

Introduction to Stochastic Gradient Markov Chain Monte Carlo Methods

Changyou Chen

Department of Electrical and Computer Engineering, Duke University
cc448@duke.edu

Duke-Tsinghua Machine Learning Summer School
August 10, 2016

Stochastic gradient Markov chain Monte Carlo (SG-MCMC):

- A new technique for approximate Bayesian sampling.
- It is about **scalable** Bayesian learning for **big data**.
- It draws samples $\{\theta\}$'s from $p(\theta; \mathbf{D})$ where $p(\theta; \mathbf{D})$ is too expensive to be evaluated in each iteration.

This lecture:

- Will cover: basic ideas behind SG-MCMC.
- Will not cover: different kinds of SG-MCMC algorithms, applications, and the corresponding convergence theory.

- 1 Markov Chain Monte Carlo Methods
 - Monte Carlo methods
 - Markov chain Monte Carlo
- 2 Stochastic Gradient Markov Chain Monte Carlo Methods
 - Introduction
 - Stochastic gradient Langevin dynamics
 - Stochastic gradient Hamiltonian Monte Carlo
 - Application in Latent Dirichlet allocation

1 Markov Chain Monte Carlo Methods

- Monte Carlo methods
- Markov chain Monte Carlo

2 Stochastic Gradient Markov Chain Monte Carlo Methods

- Introduction
- Stochastic gradient Langevin dynamics
- Stochastic gradient Hamiltonian Monte Carlo
- Application in Latent Dirichlet allocation

Monte Carlo methods

- Monte Carlo method is about drawing a set of samples from $p(\theta)$:

$$\theta_l \sim p(\theta), \quad l = 1, 2, \dots, L$$

- Approximate the target distribution $p(\theta)$ as count frequency:

$$p(\theta) \approx \frac{1}{L} \sum_{l=1}^L \delta(\theta, \theta_l)$$

- An intractable integration is approximated as:

$$\int f(\theta) p(\theta) \approx \frac{1}{L} \sum_{l=1}^L f(\theta_l)$$

- In Bayesian modeling, $p(\theta)$ is usually a posterior distribution, the integral is a predicted quantity.

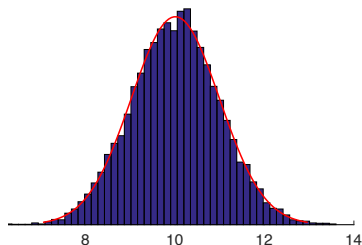
Monte Carlo methods

- Monte Carlo method is about drawing a set of samples from $p(\theta)$:

$$\theta_l \sim p(\theta), \quad l = 1, 2, \dots, L$$

- Approximate the target distribution $p(\theta)$ as count frequency:

$$p(\theta) \approx \frac{1}{L} \sum_{l=1}^L \delta(\theta, \theta_l)$$



- An intractable integration is approximated as:

$$\int f(\theta) p(\theta) \approx \frac{1}{L} \sum_{l=1}^L f(\theta_l)$$

- In Bayesian modeling, $p(\theta)$ is usually a posterior distribution, the integral is a predicted quantity.

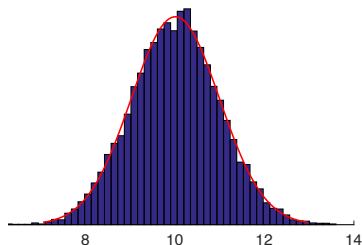
Monte Carlo methods

- Monte Carlo method is about drawing a set of samples from $p(\theta)$:

$$\theta_l \sim p(\theta), \quad l = 1, 2, \dots, L$$

- Approximate the target distribution $p(\theta)$ as count frequency:

$$p(\theta) \approx \frac{1}{L} \sum_{l=1}^L \delta(\theta, \theta_l)$$



- An intractable integration is approximated as:

$$\int f(\theta) p(\theta) \approx \frac{1}{L} \sum_{l=1}^L f(\theta_l)$$

- In Bayesian modeling, $p(\theta)$ is usually a posterior distribution, the integral is a predicted quantity.

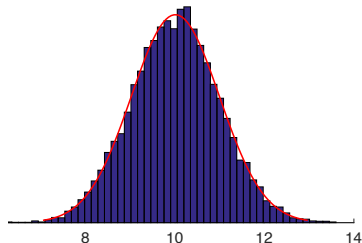
Monte Carlo methods

- Monte Carlo method is about drawing a set of samples from $p(\theta)$:

$$\theta_l \sim p(\theta), \quad l = 1, 2, \dots, L$$

- Approximate the target distribution $p(\theta)$ as count frequency:

$$p(\theta) \approx \frac{1}{L} \sum_{l=1}^L \delta(\theta, \theta_l)$$



- An intractable integration is approximated as:

$$\int f(\theta) p(\theta) \approx \frac{1}{L} \sum_{l=1}^L f(\theta_l)$$

- In Bayesian modeling, $p(\theta)$ is usually a posterior distribution, the integral is a predicted quantity.

How does the approximation work?

- 1 An intractable integration is approximated as:

$$\int f(\theta)p(\theta) \approx \frac{1}{L} \sum_{l=1}^L f(\theta_l) \triangleq \tilde{f}$$

- 2 If $\{\theta_l\}$'s are independent:

$$\mathbb{E}\tilde{f} = \mathbb{E}\left[\frac{1}{L} \sum_{l=1}^L f(\theta_l)\right] = \mathbb{E}f, \quad \text{Var}(\tilde{f}) = \text{Var}\left(\frac{1}{L} \sum_{l=1}^L f(\theta_l)\right) = \frac{1}{L} \text{Var}(f)$$

- ▶ the variance decreases linearly w.r.t. the number of samples, and independent of the dimension of θ

- 3 However, obtaining independent samples is hard:

- ▶ usually resort to drawing dependent samples with **Markov chain Monte Carlo** (MCMC)

How does the approximation work?

- ① An intractable integration is approximated as:

$$\int f(\theta)p(\theta) \approx \frac{1}{L} \sum_{l=1}^L f(\theta_l) \triangleq \tilde{f}$$

- ② If $\{\theta_l\}$'s are independent:

$$\mathbb{E}\tilde{f} = \mathbb{E}\left[\frac{1}{L} \sum_{l=1}^L f(\theta_l)\right] = \mathbb{E}f, \quad \text{Var}(\tilde{f}) = \text{Var}\left(\frac{1}{L} \sum_{l=1}^L f(\theta_l)\right) = \frac{1}{L} \text{Var}(f)$$

- ▶ the variance decreases linearly w.r.t. the number of samples, and independent of the dimension of θ

- ③ However, obtaining independent samples is hard:

- ▶ usually resort to drawing dependent samples with Markov chain Monte Carlo (MCMC)

How does the approximation work?

- ① An intractable integration is approximated as:

$$\int f(\theta)p(\theta) \approx \frac{1}{L} \sum_{l=1}^L f(\theta_l) \triangleq \tilde{f}$$

- ② If $\{\theta_l\}$'s are independent:

$$\mathbb{E}\tilde{f} = \mathbb{E}\left[\frac{1}{L} \sum_{l=1}^L f(\theta_l)\right] = \mathbb{E}f, \quad \text{Var}(\tilde{f}) = \text{Var}\left(\frac{1}{L} \sum_{l=1}^L f(\theta_l)\right) = \frac{1}{L} \text{Var}(f)$$

- ▶ the variance decreases linearly w.r.t. the number of samples, and independent of the dimension of θ

- ③ However, obtaining independent samples is hard:

- ▶ usually resort to drawing dependent samples with **Markov chain Monte Carlo (MCMC)**

1 Markov Chain Monte Carlo Methods

- Monte Carlo methods
- Markov chain Monte Carlo

2 Stochastic Gradient Markov Chain Monte Carlo Methods

- Introduction
- Stochastic gradient Langevin dynamics
- Stochastic gradient Hamiltonian Monte Carlo
- Application in Latent Dirichlet allocation

MCMC example: a Gaussian model

- 1 Assume the following generative process (with $\alpha = 5, \beta = 1$):

$$\begin{aligned}x_i | \mu, \tau &\sim N(\mu, 1/\tau), \quad i = 1, \dots, n = 1000 \\ \mu | \tau, \{x_i\} &\sim N(\mu_0, 1/\tau), \\ \tau &\sim \text{Gamma}(\alpha, \beta)\end{aligned}$$

- 2 Posterior distribution:

$$p(\mu, \tau | \{x_i\}) \propto \left[\prod_{i=1}^n N(x_i; \mu, 1/\tau) \right] N(\mu; \mu_0, 1/\tau) \text{Gamma}(\tau; \alpha, \beta)$$

- 3 Marginal posterior distributions for μ and τ are available:

$$\begin{aligned}p(\mu | \{x_i\}) &\propto \left(2\beta + (\mu - \mu_0)^2 + \sum_i (x_i - \mu)^2 \right)^{-\alpha - (n+1)/2} \\ p(\tau | \{x_i\}) &= \text{Gamma} \left(\alpha + \frac{n}{2}, \beta + \frac{1}{2} \sum_i (x_i - \bar{x})^2 + \frac{n}{2(n+1)} (\bar{x} - \mu_0)^2 \right)\end{aligned}$$

- $p(\mu | \{x_i\})$ is a non-standardized Student's t -distribution with mean $(\sum_i x_i + \mu_0)/(n+1)$

MCMC example: a Gaussian model

- 1 Assume the following generative process (with $\alpha = 5, \beta = 1$):

$$\begin{aligned}x_i | \mu, \tau &\sim N(\mu, 1/\tau), \quad i = 1, \dots, n = 1000 \\ \mu | \tau, \{x_i\} &\sim N(\mu_0, 1/\tau), \\ \tau &\sim \text{Gamma}(\alpha, \beta)\end{aligned}$$

- 2 Posterior distribution:

$$p(\mu, \tau | \{x_i\}) \propto \left[\prod_{i=1}^n N(x_i; \mu, 1/\tau) \right] N(\mu; \mu_0, 1/\tau) \text{Gamma}(\tau; \alpha, \beta)$$

