

On Markov Chain Gradient Descent

Tao Sun (NUDT), Yuejiao Sun (UCLA), and Wotao Yin (UCLA)

nudtsuntao@163.com, sunyj@math.ucla.edu, wotaoyin@math.ucla.edu

PROBLEM

$$\underset{x \in X \subseteq \mathbb{R}^n}{\text{minimize}} \mathbb{E}_{\xi} (F(x; \xi)) \quad (1)$$

where Π is the distribution of a sample space Ξ , X is closed and convex, $F(\cdot, \xi) : X \rightarrow \mathbb{R}$ is convex or nonconvex but differentiable, associated with $\xi \in \Xi$.

MOTIVATION

Instead of **Stochastic Gradient Descent (SGD)** (i.i.d sample ξ^k):

$$x^{k+1} = \mathbf{Proj}_X (x^k - \gamma_k \partial F(x^k; \xi^k))$$

We use ξ^k on a Markov-chain trajectory, and call this method **Markov Chain Gradient Descent (MCGD)**. Benefits:

- some distributions (e.g., $\Xi := \{x \in \{0, 1\}^n | \langle a, x \rangle \leq b\}$) cannot be sampled directly, but they have Markov-chain samples.
- Markov chains naturally arise in some applications, e.g. linear dynamic systems with random transitions or errors, and distributed systems in which each node stores a subset of training samples.

ADVANTAGE OF MCGD OVER SGD

We use a numerical example to illustrate the advantage of MCGD over SGD. Consider an auto regressive model:

$$\begin{aligned} \xi_t^1 &= A\xi_{t-1}^1 + e_1 W_t, \quad W_t \stackrel{\text{i.i.d}}{\sim} N(0, 1) \\ \xi_t^2 &= \begin{cases} 1, & \text{if } \langle u, \xi_t^1 \rangle > 0, \\ 0, & \text{otherwise;} \end{cases} \quad \xi_t^2 = \begin{cases} \bar{\xi}_t^2, & \text{with probability 0.8,} \\ 1 - \bar{\xi}_t^2, & \text{with probability 0.2.} \end{cases} \end{aligned}$$

Clearly, $(\xi_t^1, \xi_t^2)_{t=1}^{\infty}$ forms a Markov chain. Let Π denote the stationary distribution of this Markov chain. We recover u as the solution to the following problem:

