

# Markov Chain Monte Carlo (MCMC)

*Dmitry Kropotov*  
*Lomonosov Moscow State University*



# Contents

Ising model

Gibbs sampling

Metropolis-Hastings sampling

Hamiltonian Monte Carlo

## Recap: Variational inference

Probabilistic model:  $p(\mathbf{x}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\theta})$ .

Goal: find approximant  $q(\boldsymbol{\theta} | \boldsymbol{\lambda})$  for posterior  $p(\boldsymbol{\theta} | \mathbf{x})$ .

Method:  $\log p(\mathbf{x}) \geq \text{ELBO}(\boldsymbol{\lambda}) = \sum_{i=1}^n \mathbb{E}_{q(\boldsymbol{\theta} | \boldsymbol{\lambda})} \log p(\mathbf{x}_i | \boldsymbol{\theta}) - \mathbb{E}_{q(\boldsymbol{\theta} | \boldsymbol{\lambda})} \log \frac{q(\boldsymbol{\theta} | \boldsymbol{\lambda})}{p(\boldsymbol{\theta})} \rightarrow \max_{\boldsymbol{\lambda}}.$

## Recap: Variational inference

Probabilistic model:  $p(\mathbf{x}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\theta})$ .

Goal: find approximant  $q(\boldsymbol{\theta} | \boldsymbol{\lambda})$  for posterior  $p(\boldsymbol{\theta} | \mathbf{x})$ .

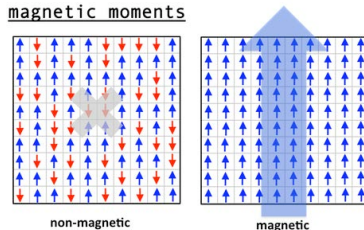
Method:  $\log p(\mathbf{x}) \geq \text{ELBO}(\boldsymbol{\lambda}) = \sum_{i=1}^n \mathbb{E}_{q(\boldsymbol{\theta} | \boldsymbol{\lambda})} \log p(\mathbf{x}_i | \boldsymbol{\theta}) - \mathbb{E}_{q(\boldsymbol{\theta} | \boldsymbol{\lambda})} \log \frac{q(\boldsymbol{\theta} | \boldsymbol{\lambda})}{p(\boldsymbol{\theta})} \rightarrow \max_{\boldsymbol{\lambda}}$ .

Problem:  $q(\boldsymbol{\theta} | \boldsymbol{\lambda})$  could be quite poor approximant for the true posterior!

# Ising model

**Ising model** – a model from statistical physics describing magnetic properties of ferromagnetic substances depending on temperature. It reflects **phase transition effect** – a drastic magnetization changing within narrow temperature interval.

- ▶ Consider 2-dimensional atomic lattice.
- ▶  $x_i \in \{-1, +1\}$  – magnetic moment for one atom.
- ▶ Energy  $E(\mathbf{x}) = -\sum_{(i,j) \in \mathcal{E}} x_i x_j - \sum_{i=1}^n h_i x_i$ .
- ▶ Boltzmann probability distribution  $p(\mathbf{x}) = \frac{1}{Z} \exp\left(-\frac{1}{T} E(\mathbf{x})\right)$ , where  $T$  – temperature,  $Z = \sum_{\mathbf{x}} \exp\left(-\frac{1}{T} E(\mathbf{x})\right)$  – normalizing constant.



# Ising model

Ising model:

$$E(\mathbf{x}) = - \sum_{(i,j) \in \mathcal{E}} x_i x_j - \sum_{i=1}^n h_i x_i,$$
$$p(\mathbf{x}) = \frac{1}{Z} \exp \left( -\frac{1}{T} E(\mathbf{x}) \right).$$

Would like to estimate for each temperature:

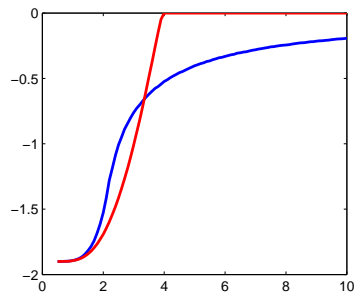
- ▶ Mean energy  $\mathbb{E}_{p(\mathbf{x})} E(\mathbf{x})$ ;
- ▶ Energy variance  $\mathbb{D}_{p(\mathbf{x})} E(\mathbf{x})$ ;
- ▶ Mean magnetization  $\sqrt{\mathbb{E}_{p(\mathbf{x})} \mu^2}$ , where  $\mu = \frac{1}{n} \sum_{i=1}^n x_i$ .