- 3 Marginal posterior distributions for μ and τ are available:

$$\begin{aligned}p(\mu | \{x_i\}) &\propto \left(2\beta + (\mu - \mu_0)^2 + \sum_i (x_i - \mu)^2 \right)^{-\alpha - (n+1)/2} \\ p(\tau | \{x_i\}) &= \text{Gamma} \left(\alpha + \frac{n}{2}, \beta + \frac{1}{2} \sum_i (x_i - \bar{x})^2 + \frac{n}{2(n+1)} (\bar{x} - \mu_0)^2 \right)\end{aligned}$$

- $p(\mu | \{x_i\})$ is a non-standardized Student's t -distribution with mean $(\sum_i x_i + \mu_0)/(n+1)$

MCMC example: a Gaussian model

- 1 Assume the following generative process (with $\alpha = 5, \beta = 1$):

$$\begin{aligned}x_i | \mu, \tau &\sim N(\mu, 1/\tau), \quad i = 1, \dots, n = 1000 \\ \mu | \tau, \{x_i\} &\sim N(\mu_0, 1/\tau), \\ \tau &\sim \text{Gamma}(\alpha, \beta)\end{aligned}$$

- 2 Posterior distribution:

$$p(\mu, \tau | \{x_i\}) \propto \left[\prod_{i=1}^n N(x_i; \mu, 1/\tau) \right] N(\mu; \mu_0, 1/\tau) \text{Gamma}(\tau; \alpha, \beta)$$

- 3 Marginal posterior distributions for μ and τ are available:

$$\begin{aligned}p(\mu | \{x_i\}) &\propto \left(2\beta + (\mu - \mu_0)^2 + \sum_i (x_i - \mu)^2 \right)^{-\alpha - (n+1)/2} \\ p(\tau | \{x_i\}) &= \text{Gamma} \left(\alpha + \frac{n}{2}, \beta + \frac{1}{2} \sum_i (x_i - \bar{x})^2 + \frac{n}{2(n+1)} (\bar{x} - \mu_0)^2 \right)\end{aligned}$$

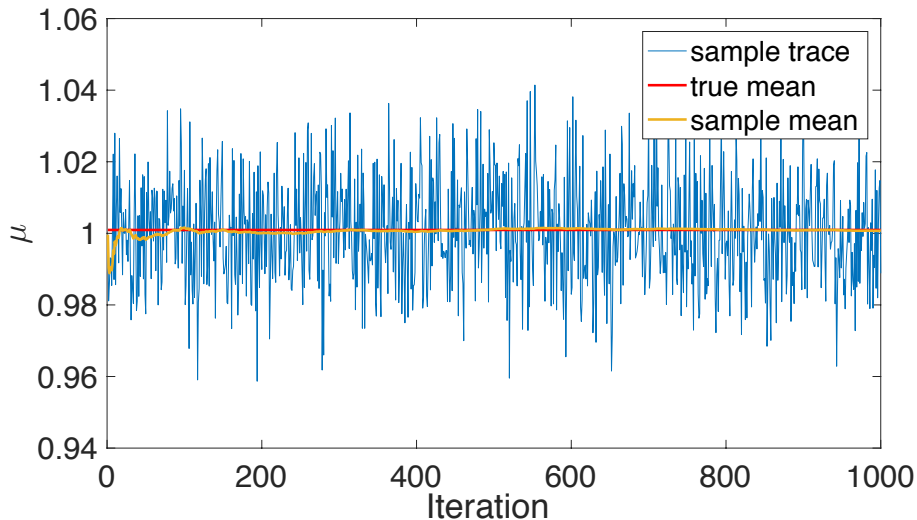
- $p(\mu | \{x_i\})$ is a non-standardized Student's t -distribution with mean $(\sum_i x_i + \mu_0)/(n+1)$

1 Conditional distributions:

$$\mu|\tau, \{x_i\} \sim N\left(\frac{n}{n+1}\bar{x} + \frac{1}{n+1}\mu_0, \frac{1}{(n+1)\tau}\right)$$

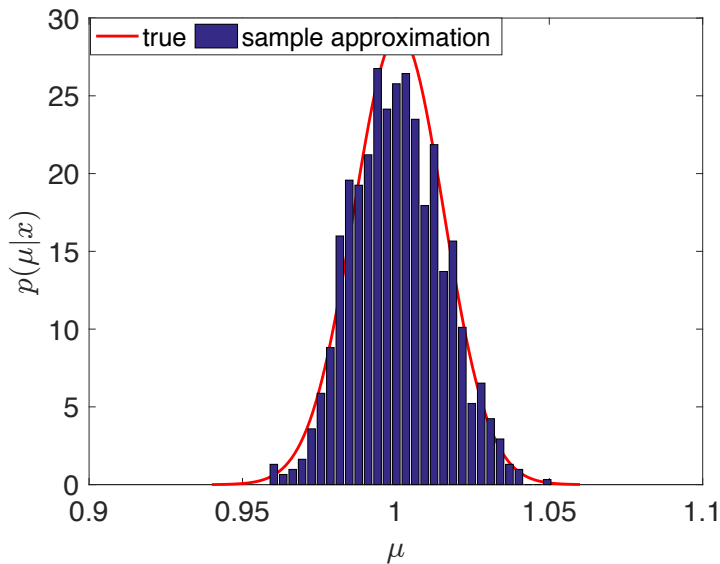
$$\tau|\mu, \{x_i\} \sim \text{Gamma}\left(\alpha + \frac{n+1}{2}, \beta + \frac{\sum_i (x_i - \mu)^2 + (\mu - \mu_0)^2}{2}\right)$$

Trace plot for μ

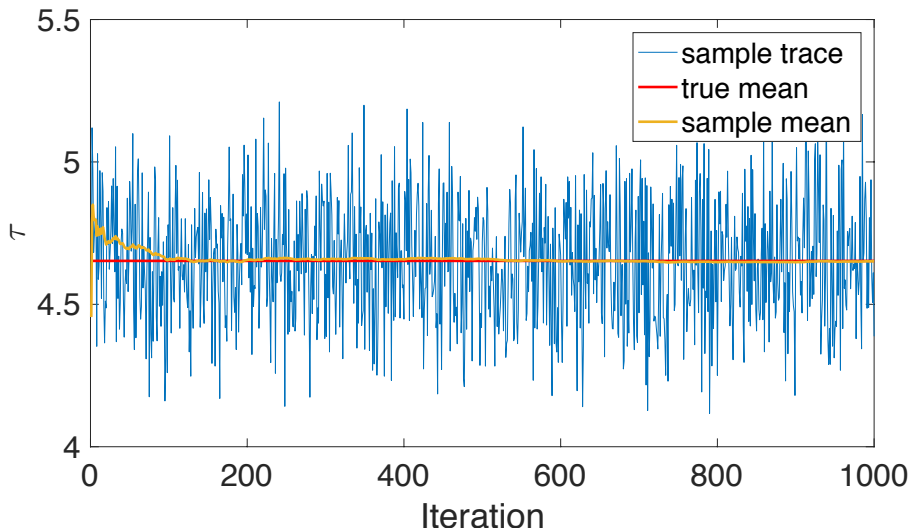


Sample approximation for μ

- True posterior is a non-standardized Student's t -distribution.

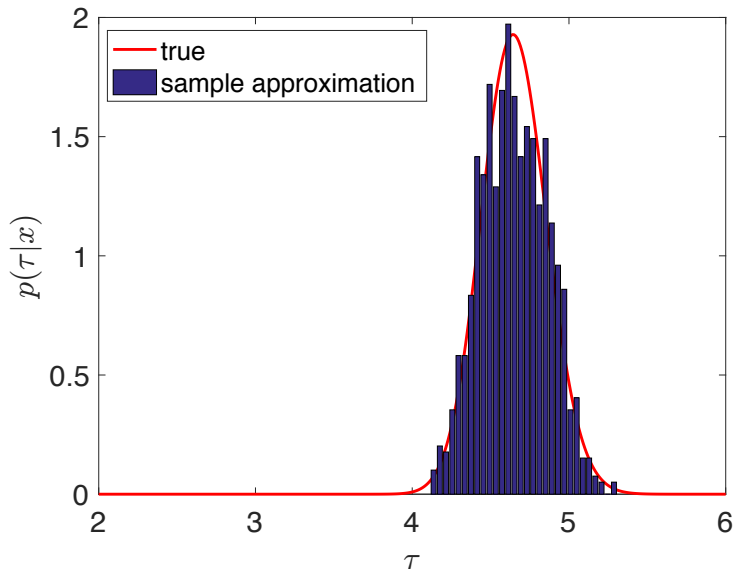


Trace plot for τ



Sample approximation for τ

- True posterior is a Gamma distribution.



Markov chain Monte Carlo methods

- 1 We are interested in drawing samples from some desired distribution $p^*(\theta) = \frac{1}{Z}\tilde{p}^*(\theta)$.
- 2 Define a Markov chain:

$$\theta_0 \rightarrow \theta_1 \rightarrow \theta_2 \rightarrow \theta_3 \rightarrow \theta_4 \rightarrow \theta_5 \rightarrow \cdots$$

where $\theta_0 \sim p_0(\theta)$, $\theta_1 \sim p_1(\theta)$, \cdots , satisfying

$$p_t(\theta') = \int p_{t-1}(\theta) T(\theta \rightarrow \theta') d\theta,$$

where $T(\theta \rightarrow \theta')$ is the Markov chain transition probability from θ to θ' .

- 3 We say $p^*(\theta)$ is an invariant (stationary) distribution of the Markov chain iff:

$$p^*(\theta') = \int p^*(\theta) T(\theta \rightarrow \theta') d\theta$$

Markov chain Monte Carlo methods

- 1 We are interested in drawing samples from some desired distribution $p^*(\theta) = \frac{1}{Z} \tilde{p}^*(\theta)$.
- 2 Define a Markov chain:

$$\theta_0 \rightarrow \theta_1 \rightarrow \theta_2 \rightarrow \theta_3 \rightarrow \theta_4 \rightarrow \theta_5 \rightarrow \dots$$

where $\theta_0 \sim p_0(\theta)$, $\theta_1 \sim p_1(\theta)$, \dots , satisfying

$$p_t(\theta') = \int p_{t-1}(\theta) T(\theta \rightarrow \theta') d\theta,$$

where $T(\theta \rightarrow \theta')$ is the Markov chain transition probability from θ to θ' .