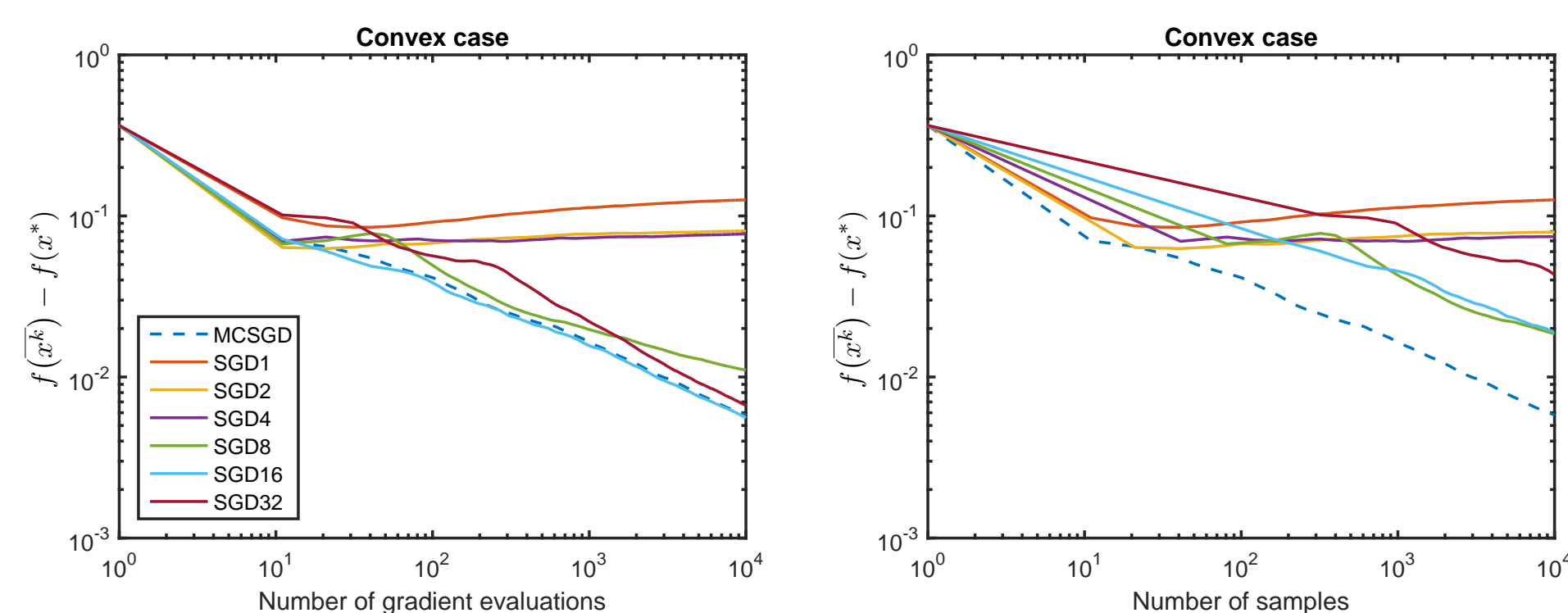
$$\underset{x}{\text{minimize}} \quad \mathbb{E}_{(\xi^1, \xi^2) \sim \Pi} \ell(x; \xi^1, \xi^2).$$

Compare:

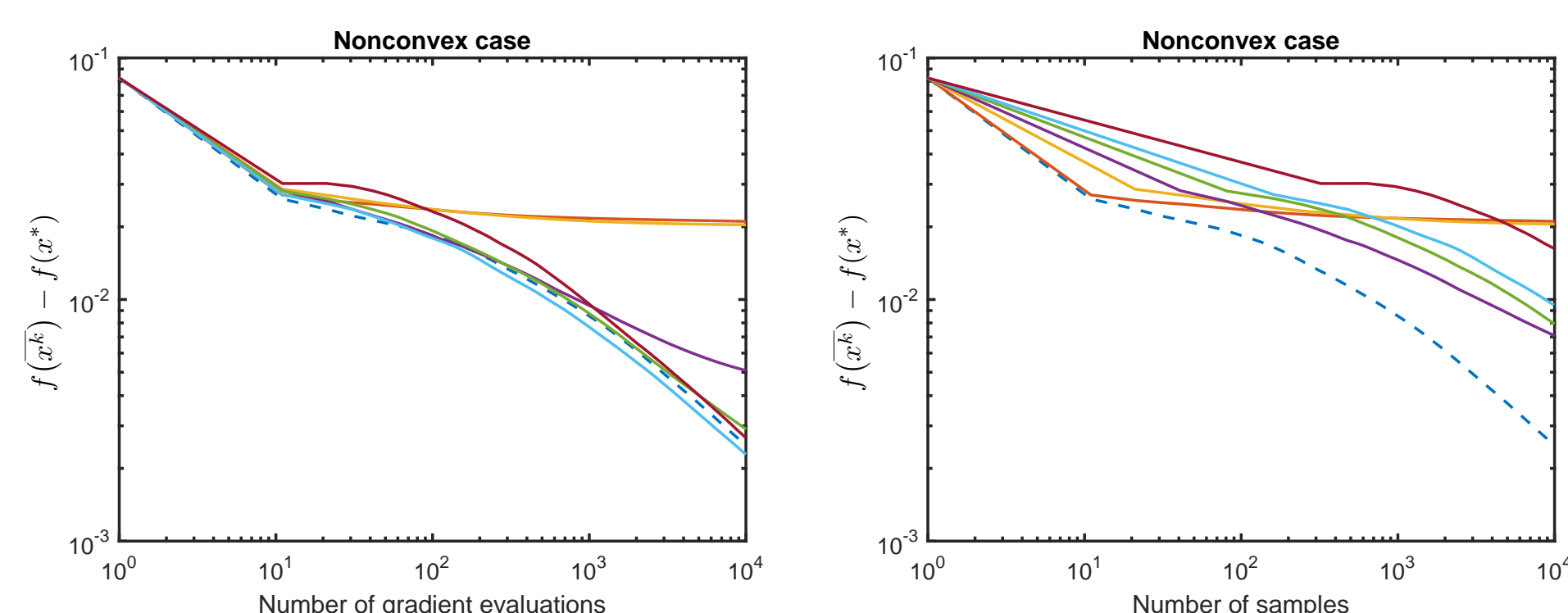
- MCGD, where samples are taken from one trajectory of the Markov chain;
- SGDT, where T is the burn-in time and each sample is the T th sample of a fresh, independent trajectory.

