

# Langevin Monte Carlo 法の棄却ステップについて

南 賢太郎

東京大学大学院 情報理工学系研究科

2016 年 8 月 25 日

## 概要

MCMC の一種で, Langevin Monte Carlo 法 (LMC) と呼ばれているアルゴリズムについて概観する.

## 1 はじめに

近年, 機械学習系のベイズ的なモデルにおける推論 (特にトピックモデル系のグラフィカルモデル) で, SGLD や SGRLD といったサンプリング手法が用いられることがある. これらはマルコフ連鎖モンテカルロ法 (MCMC) の一種である. MCMC 中での位置づけは, Langevin Monte Carlo (LMC) とか Langevin Dynamics などという名前と呼ばれている既存アルゴリズムがまずあり, それに伴う勾配計算をサブサンプリングを利用して簡略化したものである. 拡張の方向性としては, 最適化における勾配降下法を確率的勾配降下法に拡張すること ( $GD \rightarrow SGD$ ) に似ている.

LMC 法を導入するモチベーションの説明として, Metropolis-Hastings 法に見られるような棄却ステップが必要ないということがよく言われる. MH 法は提案分布がうまく選べていないと棄却がたくさん発生してしまうことが知られていて, サンプリングに無駄な時間がかかってしまう. MH 法の採択確率を (ほとんど) 1 にするというのは全ペイジアン(ペイジアン: ページの)の夢であるが, その目的をある意味で達成しているのが LMC 系のアルゴリズムである. 例えば, SGLD の元論文 [1] などでは, 棄却率を低くしたいという動機をかなり強調して書いてあったりする.

この動機の部分だけ聞くと「LMC は棄却が要らない」という誤解が生じる可能性がある. しかし, 残念ながらそういうわけでもない. 実際には, 見た目上は同じアルゴリズムでも, サンプルを得たい分布の性質によって棄却が必要な場合とそうでない場合が存在するが, その境界線をひとことで説明するのは難しい.

このノートの目的は, LMC 系のアルゴリズムを概観し, 数値例を通して棄却ステップがどのように働くかを確認することである. 個人的には, 「本来は棄却しなければならない場合に棄却をせずに強行突破するとどのような不都合が起きるか」ということを理解していなかったのが, ちゃんとスクリプトを書いて確認してみた.

本ノートでは次のことを確認する.

- サンプリングしたい分布の密度関数から形式的に LMC の更新式を作ることができるが, 対応する連続時間の確率過程が定常分布を持たない場合はパスが発散するので動かない.
- 対応する連続時間の確率過程が定常分布を持つ場合でも, 棄却ステップのない LMC は収束しないことがある.

- 棄却ステップを導入すると動く場合がある。

## 2 LMC

### 2.1 確率微分方程式の離散化

LMC の基本的なアイデアは「確率微分方程式の離散化」である。次のような確率微分方程式にしたがう  $\mathbb{R}^d$ -値の確率過程をランジュバン拡散 (Langevin diffusion) という。

$$dX_t = -\nabla U(X_t)dt + \sqrt{2}dW_t, \quad X_0 = x \in \mathbb{R}^d \quad (1)$$

ただし,  $U: \mathbb{R}^d \rightarrow \mathbb{R}$  は (なめらかな) ポテンシャル関数,  $(W_t)_{t \geq 0}$  は  $d$  次元の標準ブラウン運動とする。(1) は  $t \rightarrow \infty$  である定常分布  $\pi$  に収束することが知られていて, さらに  $\pi$  は次のような確率密度をもつ:

$$\pi(x) = \frac{\exp(-U(x))}{\int_{\mathbb{R}^d} \exp(-U(x'))dx'}. \quad (2)$$

そこで, 逆にもし (2) のような分布からサンプリングを行いたいのであれば, (1) の SDE のシミュレーションを行えばよいのではないかと考えられる。

(1) のシミュレーションを Euler-Maruyama 法で行うことを考える。簡単のためステップ幅を定数  $h > 0$  とすると, 更新式は次のようになる。

$$\theta_{k+1} = \theta_k - h\nabla U(\theta_k) + \sqrt{2h}\xi_{k+1}, \quad k = 0, 1, 2, \dots \quad (3)$$

ただし,  $\xi_k$  ( $k = 1, 2, \dots$ ) は  $d$  次元の標準正規分布  $N(0, I_d)$  に従う独立な確率変数とする。つまり, 勾配  $\nabla U$  の方向に  $h$  だけ降下して, さらに正規分布に従うノイズを足すということを繰り返す。このアルゴリズムを Langevin Monte Carlo 法, あるいは後で考える MALA との比較の意味で unadjusted Langevin algorithm などと呼ぶ。