- 3 We say $p^*(\theta)$ is an invariant (stationary) distribution of the Markov chain iff:

$$p^*(\theta') = \int p^*(\theta) T(\theta \rightarrow \theta') d\theta$$

Markov chain Monte Carlo methods

- 1 We are interested in drawing samples from some desired distribution $p^*(\theta) = \frac{1}{Z} \tilde{p}^*(\theta)$.
- 2 Define a Markov chain:

$$\theta_0 \rightarrow \theta_1 \rightarrow \theta_2 \rightarrow \theta_3 \rightarrow \theta_4 \rightarrow \theta_5 \rightarrow \dots$$

where $\theta_0 \sim p_0(\theta)$, $\theta_1 \sim p_1(\theta)$, \dots , satisfying

$$p_t(\theta') = \int p_{t-1}(\theta) T(\theta \rightarrow \theta') d\theta,$$

where $T(\theta \rightarrow \theta')$ is the Markov chain transition probability from θ to θ' .

- 3 We say $p^*(\theta)$ is an invariant (stationary) distribution of the Markov chain iff:

$$p^*(\theta') = \int p^*(\theta) T(\theta \rightarrow \theta') d\theta$$

Metropolis-Hasting algorithm

- 1 Design $T(\theta \rightarrow \theta')$ as the composition of a proposal distribution $q_t(\theta' | \theta)$ and an accept-reject mechanism.
- 2 At step t , draw a sample¹ $\theta^* \sim q_t(\theta | \theta_{t-1})$, and accept it with probability:

$$A_t(\theta^*, \theta_{t-1}) = \min \left(1, \frac{\tilde{p}(\theta^*) q_t(\theta_{t-1} | \theta^*)}{\tilde{p}(\theta_{t-1}) q_t(\theta^* | \theta_{t-1})} \right)$$

- 3 The acceptance can be done by:
 - ▶ draw a random variable $u \sim \text{Uniform}(0, 1)$
 - ▶ accept the sample if $A_t(\theta^*, \theta_{t-1}) > u$
- 4 The corresponding transition kernel satisfies the detailed balance condition, thus has an invariant probability $p^*(\theta)$.

¹ A standard setting of $q_t(\theta | \theta_{t-1})$ is a normal distribution with mean θ_{t-1} and tunable variance.

Metropolis-Hasting algorithm

- 1 Design $T(\theta \rightarrow \theta')$ as the composition of a proposal distribution $q_t(\theta' | \theta)$ and an accept-reject mechanism.
- 2 At step t , draw a sample¹ $\theta^* \sim q_t(\theta | \theta_{t-1})$, and accept it with probability:

$$A_t(\theta^*, \theta_{t-1}) = \min \left(1, \frac{\tilde{p}(\theta^*) q_t(\theta_{t-1} | \theta^*)}{\tilde{p}(\theta_{t-1}) q_t(\theta^* | \theta_{t-1})} \right)$$

- 3 The acceptance can be done by:
 - ▶ draw a random variable $u \sim \text{Uniform}(0, 1)$
 - ▶ accept the sample if $A_t(\theta^*, \theta_{t-1}) > u$
- 4 The corresponding transition kernel satisfies the detailed balance condition, thus has an invariant probability $p^*(\theta)$.

¹A standard setting of $q_t(\theta | \theta_{t-1})$ is a normal distribution with mean θ_{t-1} and tunable variance.

Metropolis-Hasting algorithm

- 1 Design $T(\theta \rightarrow \theta')$ as the composition of a proposal distribution $q_t(\theta' | \theta)$ and an accept-reject mechanism.
- 2 At step t , draw a sample¹ $\theta^* \sim q_t(\theta | \theta_{t-1})$, and accept it with probability:

$$A_t(\theta^*, \theta_{t-1}) = \min \left(1, \frac{\tilde{p}(\theta^*) q_t(\theta_{t-1} | \theta^*)}{\tilde{p}(\theta_{t-1}) q_t(\theta^* | \theta_{t-1})} \right)$$

- 3 The acceptance can be done by:
 - ▶ draw a random variable $u \sim \text{Uniform}(0, 1)$
 - ▶ accept the sample if $A_t(\theta^*, \theta_{t-1}) > u$
- 4 The corresponding transition kernel satisfies the detailed balance condition, thus has an invariant probability $p^*(\theta)$.

¹A standard setting of $q_t(\theta | \theta_{t-1})$ is a normal distribution with mean θ_{t-1} and tunable variance.

Metropolis-Hasting algorithm

- 1 Design $T(\theta \rightarrow \theta')$ as the composition of a proposal distribution $q_t(\theta' | \theta)$ and an accept-reject mechanism.
- 2 At step t , draw a sample¹ $\theta^* \sim q_t(\theta | \theta_{t-1})$, and accept it with probability:

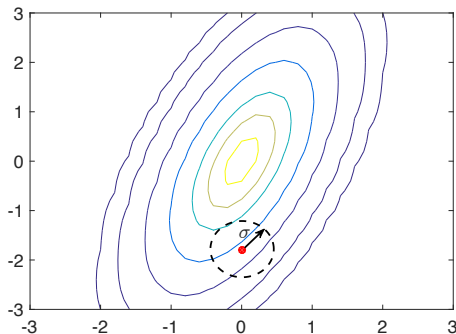
$$A_t(\theta^*, \theta_{t-1}) = \min \left(1, \frac{\tilde{p}(\theta^*) q_t(\theta_{t-1} | \theta^*)}{\tilde{p}(\theta_{t-1}) q_t(\theta^* | \theta_{t-1})} \right)$$

- 3 The acceptance can be done by:
 - ▶ draw a random variable $u \sim \text{Uniform}(0, 1)$
 - ▶ accept the sample if $A_t(\theta^*, \theta_{t-1}) > u$
- 4 The corresponding transition kernel satisfies the detailed balance condition, thus has an invariant probability $p^*(\theta)$.

¹ A standard setting of $q_t(\theta | \theta_{t-1})$ is a normal distribution with mean θ_{t-1} and tunable variance.

Discussion on the proposal distribution

- 1 Standard proposal distribution is an isotropic Gaussian center at the current state with variance σ :
 - ▶ small σ leads to high acceptance rate, but moves too slowly
 - ▶ large σ moves fast, but leads to high rejection rate
- 2 How to choose better proposals?



Gibbs sampler

- 1 Assume θ is multi-dimensional², $\theta = (\theta_1, \dots, \theta_k, \dots, \theta_K)$, denote $\theta_{-k} \triangleq \{\theta_j : j \neq k\}$.
- 2 Sample θ_k sequentially, with proposal distribution being the true conditional distribution:

$$q_k(\theta^* | \theta) = p(\theta_k^* | \theta_{-k})$$

- 3 Note $\theta_{-k}^* = \theta_{-k}$, $p(\theta) = p(\theta_k | \theta_{-k})p(\theta_{-k})$.
- 4 The MH acceptance probability is:

$$\begin{aligned} A(\theta^*, \theta) &= \frac{p(\theta^*)q_k(\theta | \theta^*)}{p(\theta)q_k(\theta^* | \theta)} = \frac{p(\theta_k^* | \theta_{-k}^*)p(\theta_{-k}^*)p(\theta_k | \theta_{-k}^*)}{p(\theta_k^* | \theta_{-k})p(\theta_{-k})p(\theta_k | \theta_{-k})} \\ &= 1 \end{aligned}$$

²One dimensional random variable is relatively easy to sample.

Gibbs sampler

- 1 Assume θ is multi-dimensional², $\theta = (\theta_1, \dots, \theta_k, \dots, \theta_K)$, denote $\theta_{-k} \triangleq \{\theta_j : j \neq k\}$.
- 2 Sample θ_k sequentially, with proposal distribution being the true **conditional distribution**:

$$q_k(\theta^* | \theta) = p(\theta_k^* | \theta_{-k})$$

- 3 Note $\theta_{-k}^* = \theta_{-k}$, $p(\theta) = p(\theta_k | \theta_{-k})p(\theta_{-k})$.
- 4 The MH acceptance probability is:

$$\begin{aligned} A(\theta^*, \theta) &= \frac{p(\theta^*)q_k(\theta | \theta^*)}{p(\theta)q_k(\theta^* | \theta)} = \frac{p(\theta_k^* | \theta_{-k}^*)p(\theta_{-k}^*)p(\theta_k | \theta_{-k}^*)}{p(\theta_k^* | \theta_{-k})p(\theta_{-k})p(\theta_k | \theta_{-k})} \\ &= 1 \end{aligned}$$

²One dimensional random variable is relatively easy to sample.

Gibbs sampler

- 1 Assume θ is multi-dimensional², $\theta = (\theta_1, \dots, \theta_k, \dots, \theta_K)$, denote $\theta_{-k} \triangleq \{\theta_j : j \neq k\}$.
- 2 Sample θ_k sequentially, with proposal distribution being the true conditional distribution:

$$q_k(\theta^* | \theta) = p(\theta_k^* | \theta_{-k})$$

- 3 Note $\theta_{-k}^* = \theta_{-k}$, $p(\theta) = p(\theta_k | \theta_{-k})p(\theta_{-k})$.
- 4 The MH acceptance probability is:

$$\begin{aligned} A(\theta^*, \theta) &= \frac{p(\theta^*)q_k(\theta | \theta^*)}{p(\theta)q_k(\theta^* | \theta)} = \frac{p(\theta_k^* | \theta_{-k}^*)p(\theta_{-k}^*)p(\theta_k | \theta_{-k}^*)}{p(\theta_k^* | \theta_{-k})p(\theta_{-k})p(\theta_k | \theta_{-k})} \\ &= 1 \end{aligned}$$

²One dimensional random variable is relatively easy to sample.

Discussion of Gibbs sampler

- 1 No accept-reject step, very efficient.
- 2 Conditional distributions are not always easy to sample.
- 3 May not mix well when in high-dimensional space with highly correlated variables.