We test on both convex and nonconvex loss functions. ($\sigma(t) = \frac{1}{1+\exp(-t)}$).

$$\ell(x; \xi^1, \xi^2) = -\xi^2 \log(\sigma(\langle x, \xi^1 \rangle)) - (1 - \xi^2) \log(1 - \sigma(\langle x, \xi^1 \rangle)).$$



$$\ell(x; \xi^1, \xi^2) = \frac{1}{2} (\sigma(\langle x, \xi^1 \rangle) - \xi^2)^2.$$



CONTRIBUTIONS

All analyses of MCGD must deal with the biased expectation. Existing MCGD work [1-4]. New results of this work:

- allow **non-reversible** Markov chain for faster convergence;
- objective can be nonconvex.
- non-ergodic convergence of the objective.

CONVERGENCE ANALYSIS

FINITE STATE SPACE

The analysis is based on the following assumptions:

- The Markov chain $(X_k)_{k \geq 0}$ is time-homogeneous, irreducible, and aperiodic. It has a transition matrix P and stationary distribution π^* .
- X is convex and compact.
- Ξ is finite. Let $f_i(x) = M \cdot \mathbf{Prob}(\xi = y^i) \cdot F(x, y^i)$, and reformulate (1) as

$$\underset{x \in X \subseteq \mathbb{R}^d}{\text{minimize}} f(x) = \frac{1}{M} \sum_{i=1}^M f_i(x),$$

where each state i has the uniform probability $1/M$. MCGD runs

$$x^{k+1} = \mathbf{Proj}_X (x^k - \gamma_k \partial f_{j_k}(x^k)),$$

where $(j_k)_{k \geq 0}$ forms a Markov chain trajectory. It can be illustrated in the following diagram:

$$\begin{array}{ccccccc} & & j_0 & \rightarrow & j_1 & \rightarrow & j_2 & \rightarrow & \dots \\ & & \downarrow & & \downarrow & & \downarrow & & \\ x^0 & \rightarrow & x^1 & \rightarrow & x^2 & \rightarrow & x^3 & \rightarrow & \dots \end{array}$$

CONVEX CASE

Assume that f_i , $i \in [M]$, are convex functions, and the stepsizes satisfy

$$\sum_k \gamma_k = +\infty, \quad \sum_k \ln k \cdot \gamma_k^2 < +\infty. \quad (2)$$

Let $\bar{x}^k = \frac{\sum_{i=1}^k \gamma_i x^i}{\sum_{i=1}^k \gamma_i}$, and we have

$$\begin{aligned} \lim_k \mathbb{E} f(x^k) &= f^*, \\ \mathbb{E} (f(\bar{x}^k) - f^*) &= O\left(\frac{\Phi(P)}{\sum_{i=1}^k \gamma_i}\right). \end{aligned}$$

When choosing $\gamma_k = O(\frac{1}{k^q})$, $\frac{1}{2} < q < 1$,

$$\mathbb{E} (f(\bar{x}^k) - f^*) = O\left(\frac{\Phi(P)}{k^{1-q}}\right).$$

Note that the same stepsize and convergence rate can hold for SGD and subgradient algorithms.