Key challenge: normalizing constant  $Z$  is a sum over  $2^n$  components.

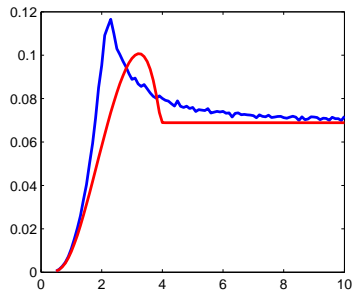
# Variational inference for Ising model

Mean field approximation:  $p(\mathbf{x}) = \frac{1}{Z} \exp\left(-\frac{1}{T} E(\mathbf{x})\right) \approx q(\mathbf{x}) = \prod_{i=1}^n q_i(x_i)$ .

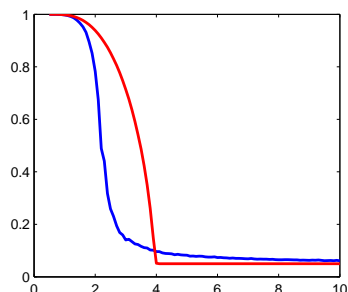
Recalculation:  $q_i(x_i = 1) = 1/(1 + \exp(-(2/T)(h_i + \sum_{j:(i,j) \in \mathcal{E}} \mathbb{E}_{q_j} x_j)))$ .



$\mathbb{E}_{q(\mathbf{x})} E(\mathbf{x})$



$\mathbb{D}_{q(\mathbf{x})} E(\mathbf{x})$



$\sqrt{\mathbb{E}_{q(\mathbf{x})} \mu^2}$

# Variational inference for Ising model

## Variational approximations:



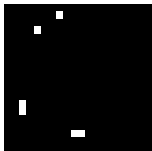
Low temperature

Critical temperature



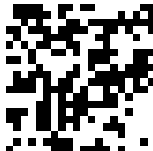
High temperature

## Samples from true posterior:



Low temperature

Critical temperature



High temperature



# MCMC

Probabilistic model:  $p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x})$ ,  $Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$ .

MCMC methods construct samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$  from  $p(\mathbf{x})$  using only  $\tilde{p}(\mathbf{x})$ .

Usage:  $\mathbb{E}_{p(\mathbf{x})} f(\mathbf{x}) = \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \approx \frac{1}{m} \sum_{i=1}^m f(\mathbf{x}_i)$ .

# MCMC

MCMC generate samples  $\mathbf{x}_1, \dots, \mathbf{x}_m$  using Markov Chain with some transition probability  $q(\mathbf{x}|\mathbf{y})$  in the following way:

- ▶ Generate  $\mathbf{x}_1$  from some initial distribution  $p_0(\mathbf{x})$ ;
- ▶ Generate  $\mathbf{x}_2$  from  $q(\mathbf{x}|\mathbf{x}_1)$ ;
- ▶ Generate  $\mathbf{x}_3$  from  $q(\mathbf{x}|\mathbf{x}_2)$ ;
- ▶ etc.

Samples  $\mathbf{x}_1, \dots, \mathbf{x}_m$  are not independent, but still can be used for estimation  $\mathbb{E}_{p(\mathbf{x})}f(\mathbf{x})$ .

# Markov chain properties

Consider a Markov chain with transition probability  $q(\mathbf{x}|\mathbf{y})$ .

Probability distribution  $\pi(\mathbf{x})$  is called **invariant** under given Markov chain iff

$$\int q(\mathbf{x}|\mathbf{y})\pi(\mathbf{y})d\mathbf{y} = \pi(\mathbf{x}).$$

Consider initial distribution  $p_0(\mathbf{x})$  and denote  $p_i(\mathbf{x})$  – distribution of points after  $i$  steps of given Markov chain. Then the Markov chain is called **ergodic** iff

$$p_i(\mathbf{x}) \rightarrow \pi(\mathbf{x}), \quad i \rightarrow \infty, \quad \forall p_0(\mathbf{x}),$$

where  $\pi(\mathbf{x})$  is invariant distribution for this chain.

