

A Contour Stochastic Gradient Langevin Dynamics Algorithm for Simulations of Multi-modal Distributions

A scalable dynamic importance sampling algorithm

Wei Deng¹ Guang Lin¹ Faming Liang¹

October 23, 2020

¹Purdue University

Non-convex energy function leads to slow mixing

Given a non-convex energy function $U(\mathbf{x})$ of a density $\pi(\mathbf{x}) \propto e^{-\frac{U(\mathbf{x})}{\tau}}$, the standard sampling algorithm converges quite slowly.

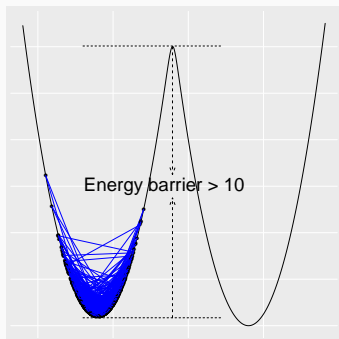


Figure 1: An example of a non-convex energy function $U(\mathbf{x})$.

Acceleration via importance sampling

To speed-up the simulation, we consider the importance sampling:

Acceleration via importance sampling

To speed-up the simulation, we consider the importance sampling:

$$\mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x}) \underbrace{\pi(\mathbf{x})}_{\text{original density}} d\mathbf{x} = \int f(\mathbf{x}) \underbrace{\varpi_{\psi_{\theta}}(\mathbf{x})}_{\text{a new density}} \underbrace{\frac{\pi(\mathbf{x})}{\varpi_{\psi_{\theta}}(\mathbf{x})}}_{\text{importance weight}} d\mathbf{x},$$

where $f(\cdot)$ is a test function and $\pi(\cdot)$ is a multi-modal distribution.

Acceleration via importance sampling

To speed-up the simulation, we consider the importance sampling:

$$\mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x}) \underbrace{\pi(\mathbf{x})}_{\text{original density}} d\mathbf{x} = \int f(\mathbf{x}) \underbrace{\varpi_{\Psi_{\theta}}(\mathbf{x})}_{\text{a new density}} \underbrace{\frac{\pi(\mathbf{x})}{\varpi_{\Psi_{\theta}}(\mathbf{x})}}_{\text{importance weight}} d\mathbf{x},$$

where $f(\cdot)$ is a test function and $\pi(\cdot)$ is a multi-modal distribution.

We can try to simulate from an easier distribution $\varpi_{\Psi_{\theta}}(\mathbf{x})$ indirectly.

Why not simulate from a flattened distribution

A flattened distribution $\varpi_{\Psi_{\theta}}(\cdot)$ reduces the energy barrier.

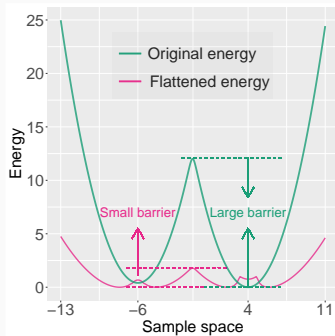


Figure 2: Energy function: original v.s. flattened.

How to construct the flattened distribution?

We partition the sample space \mathcal{X} into m subregions based on $U(\mathbf{x})$:
 $\mathcal{X}_1 = \{\mathbf{x} : U(\mathbf{x}) \leq u_1\}$, $\mathcal{X}_2 = \{\mathbf{x} : u_1 < U(\mathbf{x}) \leq u_2\}$, \dots , for some $\{u_i\}_{i=1}^m$.

How to construct the flattened distribution?

We partition the sample space \mathcal{X} into m subregions based on $U(\mathbf{x})$:
 $\mathcal{X}_1 = \{\mathbf{x} : U(\mathbf{x}) \leq u_1\}$, $\mathcal{X}_2 = \{\mathbf{x} : u_1 < U(\mathbf{x}) \leq u_2\}$, \dots , for some $\{u_i\}_{i=1}^m$.

We propose to simulate from

$$\varpi_{\Psi_{\theta}}(\mathbf{x}) \propto \frac{\pi(\mathbf{x})}{\Psi_{\theta}^{\zeta}(U(\mathbf{x}))},$$

where the importance weight $\Psi_{\theta}^{\zeta}(U(\mathbf{x}))$ satisfies

How to construct the flattened distribution?