### 2.2 数値実験 1

■  $t$  分布 自由度  $\nu$  の  $t$ -分布の密度関数は

$$p(x) \propto \left(1 + \frac{x^2}{\nu}\right)^{-(\nu+1)/2}$$

である。そこで, (3) において

$$\nabla U(x) = -\nabla \log p(x) = \frac{(\nu+1)x}{\nu+x^2}$$

とすることによって,  $t$ -分布のサンプリングアルゴリズムが得られてほしい。図 1 は (3) によって拡散のシミュレーションを行ったものと, 各ステップにおける  $\theta_k$  の値のヒストグラムである。 $t$  分布の自由度は  $\nu = 1$  とし, ステップ幅は  $h = 0.2$  とした。

■ 指数的な分布  $\beta > 0$  として, 密度関数が

$$\pi_\beta(x) \propto \exp(-|x|^\beta)$$

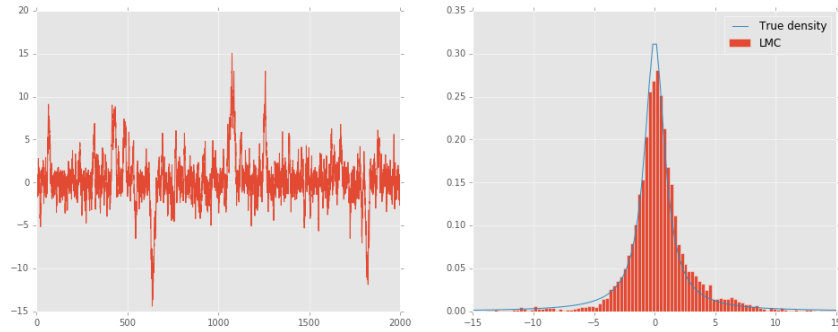


図1 左: Langevin 拡散のシミュレーション ( $h = 0.2$ ), 右:  $\theta_k$  の値のヒストグラムと真の密度関数の比較

で与えられる分布を考える. これは,  $\pi_1$  はラプラス分布,  $\pi_2$  は正規分布であり,  $\beta > 2$  のとき  $\pi_\beta$  正規分布より裾の軽い分布になる.

まず,  $\beta = 1/2$  の場合を考えてみる. このときは, 形式的に対応する確率過程 (1) はエルゴード性を持たず, 定常分布に収束しない. そのための確率過程が収束しないため, LMC 法を動かしてみてもパスが発散してしまう (図2).

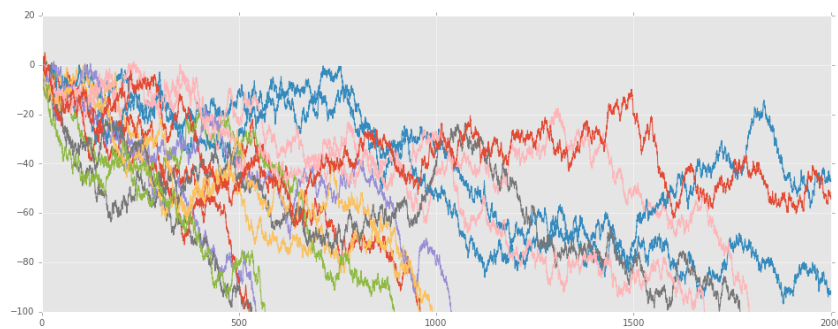


図2  $dX_t = -\frac{1}{2\sqrt{t}}dt + \sqrt{2}W_t$  のシミュレーション.

次に,  $\beta = 4$  の分布について考えてみる.

(1) は

$$dX_t = -4t^3 dt + \sqrt{2}dW_t$$

となる.  $\beta = 1/2$  の場合と違って, こちらは  $\pi_4$  をちゃんと定常分布にもつことが知られている [2]. さらにいうと,  $X_t$  の分布が  $\pi_4$  に total variation の意味で収束する速さが, ある定数  $c > 0$  を用いて  $O(e^{-ct})$  と書けることが知られている (このような性質を幾何エルゴード性という). よって, この確率過程を離散化することで  $\pi_4$  からのサンプリングアルゴリズムが得られるような気持ちになるが, 残念ながら現実にはそうはならない. 実際, この場合には (3) で得られるマルコフ連鎖は, ステップ幅をどのように設定しても一時的 (transient) になってしまうことが知られている [2].

実際にシミュレーションしてみると, 最初の数ステップは正しいサンプリングを行えているように見えるものの, どのパスも有限ステップで発散してしまっている (図2).

上の例から, 興味のある分布が,

- 確率過程 (1) が定常分布を持たない.