For generating samples from distribution  $p(\mathbf{x})$  using MCMC, the corresponding Markov chain should be ergodic with invariant distribution  $p(\mathbf{x})$ .

# Gibbs sampling

Goal: generate samples from  $p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x})$ , where  $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k]$ .

Suppose  $\mathbf{x} \sim p(\mathbf{x})$ . Then the next point  $\mathbf{x}^{new}$  is generated as follows:

- ▶  $\mathbf{x}_1^{new} \sim p(\mathbf{x}_1 | \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_k);$
- ▶  $\mathbf{x}_2^{new} \sim p(\mathbf{x}_2 | \mathbf{x}_1^{new}, \mathbf{x}_3, \dots, \mathbf{x}_k);$
- ▶  $\dots;$
- ▶  $\mathbf{x}_k^{new} \sim p(\mathbf{x}_k | \mathbf{x}_1^{new}, \mathbf{x}_3^{new}, \dots, \mathbf{x}_{k-1}^{new});$
- ▶  $\mathbf{x}^{new} = [\mathbf{x}_1^{new}, \mathbf{x}_2^{new}, \dots, \mathbf{x}_k^{new}].$

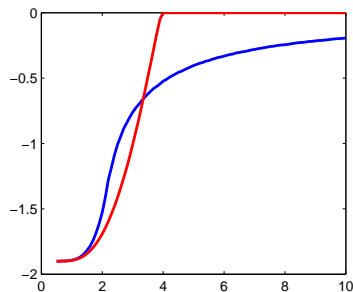
It is easy to show that  $p(\mathbf{x})$  is invariant under such transition probability. If all conditionals  $p(\mathbf{x}_i | \mathbf{x}_{\setminus i}) > 0$ , then the corresponding Markov chain would be ergodic.

# Gibbs sampling for Ising model

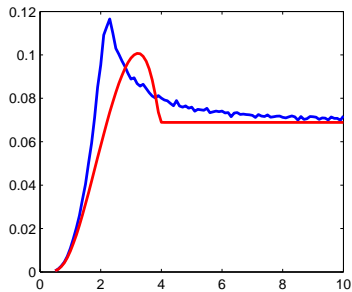
$$\text{Ising model: } p(\mathbf{x}) = \frac{1}{Z} \exp \left( \frac{1}{T} \left( \sum_{i=1}^n h_i x_i + \sum_{(i,j) \in \mathcal{E}} x_i x_j \right) \right).$$

For Gibbs sampling we need to calculate all conditionals:

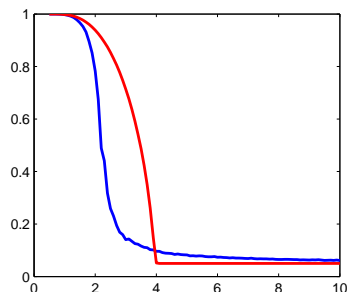
$$p(x_i = 1 | \mathbf{x}_{\setminus i}) = 1 / (1 + \exp(-(2/T)(h_i + \sum_{j:(i,j) \in \mathcal{E}} x_j))).$$



$\mathbb{E}_{q(\mathbf{x})} E(\mathbf{x})$



$\mathbb{D}_{q(\mathbf{x})} E(\mathbf{x})$



$\sqrt{\mathbb{E}_{q(\mathbf{x})} \mu^2}$

# Metropolis-Hastings sampling

Goal: generate samples from  $p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x})$ .