We partition the sample space \mathcal{X} into m subregions based on $U(\mathbf{x})$:
 $\mathcal{X}_1 = \{\mathbf{x} : U(\mathbf{x}) \leq u_1\}$, $\mathcal{X}_2 = \{\mathbf{x} : u_1 < U(\mathbf{x}) \leq u_2\}$, \dots , for some $\{u_i\}_{i=1}^m$.

We propose to simulate from

$$\varpi_{\Psi_{\theta}}(\mathbf{x}) \propto \frac{\pi(\mathbf{x})}{\Psi_{\theta}^{\zeta}(U(\mathbf{x}))},$$

where the importance weight $\Psi_{\theta}^{\zeta}(U(\mathbf{x}))$ satisfies

$$(i) \ \Psi_{\theta}(U(\mathbf{x})) = \sum_{i=1}^m \theta(i) 1_{u_{i-1} < U(\mathbf{x}) \leq u_i},$$

$$(ii) \ \theta(i) = \theta_{\star}(i), \text{ where } \theta_{\star}(i) = \int_{\mathcal{X}_i} \pi(\mathbf{x}) d\mathbf{x} \text{ for } i \in \{1, 2, \dots, m\}.$$

The sampler moves like a “random walk” in the space of energy.

How to learn θ_\star on the fly (I)

However, the extension of that idea to the Langevin kernel is not straightforward. The naive setup $\Psi_\theta(U(x)) = \sum_{i=1}^m \theta(i) 1_{u_{i-1} < U(x) \leq u_i}$ leads to $\frac{\partial \log \Psi_\theta(u)}{\partial u} = 0$ a.e and fails to simulate from a new density.

How to learn θ_\star on the fly (I)

However, the extension of that idea to the Langevin kernel is not straightforward. The naive setup $\Psi_\theta(U(x)) = \sum_{i=1}^m \theta(i) 1_{u_{i-1} < U(x) \leq u_i}$ leads to $\frac{\partial \log \Psi_\theta(u)}{\partial u} = 0$ a.e and fails to simulate from a new density.

To tackle this issue, we set $\Psi_\theta(U(x))$ as piecewise continuous:

$$\Psi_\theta(U(x)) = \sum_{i=1}^m \left(\theta(i-1) e^{(\log \theta(i) - \log \theta(i-1)) \frac{U(x) - u_{i-1}}{\Delta u}} \right) 1_{u_{i-1} < U(x) \leq u_i},$$

How to learn θ_* on the fly (I)

However, the extension of that idea to the Langevin kernel is not straightforward. The naive setup $\Psi_{\theta}(U(x)) = \sum_{i=1}^m \theta(i) 1_{u_{i-1} < U(x) \leq u_i}$ leads to $\frac{\partial \log \Psi_{\theta}(u)}{\partial u} = 0$ a.e and fails to simulate from a new density.

To tackle this issue, we set $\Psi_{\theta}(U(x))$ as piecewise continuous:

$$\Psi_{\theta}(U(x)) = \sum_{i=1}^m (\theta(i-1) e^{(\log \theta(i) - \log \theta(i-1)) \frac{U(x) - u_{i-1}}{\Delta u}}) 1_{u_{i-1} < U(x) \leq u_i},$$

which leads to the desired gradient for the flattened distribution

$$\nabla_x \log \varpi_{\Psi_{\theta}}(x) = - \left[1 + \zeta \tau \frac{\log \theta(J(x)) - \log \theta((J(x) - 1) \vee 1)}{\Delta u} \right] \frac{\nabla_x U(x)}{\tau}.$$

Contour SGLD: A scalable adaptive importance sampling

Sampling step Sample \mathbf{x}_{k+1} using the SGLD algorithm

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \epsilon_{k+1} \frac{N}{n} \left[1 + \zeta \tau \frac{\log \theta_k(\tilde{J}(\mathbf{x}_k)) - \log \theta_k((\tilde{J}(\mathbf{x}_k) - 1) \vee 1)}{\Delta u} \right] \nabla_{\mathbf{x}} \tilde{U}(\mathbf{x}_k) + \sqrt{2\tau\epsilon_{k+1}} \mathbf{w}_{k+1},$$

where $\mathbf{w}_{k+1} \sim N(0, I_d)$, d is the dimension, ϵ_{k+1} is the learning rate.