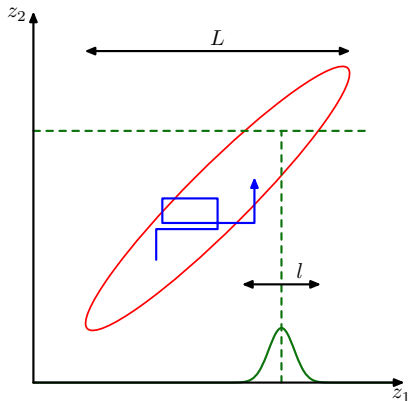
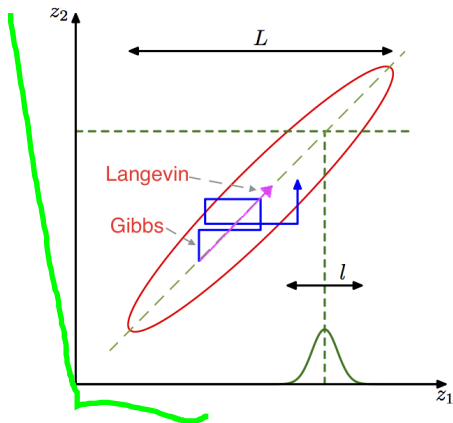


Figure: Sample path does not follow gradients. Figure from PRML, Bishop (2006)

The Metropolis-adjusted Langevin: a better proposal

- 1 Gibbs sampling travels the parameter space following a zipzag curve, which might be slow in high-dimensional space.
- 2 The Metropolis-adjusted Langevin uses a proposal that points directly to the center of the probabilistic contour.



The Metropolis-adjusted Langevin: a better proposal

- 1 Let $E(\theta) \triangleq -\log \tilde{p}(\theta)$, the direction of the contour is just the gradient: $-\nabla_{\theta} E(\theta)$.
- 2 In iteration l , define the proposal as a Gaussian centering at $\theta^* = \theta_{l-1} - \nabla_{\theta} E(\theta_{l-1})h_l$, where h_l is a small stepsize:

$$q(\theta_l | \theta_{l-1}) = N(\theta_l; \theta^*, \sigma^2) .$$


- 3 Need to do an accept-reject step:
 - ▶ calculate the acceptance probability:

$$A(\theta^*, \theta_{l-1}) = \frac{\tilde{p}(\theta^*)q(\theta_{l-1} | \theta^*)}{\tilde{p}(\theta)q(\theta^* | \theta_{l-1})}$$

- ▶ accept θ^* with probability $A(\theta^*, \theta_{l-1})$, otherwise set $\theta_l = \theta_{l-1}$

The Metropolis-adjusted Langevin: a better proposal

- 1 Let $E(\theta) \triangleq -\log \tilde{p}(\theta)$, the direction of the contour is just the gradient: $-\nabla_{\theta} E(\theta)$.
- 2 In iteration l , define the proposal as a Gaussian centering at $\theta^* = \theta_{l-1} - \nabla_{\theta} E(\theta_{l-1})h_l$, where h_l is a small stepsize:

$$q(\theta_l | \theta_{l-1}) = N(\theta_l; \theta^*, \sigma^2) .$$


- 3 Need to do an accept-reject step:

- calculate the acceptance probability:

Basically: SGD step with a

$$A(\theta^*, \theta_{l-1}) = \frac{\tilde{p}(\theta^*)q(\theta_{l-1} | \theta^*)}{\tilde{p}(\theta)q(\theta^* | \theta_{l-1})}$$

- accept θ^* with probability $A(\theta^*, \theta_{l-1})$, otherwise set $\theta_l = \theta_{l-1}$

Hamiltonian Monte Carlo

Frictionless ball rolling:

- 1 A dynamic system with total energy or Hamiltonian:

$$H = E(\theta) + K(\mathbf{v}), \text{ where}$$

$$E(\theta) \triangleq -\log \tilde{p}(\theta),$$

$$K(\mathbf{v}) \triangleq \mathbf{v}^T \mathbf{v} / 2.$$

- 2 Hamiltonian's equation describes the equations of motion of the ball:

$$\frac{d\theta}{dt} = \frac{\partial H}{\partial \mathbf{v}} = \mathbf{v}$$

$$\frac{d\mathbf{v}}{dt} = -\frac{\partial H}{\partial \theta} = \frac{\partial \log \tilde{p}(\theta)}{\partial \theta}$$

- 3 Joint distribution:
 $p(\theta, \mathbf{v}) \propto e^{-H(\theta, \mathbf{v})}.$

Figure: Rolling ball. Movie from Matthias Liepe

Hamiltonian Monte Carlo

Frictionless ball rolling:

- 1 A dynamic system with total energy or Hamiltonian:

$$H = E(\theta) + K(\mathbf{v}), \text{ where}$$

$$E(\theta) \triangleq -\log \tilde{p}(\theta),$$

$$K(\mathbf{v}) \triangleq \mathbf{v}^T \mathbf{v} / 2.$$

- 2 Hamiltonian's equation describes the equations of motion of the ball:

$$\frac{d\theta}{dt} = \frac{\partial H}{\partial \mathbf{v}} = \mathbf{v}$$

$$\frac{d\mathbf{v}}{dt} = -\frac{\partial H}{\partial \theta} = \frac{\partial \log \tilde{p}(\theta)}{\partial \theta}$$

- 3 Joint distribution:
 $p(\theta, \mathbf{v}) \propto e^{-H(\theta, \mathbf{v})}.$

Figure: Rolling ball. Movie from Matthias Liepe

Hamiltonian Monte Carlo

Frictionless ball rolling:

- 1 A dynamic system with total energy or Hamiltonian:

$$H = E(\theta) + K(\mathbf{v}), \text{ where}$$
$$E(\theta) \triangleq -\log \tilde{p}(\theta),$$
$$K(\mathbf{v}) \triangleq \mathbf{v}^T \mathbf{v} / 2.$$

- 2 Hamiltonian's equation describes the equations of motion of the ball:

$$\frac{d\theta}{dt} = \frac{\partial H}{\partial \mathbf{v}} = \mathbf{v}$$
$$\frac{d\mathbf{v}}{dt} = -\frac{\partial H}{\partial \theta} = \frac{\partial \log \tilde{p}(\theta)}{\partial \theta}$$

- 3 Joint distribution:
 $p(\theta, \mathbf{v}) \propto e^{-H(\theta, \mathbf{v})}.$

Figure: Rolling ball. Movie from Matthias Liepe

Solving Hamiltonian dynamics

- 1 Solving the continuous-time differential equation with discretized-time approximation:

$$\begin{cases} d\boldsymbol{\theta} = \mathbf{v} dt \\ d\mathbf{v} = \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\boldsymbol{\theta}) dt \end{cases} \implies \begin{cases} \boldsymbol{\theta}_l = \boldsymbol{\theta}_{l-1} + \mathbf{v}_{l-1} h_l \\ \mathbf{v}_l = \mathbf{v}_{l-1} + \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\boldsymbol{\theta}_l) h_l \end{cases}$$

- ▶ proposals follow historical gradients of the distribution contour

- 2 Need an accept-reject test to design whether accept the proposal, because of the discretization error:

- ▶ proposal is deterministic
- ▶ acceptance probability: $\min(1, \exp\{H(\boldsymbol{\theta}_l, \mathbf{v}_l) - H(\boldsymbol{\theta}_{l+1}, \mathbf{v}_{l+1})\})$

- 3 Almost identical to SGD with momentum:

- ▶
$$\begin{cases} \boldsymbol{\theta}_l = \boldsymbol{\theta}_{l-1} + \mathbf{p}_{l-1} \\ \mathbf{p}_l = (1 - m) \mathbf{p}_{l-1} + \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\boldsymbol{\theta}_l) \epsilon_l \end{cases}$$
- ▶ they will be make equivalent in the context of stochastic gradient MCMC

Solving Hamiltonian dynamics

- 1 Solving the continuous-time differential equation with discretized-time approximation:

$$\begin{cases} d\boldsymbol{\theta} = \mathbf{v} dt \\ d\mathbf{v} = \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\boldsymbol{\theta}) dt \end{cases} \implies \begin{cases} \boldsymbol{\theta}_l = \boldsymbol{\theta}_{l-1} + \mathbf{v}_{l-1} h_l \\ \mathbf{v}_l = \mathbf{v}_{l-1} + \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\boldsymbol{\theta}_l) h_l \end{cases}$$

- ▶ proposals follow historical gradients of the distribution contour
- 2 Need an accept-reject test to design whether accept the proposal, because of the discretization error:
 - ▶ proposal is deterministic
 - ▶ acceptance probability: $\min(1, \exp\{H(\boldsymbol{\theta}_l, \mathbf{v}_l) - H(\boldsymbol{\theta}_{l+1}, \mathbf{v}_{l+1})\})$
 - 3 Almost identical to SGD with momentum:
 - ▶
$$\begin{cases} \boldsymbol{\theta}_l = \boldsymbol{\theta}_{l-1} + \mathbf{p}_{l-1} \\ \mathbf{p}_l = (1 - m) \mathbf{p}_{l-1} + \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\boldsymbol{\theta}_l) \epsilon_l \end{cases}$$
 - ▶ they will be make equivalent in the context of stochastic gradient MCMC



- 1 Solving the continuous-time differential equation with discretized-time approximation:

$$\begin{cases} d\boldsymbol{\theta} = \mathbf{v} dt \\ d\mathbf{v} = \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\boldsymbol{\theta}) dt \end{cases} \implies \begin{cases} \boldsymbol{\theta}_l = \boldsymbol{\theta}_{l-1} + \mathbf{v}_{l-1} h_l \\ \mathbf{v}_l = \mathbf{v}_{l-1} + \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\boldsymbol{\theta}_l) h_l \end{cases}$$

- ▶ proposals follow historical gradients of the distribution contour
- 2 Need an accept-reject test to design whether accept the proposal, because of the discretization error:
 - ▶ proposal is deterministic
 - ▶ acceptance probability: $\min(1, \exp\{H(\boldsymbol{\theta}_l, \mathbf{v}_l) - H(\boldsymbol{\theta}_{l+1}, \mathbf{v}_{l+1})\})$
 - 3 Almost identical to SGD with momentum:
 - ▶ $\begin{cases} \boldsymbol{\theta}_l = \boldsymbol{\theta}_{l-1} + \mathbf{p}_{l-1} \\ \mathbf{p}_l = (1 - m) \mathbf{p}_{l-1} + \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\boldsymbol{\theta}_l) \epsilon_l \end{cases}$
 - ▶ they will be make equivalent in the context of stochastic gradient MCMC

Demo: MH vs. HMC

- 1 Nine mixtures of Gaussians³.
- 2 Sequential of samples connected by yellow lines.