Suppose we have some proposal distribution  $r(\mathbf{x}|\mathbf{y})$ , from which we can generate samples. Sampling scheme:

- ▶ Generate  $\mathbf{x}_{trial} \sim r(\mathbf{x}|\mathbf{x}_{old})$ ;
- ▶ Calculate  $A = \min \left( 1, \frac{\tilde{p}(\mathbf{x}_{trial})r(\mathbf{x}_{old}|\mathbf{x}_{trial})}{\tilde{p}(\mathbf{x}_{old})r(\mathbf{x}_{trial}|\mathbf{x}_{old})} \right)$ ;
- ▶  $\mathbf{x}_{new} = \mathbf{x}_{trial}$  with probability  $A$  and  $\mathbf{x}_{new} = \mathbf{x}_{old}$  with probability  $1 - A$ .

# Metropolis sampling

Suppose that proposal is symmetric, i.e.  $r(\mathbf{x}|\mathbf{y}) = r(\mathbf{y}|\mathbf{x})$ . Then probability  $A$  is reduced to:

$$A = \min \left( 1, \frac{\tilde{p}(\mathbf{x}_{trial})}{\tilde{p}(\mathbf{x}_{old})} \right).$$

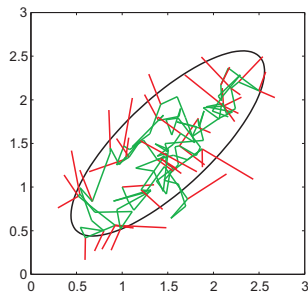
Sampling scheme:

- ▶ Generate  $\mathbf{x}_{trial} \sim r(\mathbf{x}|\mathbf{x}_{old})$ ;
- ▶ If  $\tilde{p}(\mathbf{x}_{trial}) \geq \tilde{p}(\mathbf{x}_{old})$ , then  $\mathbf{x}_{new} = \mathbf{x}_{trial}$ ;
- ▶ Else  $\mathbf{x}_{new} = \mathbf{x}_{trial}$  with probability  $\tilde{p}(\mathbf{x}_{trial})/\tilde{p}(\mathbf{x}_{old})$  and  $\mathbf{x}_{new} = \mathbf{x}_{old}$  otherwise.

# Metropolis sampling

Suppose we would like to sample from correlated 2-dimensional Gaussian distribution using MH sampling with the following symmetric proposal:

$$r(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{y}, \sigma^2 I).$$



For poor proposals MH sampling shows random walk behaviour!



# Hamiltonian equations

Suppose we have a particle moving in some potential field.

- ▶  $\mathbf{x}$  – particle coordinates;  $\mathbf{p}$  – particle momentum,  $\mathbf{p} = m \frac{d\mathbf{x}}{dt}$ ;
- ▶  $U(\mathbf{x})$  – potential energy;  $K(\mathbf{p}) = \mathbf{p}^T \mathbf{p} / m$  – kinetic energy;
- ▶  $H(\mathbf{x}, \mathbf{p}) = U(\mathbf{x}) + K(\mathbf{p})$  – Hamiltonian.

Hamiltonian equations:

$$\begin{aligned}\frac{dx_i}{dt} &= \frac{\partial H}{\partial p_i}, \\ \frac{dp_i}{dt} &= -\frac{\partial H}{\partial x_i}.\end{aligned}$$

# Hamiltonian Monte Carlo (HMC)

Goal: generate samples from  $p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x})$ .

Let's introduce the following probabilistic model:

$$p(\mathbf{x}, \mathbf{p}) = \frac{1}{Z} \exp(-H(\mathbf{x}, \mathbf{p})) = \frac{1}{Z} \exp(-U(\mathbf{x}) - \frac{1}{2} \mathbf{p}^T \mathbf{p})$$

$$U(\mathbf{x}) = -\log \tilde{p}(\mathbf{x}),$$

$\mathbf{p}$  – auxiliary variable.

HMC generation scheme:

- ▶ Generate  $\mathbf{p}$  from  $\frac{1}{Z} \exp(-(1/2) \mathbf{p}^T \mathbf{p})$ ;
- ▶ Solve Hamiltonian equation starting from  $(\mathbf{x}_{old}, \mathbf{p})$  and obtain  $(\mathbf{x}_{new}, \mathbf{p}_{new})$ ;
- ▶ Accept the new point  $\mathbf{x}_{new}$  with probability  $\min(1, \exp(-H(\mathbf{x}_{new}, \mathbf{p}_{new}) + H(\mathbf{x}_{old}, \mathbf{p}_{old})))$ .