$\tilde{J}(\mathbf{x})$ corresponds to the index that satisfies $u_{\tilde{J}(\mathbf{x})-1} < \frac{N}{n} \tilde{U}(\mathbf{x}) \leq u_{\tilde{J}(\mathbf{x})}$.

Contour SGLD: A scalable adaptive importance sampling

Sampling step Sample \mathbf{x}_{k+1} using the SGLD algorithm

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \epsilon_{k+1} \frac{N}{n} \left[1 + \zeta \tau \frac{\log \theta_k(\tilde{j}(\mathbf{x}_k)) - \log \theta_k((\tilde{j}(\mathbf{x}_k) - 1) \vee 1)}{\Delta u} \right] \nabla_{\mathbf{x}} \tilde{U}(\mathbf{x}_k) + \sqrt{2\tau\epsilon_{k+1}} \mathbf{w}_{k+1},$$

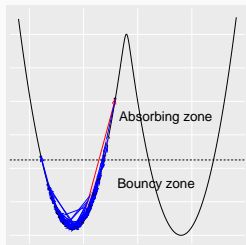
where $\mathbf{w}_{k+1} \sim N(0, I_d)$, d is the dimension, ϵ_{k+1} is the learning rate. $\tilde{j}(\mathbf{x})$ corresponds to the index that satisfies $u_{\tilde{j}(\mathbf{x})-1} < \frac{N}{n} \tilde{U}(\mathbf{x}) \leq u_{\tilde{j}(\mathbf{x})}$.

Stochastic approximation Update the estimate of θ by setting

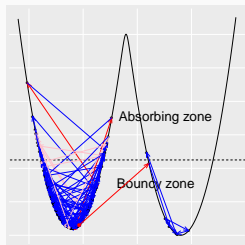
$$\theta_{k+1}(i) = \theta_k(i) + \omega_{k+1} \theta_k^\zeta(\tilde{j}(\mathbf{x}_{k+1})) \left(1_{i=\tilde{j}(\mathbf{x}_{k+1})} - \theta_k(i) \right),$$

where $1_{i=\tilde{j}(\mathbf{x}_{k+1})}$ is an indicator function for $i = 1, 2, \dots, m$.

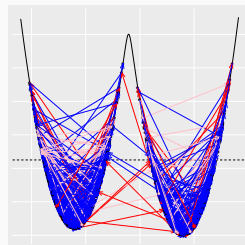
A demo of Contour SGLD



(a) CSGLD paths (early)



(b) CSGLD paths (mid)



(c) CSGLD paths (late)

Figure 3: Sample trajectories of CSGLD.

Stability for the mean-field system

Lemma (Stability)

Given a small enough ϵ (learning rate), a large enough n (batch size) and m (partition number), there is a constant $\phi > 0$ such that the mean-field system $h(\boldsymbol{\theta})$ satisfies

$$\forall \boldsymbol{\theta} \in \boldsymbol{\Theta}, \langle h(\boldsymbol{\theta}), \boldsymbol{\theta} - \boldsymbol{\theta}_\star \rangle \leq -\phi \|\boldsymbol{\theta} - \boldsymbol{\theta}_\star\|^2 + \mathcal{O}\left(\epsilon + \frac{1}{m} + \delta_n(\boldsymbol{\theta})\right),$$

where $\delta_n(\cdot)$ is a bias, $\boldsymbol{\theta}_\star = (\int_{\mathcal{X}_1} \pi(\mathbf{x}) d\mathbf{x}, \int_{\mathcal{X}_2} \pi(\mathbf{x}) d\mathbf{x}, \dots, \int_{\mathcal{X}_m} \pi(\mathbf{x}) d\mathbf{x})$.