³Demo by T. Broderick and D. Duvenaud.

Recap

- 1 Bayesian sampling with traditional MCMC methods, in each iteration:
 - ▶ generate a candidate sample from a proposal distribution
 - ▶ calculate the acceptance probability
 - ▶ accept or reject the proposed sample



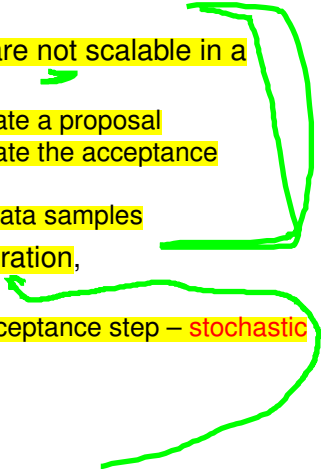
Discussion

- ❶ All the above traditional MCMC methods are not scalable in a big-data setting⁴, in each iteration:
 - ▶ the whole data need to be used to generate a proposal
 - ▶ the whole data need to be used to calculate the acceptance probability
 - ▶ scales $O(N)$, where N is the number of data samples
- ❷ Scalable MCMC uses sub-data in each iteration,
 - ▶ to calculate the acceptance probability⁵
 - ▶ to generate proposals, and ignore the acceptance step – **stochastic gradient MCMC methods (SG-MCMC)**

⁴when the number of data samples are large.

⁵A. Korattikara, Y. Chen, and M. Welling. "Austerity in MCMC Land: Cutting the Metropolis-Hastings Budget". In: *ICML*. 2014; R. Bardenet, A. Doucet, and C. Holmes. "Towards scaling up Markov chain Monte Carlo: an adaptive subsampling approach". In: *ICML*. 2014.

Discussion

- ① All the above traditional MCMC methods are not scalable in a big-data setting⁴, in each iteration:
 - ▶ the whole data need to be used to generate a proposal
 - ▶ the whole data need to be used to calculate the acceptance probability
 - ▶ scales $O(N)$, where N is the number of data samples
 - ② Scalable MCMC uses sub-data in each iteration,
 - ▶ to calculate the acceptance probability⁵
 - ▶ to generate proposals, and ignore the acceptance step – **stochastic gradient MCMC methods (SG-MCMC)**
- 

mini-batches

⁴when the number of data samples are large.

⁵A. Korattikara, Y. Chen, and M. Welling. "Austerity in MCMC Land: Cutting the Metropolis-Hastings Budget". In: *ICML. 2014*; R. Bardenet, A. Doucet, and C. Holmes. "Towards scaling up Markov chain Monte Carlo: an adaptive subsampling approach". In: *ICML. 2014*.

1 Markov Chain Monte Carlo Methods

- Monte Carlo methods
- Markov chain Monte Carlo

2 Stochastic Gradient Markov Chain Monte Carlo Methods

- Introduction
- Stochastic gradient Langevin dynamics
- Stochastic gradient Hamiltonian Monte Carlo
- Application in Latent Dirichlet allocation

Two key steps in SG-MCMC

- 1 Proposals typically follow stochastic gradients of log-posteriors:
 - ▶ make samples concentrate on the modes
- 2 Adding random Gaussian noise to proposals.
 - ▶ encourage algorithms to jump out of local modes, and to explore the parameter space
 - ▶ the noise in stochastic gradients not sufficient to make the algorithm move around parameter space

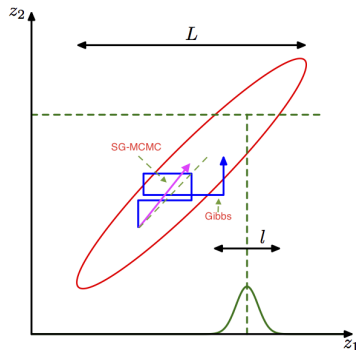


Figure: Proposals of Gibbs and SG-MCMC.

Two key steps in SG-MCMC

1 Proposals typically follow stochastic gradients of log-posteriors:

- ▶ make samples concentrate on the modes

2 Adding random Gaussian noise to proposals.

- ▶ encourage algorithms to jump out of local modes, and to explore the parameter space
- ▶ the noise in stochastic gradients not sufficient to make the algorithm move around parameter space

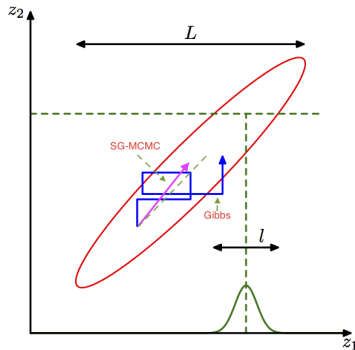


Figure: Proposals of Gibbs and SG-MCMC.

Basic setup

- 1 Given data $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, a generative model (likelihood) $p(\mathbf{X} | \theta) = \prod_{i=1}^N p(\mathbf{x}_i | \theta)$ and prior $p(\theta)$, we want to sample from the posterior:

$$p(\theta | \mathbf{X}) \propto p(\theta) p(\mathbf{X} | \theta) = p(\theta) \prod_{i=1}^N p(\mathbf{x}_i | \theta)$$

- 2 We are interested in the case when N is extremely large, so that computing $p(\mathbf{X} | \theta)$ is prohibitively expensive.
- 3 Define the following two quantities (unnormalized log-posterior and stochastic unnormalized log-posterior):

$$U(\theta) \triangleq - \sum_{i=1}^N \log p(\mathbf{x}_i | \theta) - \log p(\theta)$$

$$\tilde{U}(\theta) \triangleq - \frac{N}{n} \sum_{i=1}^n \log p(\mathbf{x}_{\pi_i} | \theta) - \log p(\theta)$$

where (π_1, \dots, π_N) is a random permutation of $(1, \dots, N)$.

Basic setup

- 1 Given data $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, a generative model (likelihood) $p(\mathbf{X} | \theta) = \prod_{i=1}^N p(\mathbf{x}_i | \theta)$ and prior $p(\theta)$, we want to sample from the posterior:

$$p(\theta | \mathbf{X}) \propto p(\theta) p(\mathbf{X} | \theta) = p(\theta) \prod_{i=1}^N p(\mathbf{x}_i | \theta)$$

- 2 We are interested in the case when N is extremely large, so that computing $p(\mathbf{X} | \theta)$ is prohibitively expensive.
- 3 Define the following two quantities (unnormalized log-posterior and **stochastic unnormalized log-posterior**):

$$U(\theta) \triangleq - \sum_{i=1}^N \log p(\mathbf{x}_i | \theta) - \log p(\theta)$$

$$\tilde{U}(\theta) \triangleq - \frac{N}{n} \sum_{i=1}^n \log p(\mathbf{x}_{\pi_i} | \theta) - \log p(\theta)$$

where (π_1, \dots, π_N) is a random permutation of $(1, \dots, N)$.

Basic setup

- 1 SG-MCMC relies on the following quantity (stochastic gradient):

$$\nabla_{\theta} \tilde{U}(\theta) \triangleq -\frac{N}{n} \sum_{i=1}^n \nabla_{\theta} \log p(\mathbf{x}_{\pi_i} | \theta) - \nabla_{\theta} \log p(\theta) ,$$

- 2 $\nabla_{\theta} \tilde{U}(\theta)$ is an unbiased estimate of $\nabla_{\theta} U(\theta)$:

- ▶ SG-MCMC samples parameters based on $\nabla_{\theta} \tilde{U}(\theta)$
- ▶ very cheap to compute
- ▶ bringing the name “stochastic gradient MCMC”

Basic setup

- ① SG-MCMC relies on the following quantity (stochastic gradient):

$$\nabla_{\theta} \tilde{U}(\theta) \triangleq -\frac{N}{n} \sum_{i=1}^n \nabla_{\theta} \log p(\mathbf{x}_{\pi_i} | \theta) - \nabla_{\theta} \log p(\theta),$$

- ② $\nabla_{\theta} \tilde{U}(\theta)$ is an unbiased estimate of $\nabla_{\theta} U(\theta)$: 

- ▶ SG-MCMC samples parameters based on $\nabla_{\theta} \tilde{U}(\theta)$
- ▶ very cheap to compute
- ▶ bringing the name “stochastic gradient MCMC”

Comparing with traditional MCMC

- 1 Ignore the acceptance step:
 - ▶ the detailed balance condition typically not hold, and the algorithm is not reversible⁶
 - ▶ typically leads to biased, but controllable estimations
- 2 Use sub-data in each iteration:
 - ▶ yielding stochastic gradients
 - ▶ does not affect the convergence properties (e.g., convergence rates), compared to using the whole data in each iteration

⁶These are sufficient conditions for a valid MCMC method, but not necessary conditions.

Comparing with traditional MCMC

1 Ignore the acceptance step:

- ▶ the detailed balance condition typically not hold, and the algorithm is not reversible⁶
- ▶ typically leads to biased, but controllable estimations

2 Use sub-data in each iteration:

- ▶ yielding stochastic gradients
- ▶ does not affect the convergence properties (e.g., convergence rates), compared to using the whole data in each iteration

⁶These are sufficient conditions for a valid MCMC method, but not necessary conditions.

Demo: the two key steps

- 1 Proposals follow stochastic gradients of log-posteriors:
 - ▶ stuck in a local mode

Demo: the two key steps

- 1 After adding random Gaussian noise:
 - it works !!

1 Markov Chain Monte Carlo Methods

- Monte Carlo methods
- Markov chain Monte Carlo

2 Stochastic Gradient Markov Chain Monte Carlo Methods

- Introduction
- Stochastic gradient Langevin dynamics
- Stochastic gradient Hamiltonian Monte Carlo
- Application in Latent Dirichlet allocation

First attempt

- 1 A 1st-order method: stochastic gradients directly applied on the model parameter θ .
- 2 Use a proposal that follows the stochastic gradient of the log-posterior:

$$\theta_{l+1} = \theta_l - h_{l+1} \nabla_{\theta} \tilde{U}(\theta_l)$$

- ▶ h_l 's are the stepsizes, could be fixed ($\forall l, h_l = h$) or decreasing ($\forall l, h_l > h_{l+1}$)

- 3 Ignore the acceptance step.
- 4 Resulting in Stochastic Gradient Descent (SGD).