# Solving Hamiltonian equations

## Hamiltonian equations:

$$\begin{aligned}\frac{dx_i}{dt} &= \frac{\partial H}{\partial p_i}, \\ \frac{dp_i}{dt} &= -\frac{\partial H}{\partial x_i}.\end{aligned}$$

## Euler's method:

$$\begin{aligned}p_i(t + \varepsilon) &= p_i(t) + \varepsilon \frac{dp_i(t)}{dt} = p_i(t) - \varepsilon \frac{\partial U(\mathbf{x}(t))}{\partial x_i}; \\ x_i(t + \varepsilon) &= x_i(t) + \varepsilon \frac{dx_i(t)}{dt} = x_i(t) + \frac{p_i(t)}{m_i}.\end{aligned}$$

# Solving Hamiltonian equations

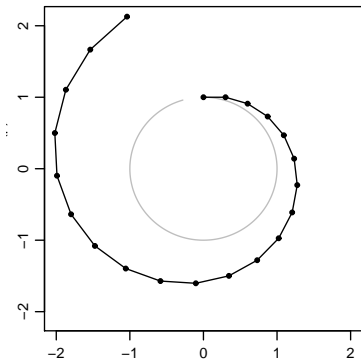
## Hamiltonian equations:

$$\begin{aligned}\frac{dx_i}{dt} &= \frac{\partial H}{\partial p_i}, \\ \frac{dp_i}{dt} &= -\frac{\partial H}{\partial x_i}.\end{aligned}$$

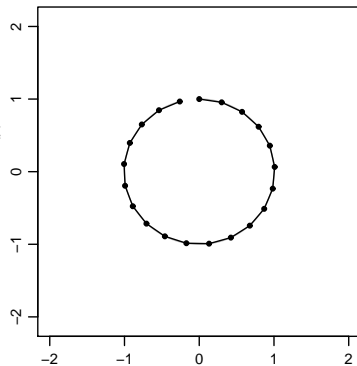
## Leapfrog method:

$$\begin{aligned}p_i(t + \varepsilon/2) &= p_i(t) - (\varepsilon/2) \frac{\partial U(\mathbf{x}(t))}{\partial x_i}; \\ x_i(t + \varepsilon) &= x_i(t) + \frac{p_i(t + \varepsilon/2)}{m_i}, \\ p_i(t + \varepsilon) &= p_i(t + \varepsilon/2) - (\varepsilon/2) \frac{\partial U(\mathbf{x}(t + \varepsilon))}{\partial x_i};\end{aligned}$$

# Solving Hamiltonian equations

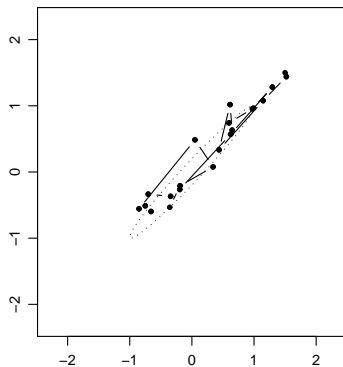


Euler's method

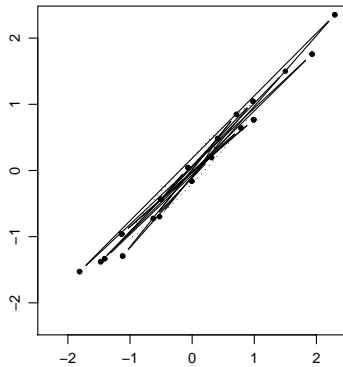


Leapfrog

# HMC Example



Random walk Metropolis



HMC

# Conclusion

- ▶ MCMC algorithms generate samples from desired distributions. Due to big amount of sampling they are usually much slower comparing to Variational Inference;
- ▶ In the limit MCMC algorithms give exact estimation of statistics  $\mathbb{E}_{p(\mathbf{x})}f(\mathbf{x})$ ;
- ▶ There is no explicit criterion function for tuning MCMC parameters;
- ▶ There exist scalable MCMC algorithms. See the next lecture.