Convergence to a unique fixed point

Theorem (L^2 convergence rate)

Given mild assumptions, $\boldsymbol{\theta}_k$ converges to a unique fixed point $\boldsymbol{\theta}_\star$ regardless of the non-convexity of $U(\mathbf{x})$ due to the stability condition.

$$\mathbb{E} [\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_\star\|^2] = \mathcal{O} \left(\omega_k + \sup_{i \geq k_0} \epsilon_i + \frac{1}{m} + \sup_{i \geq k_0} \delta_n(\boldsymbol{\theta}_i) \right).$$

Convergence to a unique fixed point

Theorem (L^2 convergence rate)

Given mild assumptions, $\boldsymbol{\theta}_k$ converges to a *unique fixed point* $\boldsymbol{\theta}_\star$ regardless of the non-convexity of $U(\mathbf{x})$ due to the stability condition.

$$\mathbb{E} [\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_\star\|^2] = \mathcal{O} \left(\omega_k + \sup_{i \geq k_0} \epsilon_i + \frac{1}{m} + \sup_{i \geq k_0} \delta_n(\boldsymbol{\theta}_i) \right).$$

Convergence to a unique fixed point

Theorem (L^2 convergence rate)

Given mild assumptions, $\boldsymbol{\theta}_k$ converges to a unique fixed point $\boldsymbol{\theta}_\star$ regardless of *the non-convexity of $U(\mathbf{x})$* due to the stability condition.

$$\mathbb{E} [\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_\star\|^2] = \mathcal{O} \left(\omega_k + \sup_{i \geq k_0} \epsilon_i + \frac{1}{m} + \sup_{i \geq k_0} \delta_n(\boldsymbol{\theta}_i) \right).$$

Convergence to a unique fixed point

Theorem (L^2 convergence rate)

Given mild assumptions, $\boldsymbol{\theta}_k$ converges to a unique fixed point $\boldsymbol{\theta}_\star$ regardless of the non-convexity of $U(\mathbf{x})$ due to the *stability condition*.

$$\mathbb{E} [\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_\star\|^2] = \mathcal{O} \left(\omega_k + \sup_{i \geq k_0} \epsilon_i + \frac{1}{m} + \sup_{i \geq k_0} \delta_n(\boldsymbol{\theta}_i) \right).$$

Convergence of the weighted averaging estimators

Theorem (Convergence of the weighted averaging estimators)

Given mild assumptions. For any bounded function f , we have

$$\left| \mathbb{E} \left[\frac{\sum_{i=1}^k \theta_i \tilde{J}(\mathbf{x}_i) f(\mathbf{x}_i)}{\sum_{i=1}^k \theta_i \tilde{J}(\mathbf{x}_i)} \right] - \int_{\mathbf{x}} f(\mathbf{x}) \pi(d\mathbf{x}) \right| \\ = \mathcal{O} \left(\frac{1}{k\epsilon} + \sqrt{\epsilon} + \sqrt{\frac{\sum_{i=1}^k \omega_k}{k}} + \frac{1}{\sqrt{m}} + \sup_{i \geq k_0} \sqrt{\delta_n(\theta_i)} \right).$$

Remark: The numerical error is slightly worse than $\mathcal{O} \left(\frac{1}{k\epsilon} + \epsilon \right)$ in the standard SGLD algorithm. This is necessary as simulating from the flattened distribution $\varpi_{\Psi_{\theta_*}}$ may lead to exponential accelerations.

Convergence of the weighted averaging estimators

Theorem (Convergence of the **weighted** averaging estimators)

Given mild assumptions. For any bounded function f , we have

$$\left| \mathbb{E} \left[\frac{\sum_{i=1}^k \theta_i \tilde{J}(\mathbf{x}_i) f(\mathbf{x}_i)}{\sum_{i=1}^k \theta_i \tilde{J}(\mathbf{x}_i)} \right] - \int_{\mathbf{x}} f(\mathbf{x}) \pi(d\mathbf{x}) \right| \\ = \mathcal{O} \left(\frac{1}{k\epsilon} + \sqrt{\epsilon} + \sqrt{\frac{\sum_{i=1}^k \omega_k}{k}} + \frac{1}{\sqrt{m}} + \sup_{i \geq k_0} \sqrt{\delta_n(\theta_i)} \right).$$

Remark: The numerical error is slightly worse than $\mathcal{O}(\frac{1}{k\epsilon} + \epsilon)$ in the standard SGLD algorithm. This is necessary as simulating from the flattened distribution $\varpi_{\Psi_{\theta_*}}$ may lead to exponential accelerations.