Random noise to the rescue

- 1 Need to make the algorithm explore the parameter space:
 - ▶ adding random Gaussian noise to the update⁷

$$\begin{aligned}\theta_{l+1} &= \theta_l - h_{l+1} \nabla_{\theta} \tilde{U}(\theta_l) + \sqrt{2h_{l+1}} \zeta_{l+1} \\ \zeta_{l+1} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I})\end{aligned}$$

- 2 The magnitude of the Gaussian needs to be $\sqrt{2h_{l+1}}$ in order to guarantee a correct sampler:
 - ▶ guaranteed by the Fokker-Planck Equation
- 3 This is called stochastic gradient Langevin dynamics (SGLD).

⁷In the following, we will directly use $\mathcal{N}(\mathbf{0}, \mathbf{I})$ to represent a normal random variable with zero-mean and covariance matrix \mathbf{I} .

Random noise to the rescue

Gaussian noise with a variance given by th

- 1 Need to make the algorithm explore the parameter space:

- ▶ adding random Gaussian noise to the update⁷

$$\begin{aligned}\theta_{l+1} &= \theta_l - h_{l+1} \nabla_{\theta} \tilde{U}(\theta_l) + \sqrt{2h_{l+1}} \zeta_{l+1} \\ \zeta_{l+1} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I})\end{aligned}$$

- 2 The magnitude of the Gaussian needs to be $\sqrt{2h_{l+1}}$ in order to guarantee a correct sampler:

- ▶ guaranteed by the Fokker-Planck Equation

- 3 This is called stochastic gradient Langevin dynamics (SGLD).

Basically: SGD + properly scaled Gaussian noise.

⁷ In the following, we will directly use $\mathcal{N}(\mathbf{0}, \mathbf{I})$ to represent a normal random variable with zero-mean and covariance matrix \mathbf{I} .

SGLD in algorithm

Input: Parameters $\{h_l\}$

Output: Approximate samples $\{\theta_l\}$

Initialize $\theta_0 \in \mathbb{R}^n$

for $l = 1, 2, \dots$ **do**

 Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the l -th minibatch

$\theta_l = \theta_{l-1} - \nabla \tilde{U}(\theta_{l-1}) h_l + \sqrt{2h_l} \mathcal{N}(\mathbf{0}, \mathbf{I})$

end

Return $\{\theta_l\}$

Algorithm 1: Stochastic Gradient Langevin Dynamics

Example⁸

① A simple Gaussian mixture:

$$\theta_1 \sim \mathcal{N}(0, 10), \quad \theta_2 \sim \mathcal{N}(0, 1)$$

$$x_i \sim \frac{1}{2}\mathcal{N}(\theta_1, 2) + \frac{1}{2}\mathcal{N}(\theta_1 + \theta_2, 2), \quad i = 1, \dots, 100$$

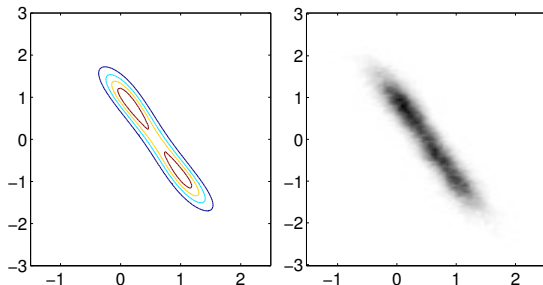


Figure: Left: true posterior; Right: sample-based estimation.

⁸M. Welling and Y. W. Teh. "Bayesian learning via stochastic gradient Langevin dynamics". In: *ICML*. 2011.

1 Markov Chain Monte Carlo Methods

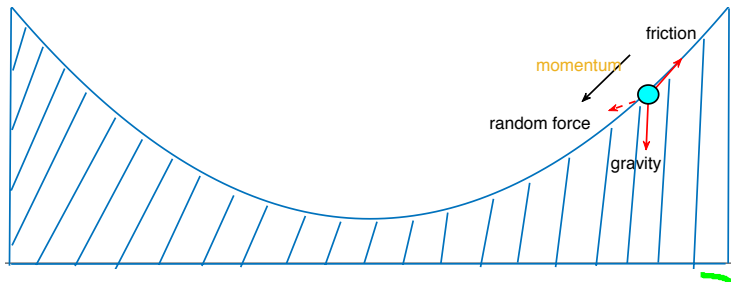
- Monte Carlo methods
- Markov chain Monte Carlo

2 Stochastic Gradient Markov Chain Monte Carlo Methods

- Introduction
- Stochastic gradient Langevin dynamics
- **Stochastic gradient Hamiltonian Monte Carlo**
- Application in Latent Dirichlet allocation

SGHMC

- 1 A 2nd-order method: stochastic gradients applied on some auxiliary parameters (momentum).
- 2 SGLD is slow when parameter space exhibits uneven curvatures.
- 3 Use the momentum idea to improve SGLD:
 - ▶ a generalization of the HMC, in that the ball is rolling on a friction surface
 - ▶ the ball follows the momentum instead of gradients, which is a summarization of historical gradients, thus could jump out local modes easier and move faster
 - ▶ needs a balance between these extra forces



Adding a friction term

- 1 Without a friction term, the random Gaussian noise would drive the ball too far away from their stationary distribution.
- 2 After adding a friction term:

$$\theta_l = \theta_{l-1} + \mathbf{v}_{l-1} h_l$$

$$\mathbf{v}_l = \mathbf{v}_{l-1} - \nabla_{\theta} \tilde{U}(\theta_l) h_l - A \mathbf{v}_{l-1} h_l + \sqrt{2A h_l} \mathcal{N}(\mathbf{0}, \mathbf{I}) ,$$

where $A > 0$ is a constant⁹, controlling the magnitude of the friction.

- 3 The friction term penalize the momentum:
 - ▶ the more momentum, the more friction it has, thus slowing down the ball

⁹In the original SGHMC paper, A is decomposed into a known variance of injected noise and an unknown variance of stochastic gradients.

Adding a friction term

- 1 Without a friction term, the random Gaussian noise would drive the ball too far away from their stationary distribution.
- 2 After adding a friction term:

$$\theta_l = \theta_{l-1} + \mathbf{v}_{l-1} h_l$$

$$\mathbf{v}_l = \mathbf{v}_{l-1} - \nabla_{\theta} \tilde{U}(\theta_l) h_l - A \mathbf{v}_{l-1} h_l + \sqrt{2A h_l} \mathcal{N}(\mathbf{0}, \mathbf{I}) ,$$

where $A > 0$ is a constant⁹, controlling the magnitude of the friction.

- 3 The friction term penalize the momentum:
 - ▶ the more momentum, the more friction it has, thus slowing down the ball

⁹In the original SGHMC paper, A is decomposed into a known variance of injected noise and an unknown variance of stochastic gradients.

Adding a friction term

- 1 Without a friction term, the random Gaussian noise would drive the ball too far away from their stationary distribution.
- 2 After adding a friction term:

$$\theta_l = \theta_{l-1} + \mathbf{v}_{l-1} h_l$$

$$\mathbf{v}_l = \mathbf{v}_{l-1} - \nabla_{\theta} \tilde{U}(\theta_l) h_l - A \mathbf{v}_{l-1} h_l + \sqrt{2A h_l} \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

where $A > 0$ is a constant⁹, controlling the magnitude of the friction.

- 3 The friction term penalize the momentum:
 - ▶ the more momentum, the more friction it has, thus slowing down the ball

friction*

⁹In the original SGHMC paper, A is decomposed into a known variance of injected noise and an unknown variance of stochastic gradients.

SGHMC in algorithm

Input: Parameters $A, \{h_l\}$

Output: Approximate samples $\{\theta_l\}$

Initialize $\theta_0 \in \mathbb{R}^n$

for $l = 1, 2, \dots$ **do**

Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the l -th minibatch

$\theta_l = \theta_{l-1} + \mathbf{v}_{l-1} h_l$

$\mathbf{v}_l = \mathbf{v}_{l-1} - \nabla \tilde{U}(\theta_l) h_l - A \mathbf{v}_{l-1} h_l + \sqrt{2A h_l} \mathcal{N}(\mathbf{0}, \mathbf{I})$

end

Return $\{\theta_l\}$

Algorithm 2: Stochastic Gradient Hamiltonian Monte Carlo

Reparametrize SGHMC

for $l = 1, 2, \dots$ **do**

 Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the
 l -th minibatch

$$\theta_l = \theta_{l-1} + \mathbf{v}_{l-1} h_l$$

$$\mathbf{v}_l = \mathbf{v}_{l-1} - \nabla \tilde{U}(\theta_l) h_l -$$

$$A \mathbf{v}_{l-1} h_l + \sqrt{2Ah_l} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

end

- Reparametrization: $\epsilon = h^2$, $m = Ah$, $\mathbf{p} = \mathbf{v} h$

Reparametrize SGHMC

for $l = 1, 2, \dots$ **do**

Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the
 l -th minibatch

$$\theta_l = \theta_{l-1} + \mathbf{v}_{l-1} h_l$$

$$\mathbf{v}_l = \mathbf{v}_{l-1} - \nabla \tilde{U}(\theta_l) h_l - A \mathbf{v}_{l-1} h_l + \sqrt{2Ah_l} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

end

for $l = 1, 2, \dots$ **do**

Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the
 l -th minibatch

$$\theta_l = \theta_{l-1} + \mathbf{p}_{l-1}$$

$$\mathbf{p}_l = (1 - m) \mathbf{p}_{l-1} - \nabla \tilde{U}(\theta_l) \epsilon_l + \sqrt{2m\epsilon_l} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

end

- Reparametrization: $\epsilon = h^2$, $m = Ah$, $\mathbf{p} = \mathbf{v} h$

Reparametrize SGHMC

for $l = 1, 2, \dots$ **do**

Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the l -th minibatch

$$\theta_l = \theta_{l-1} + \mathbf{v}_{l-1} h_l$$

$$\mathbf{v}_l = \mathbf{v}_{l-1} - \nabla \tilde{U}(\theta_l) h_l -$$

$$A \mathbf{v}_{l-1} h_l + \sqrt{2Ah_l} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

end

for $l = 1, 2, \dots$ **do**

Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the l -th minibatch

$$\theta_l = \theta_{l-1} + \mathbf{p}_{l-1}$$

$$\mathbf{p}_l = (1 - m) \mathbf{p}_{l-1} - \nabla \tilde{U}(\theta_l) \epsilon_l + \sqrt{2m\epsilon_l} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

end

- Reparametrization: $\epsilon = h^2$, $m = Ah$, $\mathbf{v} = \mathbf{p} h$
- ϵ_l : learning rate; m : momentum weight