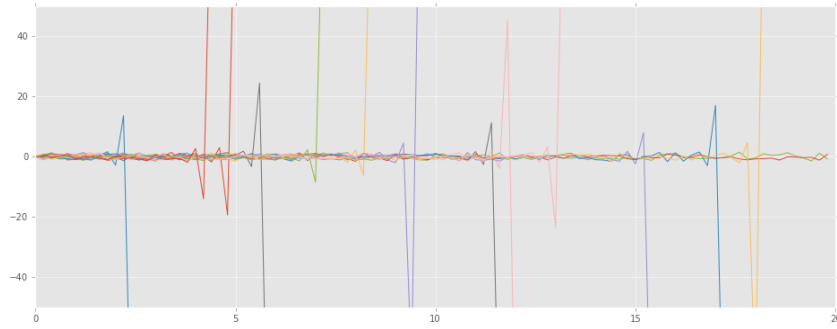


図3  $dX_t = -4t^3 dt + \sqrt{2}dW_t$  を Euler-Maruyama 法で離散化したもの. ステップ幅は  $h = 0.2$  として, パスは 15 本生成した. 定常分布は  $\pi_4(x) \propto e^{-x^4}$  であるので, だいたい  $[-1, 1]$  付近をさまよるのが「望ましい」挙動である. しかし, シミュレーションではどのパスも有限時刻で絶対値の大きさが爆発してしまう.

- 確率過程 (1) は定常分布をもつが, Euler-Maruyama 近似を行ったものは定常分布を持たない.

のどちらかのケースに該当すると, LMC 法はうまく動かないことがわかる. 大雑把にいうと,  $\lim_{|x| \rightarrow \infty} |\nabla U(x)| = 0$  となるときはそもそも (1) が定常分布を持たず,  $\nabla U(x)$  が Lipschitz でないときは Euler-Maruyama 近似がうまくいかない.

### 3 Metropolis Adjusted Langevin Algorithm

そこで, LMC に Metropolis-Hastings 法の棄却ステップを導入したものが Metropolis Adjusted Langevin Algorithm (MALA) である. (3) のマルコフ連鎖の推移確率は

$$q_h(x, y) \propto \exp\left(-\frac{|y - x + h\nabla U(x)|^2}{4h}\right)$$

と書ける. そこで, 次のようなアルゴリズムを考える.

---

**Algorithm 1** Metropolis Adjusted Langevin Algorithm (MALA)

---

**for**  $k = 0, 1, 2, \dots$  **do**

**while**  $\tilde{\theta}_{k+1}$  is not accepted **do**

$\xi_{k+1} \sim N_d(0, 1)$

$\tilde{\theta}_{k+1} \leftarrow \theta_k - h\nabla U(\theta_k) + \sqrt{2h}\xi_{k+1}$

    Accept  $\tilde{\theta}_{k+1}$  with probability

$$\min\left\{1, \frac{\pi(\tilde{\theta}_{k+1})q_h(\theta_k, \tilde{\theta}_{k+1})}{\pi(\theta_k)q_h(\tilde{\theta}_{k+1}, \theta_k)}\right\}$$

**end while**

**end for**

---

### 3.1 数値実験 2

再び  $\pi_4(x) \propto \exp(-x^4)$  を扱う. この例では, 元になる連続時間の確率過程 (1) が, 定常分布に指数的な速さで収束するにもかかわらず, 棄却のない LMC アルゴリズムは発散してしまうのだった.

図 4 は MALA によって  $\pi_4$  からのサンプリングを行ったものである.  $h = 0.01$  として,  $K = 10,000$  ステップまで計算した.

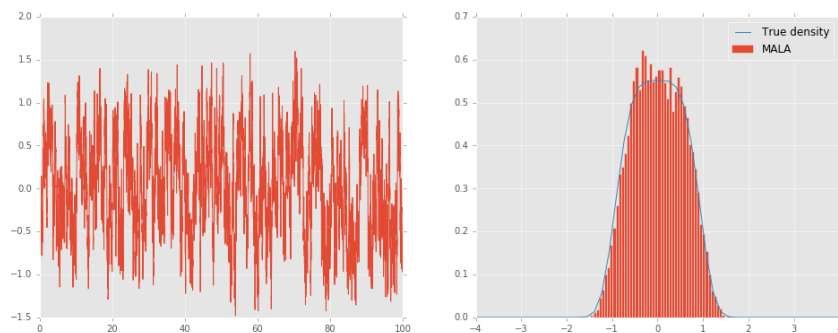


図 4 Metropolis Adjusted Langevin Algorithm による  $\pi_4(x)$  からのサンプリング.

この例では, MALA の定常分布と (1) の定常分布は一致していることがわかる. ただし, 「収束の速さ」は離散化によって保存しておらず, MALA によって得られるマルコフ連鎖は幾何エルゴード性を持たない [2].

### 参考文献

- [1] M. Welling and Y. W. Teh. Bayesian Learning via Stochastic Gradient Langevin Dynamics. In *ICML*, 2011.
- [2] G. O. Roberts and R. L. Tweedie. Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli*, 2(4):341–363, 1996.