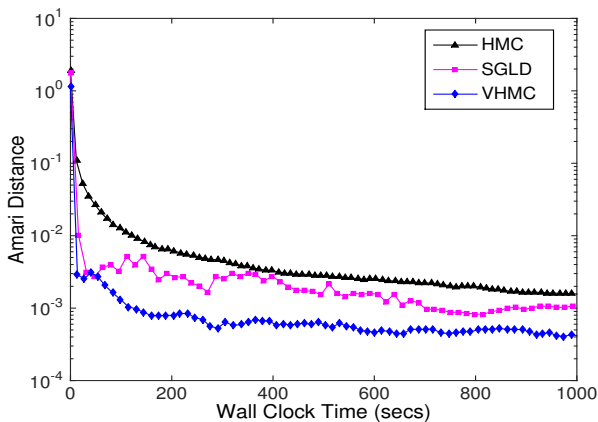


FIGURE – Convergence of Amari distance on the MEG data for HMC, SGLD and our variational HMC algorithm.

Summary and Discussion

Combine variational Bayes and MCMC via random bases surrogate :

- Construct efficient surrogate to accelerate HMC

$$\frac{d\theta}{dt} = M^{-1}r, \quad \frac{dr}{dt} = -\nabla_{\theta}z(\theta)$$








- Find good surrogate via **Free-Form** variational Bayes

$$\hat{\nu} = \arg \min_{\nu \in \Omega} \tilde{D}_{SM}(q_{\nu}(\theta) \| p(\theta, \mathcal{D}))$$








Open future directions :

- The random bases surrogate is more effective in problems with costly likelihood and a moderate number of parameters. How about really high dimensional problems ? Would more sophisticated structures (e.g., deep networks) work ?
- Evaluating full-data gradient $\nabla_{\theta}U(\theta)$ to collect training data is still expensive, how about stochastic gradient ? Be careful about overfitting !

Reference I

-  BETANCOURT, M. (2015). “The fundamental incompatibility of scalable Hamiltonian Monte Carlo and naive data subsampling”. In : **Proceedings of the 32nd International Conference on Machine Learning (ICML 2015)**.
-  CHEN, T., E. B. FOX et C. GUESTRIN (2014). “Stochastic Gradient Hamiltonian Monte Carlo”. In : **Proceedings of 31st International Conference on Machine Learning (ICML 2014)**.
-  DING, N. et al. (2014). “Bayesian Sampling Using Stochastic Gradient Thermostats”. In : **Advances in Neural Information Processing Systems 27 (NIPS 2014)**.
-  DUANE, S. et al. (1987). “Hybrid Monte Carlo”. In : **Physics Letters B** 195.2, p. 216–222.
-  GIROLAMI, M. et B. CALDERHEAD (2011). “Riemann manifold Langevin and Hamiltonian Monte Carlo methods”. In : **Journal of the Royal Statistical Society (with discussion)** 73.2, p. 123–214.
-  HOFFMAN, D. et A. GELMAN (2011). **The No-U-Turn Sampler : Adaptively Setting Path Lengths in Hamiltonian Monte Carlo**. arxiv.org/abs/1111.4246.
-  HONKELA, A. et al. (2010). “Approximate Riemannian conjugate gradient learning for fixed-form variational Bayes”. In : **Journal of Machine Learning Research** 11, p. 3235–3268.

Reference II

- 
 LAN, S. et al. (2015). **Emulation of higher-order tensors in manifold Monte Carlo methods for Bayesian inverse problems**. arxiv.org/abs/1507.06244.
- 
 MA, Y. A., T. CHEN et E. FOX (2015). “A complete recipe for stochastic gradient MCMC”. In : **Advances in Neural Information Processing Systems 28 (NIPS 2015)**.
- 
 METROPOLIS, N. et al. (1953). “Equation of State Calculations by Fast Computing Machines”. In : **The Journal of Chemical Physics** 21.6, p. 1087–1092.
- 
 NEAL, R. M. (2011). “MCMC using Hamiltonian dynamics”. In : **Handbook of Markov Chain Monte Carlo**. Sous la dir. de S. BROOKS et al. Chapman et Hall/CRC, p. 113–162.
- 
 RAHIMI, A. et B. RECHT (2008). “Uniform approximation of functions with random bases”. In : **Proc. 46th Ann. Allerton Conf. Commun., Contr. Comput.**
- 
 RASMUSSEN, C. E. (2003). “Gaussian Processes to Speed up Hybrid Monte Carlo for Expensive Bayesian Integrals”. In : **Bayesian Statistics** 7, p. 651–659.
- 
 SALIMANS, T. et D. A. KNOWLES (2013). “Fixed-form variational posterior approximation through stochastic linear regression”. In : **Bayesian Analysis** 8.4, p. 837–882.

Reference III



SAUL, L. et M. I. JORDAN (1996). “Exploiting tractable substructures in intractable networks”. In : **Advance in neural information processing systems 7 (NIPS 1996)**. Sous la dir. de G. TESAURO, D. S. TOURETZKY et T. K. LEEN. Cambridge, MA : MIT Press, p. 486–492.



STRATHMANN, H. et al. (2015). “Gradient-free Hamiltonian Monte Carlo with efficient kernel exponential families”. In : **Advances in Neural Information Processing Systems**. Cambridge, MA : MIT Press.



WANG, Z., S. MOHAMED et D. NANDO (2013). “Adaptive Hamiltonian and Riemann manifold Monte Carlo”. In : **Proceedings of the 30th International Conference on Machine Learning (ICML 2013)**, p. 1462–1470.



WELLING, M. et Y. W. TEH (2011). “Bayesian Learning via Stochastic Gradient Langevin Dynamics”. In : **Proceedings of the International Conference on Machine Learning**.



ZHANG, C., B. SHAHBABA et H. K. ZHAO (2015). **Hamiltonian Monte Carlo Acceleration Using Surrogate Functions with Random Bases**.
arxiv.org/abs/1506.05555.