SGD vs. SGLD

Just added properly scaled

$$\nabla_{\theta} \tilde{U}(\theta_{l-1}) \triangleq -\frac{N}{n} \sum_{i=1}^n \nabla_{\theta} \log p(\mathbf{x}_{\pi_i} | \theta_{l-1}) - \nabla_{\theta} \log p(\theta_{l-1}),$$

SGD:

for $l = 1, 2, \dots$ **do**

Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the
 l -th minibatch

$$\theta_l = \theta_{l-1} - \nabla \tilde{U}(\theta_l) \epsilon_l$$

end

SGLD:

for $l = 1, 2, \dots$ **do**

Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the
 l -th minibatch

$$\theta_l = \theta_{l-1} - \nabla \tilde{U}(\theta_l) \epsilon_l + \delta_l$$

$$\delta_l \sim \mathcal{N}(\mathbf{0}, 2\epsilon_l \mathbf{I})$$

end

SGD with Momentum (SGD-M) vs. SGHMC

Just added properly scaled

$$\nabla_{\theta} \tilde{U}(\theta_{l-1}) \triangleq -\frac{N}{n} \sum_{i=1}^n \nabla_{\theta} \log p(\mathbf{x}_{\pi_i} | \theta_{l-1}) - \nabla_{\theta} \log p(\theta_{l-1}),$$

SGD-M:

for $l = 1, 2, \dots$ **do**

Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the
 l -th minibatch

$$\theta_l = \theta_{l-1} + \mathbf{p}_{l-1}$$

$$\mathbf{p}_l = (1 - m) \mathbf{p}_{l-1} - \nabla \tilde{U}(\theta_l) \epsilon_l$$

end

SGHMC:

for $l = 1, 2, \dots$ **do**

Evaluate $\nabla_{\theta} \tilde{U}(\theta_{l-1})$ from the
 l -th minibatch

$$\theta_l = \theta_{l-1} + \mathbf{p}_{l-1}$$

$$\mathbf{p}_l = (1 - m) \mathbf{p}_{l-1} - \nabla \tilde{U}(\theta_l) \epsilon_l + \delta_l$$

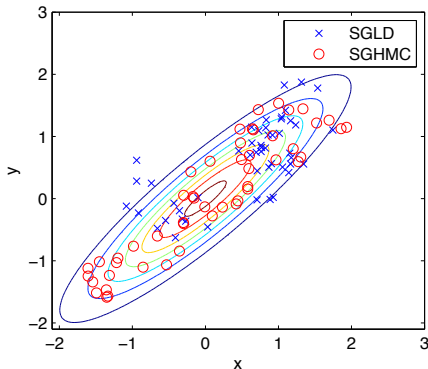
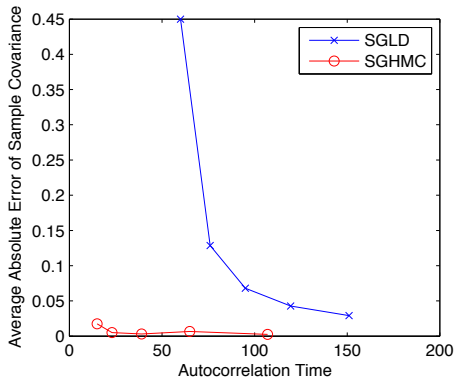
$$\delta_l \sim \mathcal{N}(\mathbf{0}, 2m\epsilon_l \mathbf{I})$$

end

Example¹⁰

1 Sample from a 2D Gaussian distribution:

► $U(\theta) = \frac{1}{2} \theta^T \Sigma^{-1} \theta$



¹⁰T. Chen, E. B. Fox, and C. Guestrin. "Stochastic Gradient Hamiltonian Monte Carlo". In: *ICML*. 2014.

Recap

1 For SG-MCMC methods, in each iteration:

- ▶ calculate the stochastic gradient based on the current parameter sample
- ▶ generate the next sample by moving the current sample (probably in an extended space) along the direction of the stochastic gradient, plus a suitable random Gaussian noise
- ▶ no need for accept-reject
- ▶ guaranteed to converge close to the true posterior in some sense

1 Markov Chain Monte Carlo Methods

- Monte Carlo methods
- Markov chain Monte Carlo

2 Stochastic Gradient Markov Chain Monte Carlo Methods

- Introduction
- Stochastic gradient Langevin dynamics
- Stochastic gradient Hamiltonian Monte Carlo
- Application in Latent Dirichlet allocation

Latent Dirichlet allocation

- 1 For each topic k , draw the **topic-word distribution**:

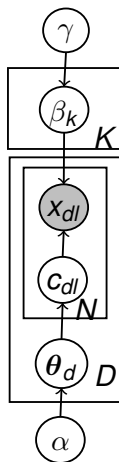
$$\beta_k \sim \text{Dir}(\gamma)$$

- 2 For each document d , draw its **topic distribution**: $\theta_d \sim \text{Dir}(\alpha)$
- ▶ For each word l , draw its **topic indicator**:

$$c_{dl} \sim \text{Discrete}(\theta_d)$$

- ▶ Draw the observed **word**:

$$x_{dl} \sim \text{Discrete}(\beta_{c_{dl}})$$



Latent Dirichlet allocation

- 1 Let $\beta \triangleq (\beta_k)_{k=1}^K$, $\theta \triangleq (\theta_d)_{d=1}^D$, $\mathbf{C} \triangleq (c_{dl})_{d,l=1}^{D,n_d}$, $\mathbf{X} \triangleq (x_{dl})_{d,l=1}^{D,n_d}$, the posterior distribution

$$p(\beta, \theta, \mathbf{C} | \mathbf{X}) \propto \left[\prod_{k=1}^K p(\beta_k | \gamma) \right] \left[\prod_{d=1}^D p(\theta_d | \alpha) \prod_{l=1}^{n_d} p(c_{dl} | \theta_d) p(x_{dl} | \beta, c_{dl}) \right]$$

- 2 From previous lectures:

$$p(c_{dl} | \theta_d) = \prod_{k=1}^K (\theta_{dk})^{1(c_{dl}=k)}$$

$$p(x_{dl} | \theta, c_{dl}) = \prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{1(x_{dl}=v)1(c_{dl}=k)}$$

- 3 Together with the fact:

$$\int_{\theta \in \Delta_{K-1}} \prod_{k=1}^K \theta_k^{\alpha_k - 1} d\theta_k = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}$$

Latent Dirichlet allocation

- 1 Let $\beta \triangleq (\beta_k)_{k=1}^K$, $\theta \triangleq (\theta_d)_{d=1}^D$, $\mathbf{C} \triangleq (c_{dl})_{d,l=1}^{D,n_d}$, $\mathbf{X} \triangleq (x_{dl})_{d,l=1}^{D,n_d}$, the posterior distribution

$$p(\beta, \theta, \mathbf{C} | \mathbf{X}) \propto \left[\prod_{k=1}^K p(\beta_k | \gamma) \right] \left[\prod_{d=1}^D p(\theta_d | \alpha) \prod_{l=1}^{n_d} p(c_{dl} | \theta_d) p(x_{dl} | \beta, c_{dl}) \right]$$

- 2 From previous lectures:

$$p(c_{dl} | \theta_d) = \prod_{k=1}^K (\theta_{dk})^{1(c_{dl}=k)}$$

$$p(x_{dl} | \theta, c_{dl}) = \prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{1(x_{dl}=v)1(c_{dl}=k)}$$

- 3 Together with the fact:

$$\int_{\theta \in \Delta_{K-1}} \prod_{k=1}^K \theta_k^{\alpha_k - 1} d\theta_k = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}$$

Latent Dirichlet allocation

- ① Integrate out the local parameters: **topic distributions** θ for each document, it results in the following semi-collapsed distribution:
 $p(\mathbf{X}, \mathbf{C}, \beta | \alpha, \gamma) =$

$$\prod_{d=1}^D \frac{\Gamma(K\alpha)}{\Gamma(K\alpha + n_{d\cdot})} \prod_{k=1}^K \frac{\Gamma(\alpha + n_{dk\cdot})}{\Gamma(\alpha)} \prod_{k=1}^K \frac{\Gamma(V\gamma)}{\Gamma(\gamma)^V} \prod_{v=1}^V \beta_{kv}^{\gamma + n_{\cdot kv} - 1},$$

where $n_{dkw} \triangleq \sum_{l=1}^{n_d} 1(c_{dl} = k)1(x_{dl} = w)$ is #word w in doc d with topic k ; \cdot means marginal sum, e.g. $n_{\cdot kw} \triangleq \sum_{d=1}^D n_{dkw}$.

- ② SG-MCMC requires parameter spaces unconstrained:
▶ reparameterization: $\beta_{kv} = \lambda_{kv} / \sum_{v'} \lambda_{kv'}$, with the following prior:

$$\lambda_{kv} \sim \text{Ga}(\lambda_{kv}; \gamma, 1)$$

$$\prod_{k=1}^K \frac{\Gamma(V\gamma)}{\Gamma(\gamma)^V} \prod_{v=1}^V \beta_{kv}^{\gamma + n_{\cdot kv} - 1} \Rightarrow \prod_{k=1}^K \prod_{v=1}^V \text{Ga}(\lambda_{kv}; \gamma, 1) \prod_{v=1}^V (\lambda_{kv} / \sum_{v'} \lambda_{kv'})^{n_{\cdot kw}}$$

Latent Dirichlet allocation

- ① Integrate out the local parameters: **topic distributions** θ for each document, it results in the following semi-collapsed distribution:
 $p(\mathbf{X}, \mathbf{C}, \beta | \alpha, \gamma) =$

$$\prod_{d=1}^D \frac{\Gamma(K\alpha)}{\Gamma(K\alpha + n_{d\cdot})} \prod_{k=1}^K \frac{\Gamma(\alpha + n_{dk\cdot})}{\Gamma(\alpha)} \prod_{k=1}^K \frac{\Gamma(V\gamma)}{\Gamma(\gamma)^V} \prod_{v=1}^V \beta_{kv}^{\gamma + n_{\cdot kv} - 1},$$

where $n_{dkw} \triangleq \sum_{l=1}^{n_d} 1(c_{dl} = k)1(x_{dl} = w)$ is #word w in doc d with topic k ; \cdot means marginal sum, e.g. $n_{\cdot kw} \triangleq \sum_{d=1}^D n_{dkw}$.

- ② SG-MCMC requires parameter spaces unconstrained:
▶ reparameterization: $\beta_{kv} = \lambda_{kv} / \sum_{v'} \lambda_{kv'}$, with the following prior:

$$\lambda_{kv} \sim \text{Ga}(\lambda_{kv}; \gamma, 1)$$

$$\prod_{k=1}^K \frac{\Gamma(V\gamma)}{\Gamma(\gamma)^V} \prod_{v=1}^V \beta_{kv}^{\gamma + n_{\cdot kv} - 1} \implies \prod_{k=1}^K \prod_{v=1}^V \text{Ga}(\lambda_{kv}; \gamma, 1) \prod_{v=1}^V (\lambda_{kv} / \sum_{v'} \lambda_{kv'})^{n_{\cdot kw}}$$

Latent Dirichlet allocation

- 1 Still need to integrate out the local parameter \mathbf{C} :

$$p(\mathbf{X}, \lambda | \alpha, \gamma) = \mathbb{E}_{\mathbf{C}} [p(\mathbf{X}, \mathbf{C}, \beta | \alpha, \gamma)] = \mathbb{E}_{\mathbf{C}} \left[\prod_{d=1}^D \frac{\Gamma(K\alpha)}{\Gamma(K\alpha + n_{d..})} \prod_{k=1}^K \frac{\Gamma(\alpha + n_{dk.})}{\Gamma(\alpha)} \prod_{v=1}^V \text{Ga}(\lambda_{kv}; \gamma, 1) \left(\frac{\lambda_{kv}}{\sum_{v'} \lambda_{kv'}} \right)^{n_{.kv}} \right]$$

- 2 The stochastic gradient with a minibatch documents \bar{D} of size $|\bar{D}| \ll D$ is:

$$\frac{\partial \log \tilde{p}(\lambda | \alpha, \gamma, \mathbf{X})}{\partial \lambda_{kw}} = \frac{\gamma - 1}{\lambda_{kw}} - 1 + \frac{D}{|\bar{D}|} \sum_{d \in \bar{D}} \mathbb{E}_{\mathbf{C}_d | \mathbf{x}_d, \lambda, \alpha} \left[\frac{n_{dkw}}{\lambda_{kw}} - \frac{n_{dk.}}{\lambda_{k.}} \right]$$

- 3 SGLD update:

$$\lambda_{kw}^{t+1} = \lambda_{kw}^t + \frac{\partial \log \tilde{p}(\lambda | \alpha, \gamma, \mathbf{X})}{\partial \lambda_{kw}} h_{t+1} + \sqrt{2h_{t+1}} N(0, \mathbf{I})$$

Latent Dirichlet allocation

- 1 Still need to integrate out the local parameter \mathbf{C} :

$$p(\mathbf{X}, \lambda | \alpha, \gamma) = \mathbb{E}_{\mathbf{C}} [p(\mathbf{X}, \mathbf{C}, \beta | \alpha, \gamma)] = \mathbb{E}_{\mathbf{C}} \left[\prod_{d=1}^D \frac{\Gamma(K\alpha)}{\Gamma(K\alpha + n_{d..})} \prod_{k=1}^K \frac{\Gamma(\alpha + n_{dk.})}{\Gamma(\alpha)} \prod_{v=1}^V \text{Ga}(\lambda_{kv}; \gamma, 1) \left(\frac{\lambda_{kv}}{\sum_{v'} \lambda_{kv'}} \right)^{n_{.kv}} \right]$$

- 2 The stochastic gradient with a minibatch documents \bar{D} of size $|\bar{D}| \ll D$ is:

$$\frac{\partial \log \tilde{p}(\lambda | \alpha, \gamma, \mathbf{X})}{\partial \lambda_{kw}} = \frac{\gamma - 1}{\lambda_{kw}} - 1 + \frac{D}{|\bar{D}|} \sum_{d \in \bar{D}} \mathbb{E}_{\mathbf{C}_d | \mathbf{x}_d, \lambda, \alpha} \left[\frac{n_{dkw}}{\lambda_{kw}} - \frac{n_{dk.}}{\lambda_{k.}} \right]$$

- 3 SGLD update:

$$\lambda_{kw}^{t+1} = \lambda_{kw}^t + \frac{\partial \log \tilde{p}(\lambda | \alpha, \gamma, \mathbf{X})}{\partial \lambda_{kw}} h_{t+1} + \sqrt{2h_{t+1}} N(0, \mathbf{I})$$

Latent Dirichlet allocation

- 1 Still need to integrate out the local parameter \mathbf{C} :

$$p(\mathbf{X}, \lambda | \alpha, \gamma) = \mathbb{E}_{\mathbf{C}} [p(\mathbf{X}, \mathbf{C}, \beta | \alpha, \gamma)] = \mathbb{E}_{\mathbf{C}} \left[\prod_{d=1}^D \frac{\Gamma(K\alpha)}{\Gamma(K\alpha + n_{d..})} \prod_{k=1}^K \frac{\Gamma(\alpha + n_{dk.})}{\Gamma(\alpha)} \prod_{v=1}^V \text{Ga}(\lambda_{kv}; \gamma, 1) \left(\frac{\lambda_{kv}}{\sum_{v'} \lambda_{kv'}} \right)^{n_{.kv}} \right]$$

- 2 The stochastic gradient with a minibatch documents \bar{D} of size $|\bar{D}| \ll D$ is:

$$\frac{\partial \log \tilde{p}(\lambda | \alpha, \gamma, \mathbf{X})}{\partial \lambda_{kw}} = \frac{\gamma - 1}{\lambda_{kw}} - 1 + \frac{D}{|\bar{D}|} \sum_{d \in \bar{D}} \mathbb{E}_{\mathbf{C}_d | \mathbf{x}_d, \lambda, \alpha} \left[\frac{n_{dkw}}{\lambda_{kw}} - \frac{n_{dk.}}{\lambda_{k.}} \right]$$

- 3 SGLD update:

$$\lambda_{kw}^{t+1} = \lambda_{kw}^t + \frac{\partial \log \tilde{p}(\lambda | \alpha, \gamma, \mathbf{X})}{\partial \lambda_{kw}} h_{t+1} + \sqrt{2h_{t+1}} N(0, \mathbf{I})$$

Latent Dirichlet allocation

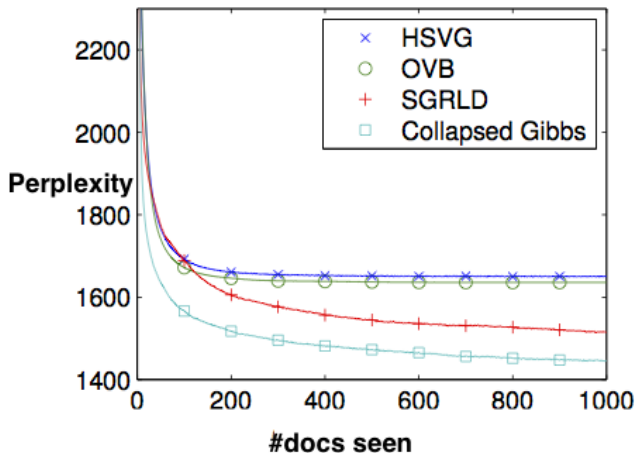
- ❶ LDA with the above SGLD update would not work well in practice because of the high dimensionality of model parameters.
- ❷ To make it work, Riemannian geometry information (2nd-order information) need to bring in SGLD:
 - ▶ leading to Stochastic Gradient Riemannian Langevin Dynamics (SGRLD) for LDA¹¹
 - ▶ it considers parameter geometry so that step sizes for each dimension of the parameter are adaptive

¹¹S. Patterson and Y. W. Teh. "Stochastic Gradient Riemannian Langevin Dynamics on the Probability Simplex". In: *NIPS*. 2013.

Experiments: SGRLD for LDA¹²

1 NIPS dataset:

- ▶ the collection of NIPS papers from 1988-2003, with 2483 documents, 50 topics

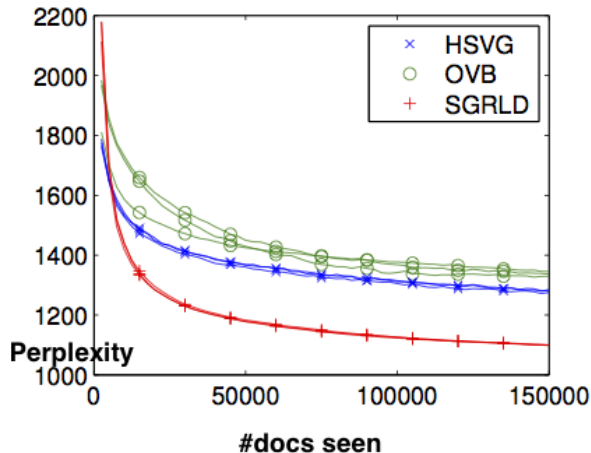


¹²S. Patterson and Y. W. Teh. "Stochastic Gradient Riemannian Langevin Dynamics on the Probability Simplex". In: *NIPS*.

Experiments: SGRLD for LDA¹³

1 Wikipedia dataset:

- ▶ a set of articles downloaded at random from Wikipedia, with 150,000 documents



¹³S. Patterson and Y. W. Teh. "Stochastic Gradient Riemannian Langevin Dynamics on the Probability Simplex". In: *NIPS*.

Conclusion

- 1 I have introduced:
 - ▶ basic concepts in MCMC
 - ▶ basic ideas in SG-MCMC, two SG-MCMC algorithms, and application in LDA
- 2 Topics not covered:
 - ▶ a general review of SG-MCMC algorithms
 - ▶ theory related to stochastic differential equations and Itô diffusions
 - ▶ convergence theory
 - ▶ various applications in deep learning, including SG-MCMC for learning weight uncertainty and SG-MCMC for deep generative models
 - ▶ interested readers should refer to related references

Thank You