NONCONVEX CASE

Assume that $X = \mathbb{R}^n$, f_i is differentiable and ∇f_i is L -Lipschitz and bounded by $D > 0$, and

$$\sum_k \gamma_k = +\infty, \quad \sum_k \ln^2 k \cdot \gamma_k^2 < +\infty. \quad (3)$$

Let $f^* = \min_{x \in X} f(x)$. Then, we have

$$\begin{aligned} \lim_k \mathbb{E} \|\nabla f(x^k)\| &= 0, \\ \mathbb{E} \left(\min_{1 \leq i \leq k} \{\|\nabla f(x^i)\|^2\} \right) &= O\left(\frac{\Phi(P)}{\sum_{i=1}^k \gamma_i}\right). \end{aligned}$$

When choosing $\gamma_k = O(\frac{1}{k^q})$, $\frac{1}{2} < q < 1$, the convergence rate is $O(\frac{\Phi(P)}{k^{1-q}})$.

CONTINUOUS STATE SPACE

Assume that state space Ξ is a continuum. Consider an infinite-state Markov chain that is time-homogeneous and reversible and solve (1) by MCGD.

CONVEX CASE

Assume that for each $\xi \in \Xi$, $F(\cdot; \xi)$ is convex, $|F(x; \xi) - F(y; \xi)| \leq L\|x - y\|$, $\sup_{x \in X, \xi \in \Xi} \{\|\hat{\nabla} F(x; \xi)\|\} \leq D$, $\mathbb{E}_{\xi} \hat{\nabla} F(x; \xi) \in \partial \mathbb{E}_{\xi} F(x; \xi)$, and $\sup_{x, y \in X, \xi \in \Xi} |F(x; \xi) - F(y; \xi)| \leq H$. Choose γ^k according to (2). Let $F^* := \min_{x \in X} \mathbb{E}_{\xi} (F(x; \xi))$. $\lambda \in (0, 1)$ is the geometric rate of the mixing time of the Markov chain. Then we have

$$\begin{aligned} \lim_k \mathbb{E} (\mathbb{E}_{\xi} (F(x^k; \xi)) - F^*) &= 0, \\ \mathbb{E} (\mathbb{E}_{\xi} (F(\bar{x}^k; \xi)) - F^*) &= O\left(\frac{\max\{1, \frac{1}{\ln(1/\lambda)}\}}{\sum_{i=1}^k \gamma_i}\right), \end{aligned}$$

NONCONVEX CASE

Let $X = \mathbb{R}^n$. Assume for any $\xi \in \Xi$, $F(x; \xi)$ is differentiable, and $\|\nabla F(x; \xi) - \nabla F(y; \xi)\| \leq L\|x - y\|$. In addition, $\sup_{x \in X, \xi \in \Xi} \{\|\nabla F(x; \xi)\|\} < +\infty$, X is the full space, and $\mathbb{E}_{\xi} \nabla F(x; \xi) = \nabla \mathbb{E}_{\xi} F(x; \xi)$. Then, we have

$$\begin{aligned} \lim_k \mathbb{E} \|\nabla \mathbb{E}_{\xi} (F(x^k; \xi))\| &= 0, \\ \mathbb{E} \left(\min_{1 \leq i \leq k} \{\|\nabla \mathbb{E}_{\xi} (F(x^i; \xi))\|^2\} \right) &= 0. \end{aligned}$$

ACCELERATION DUE TO NON-REVERSIBILITY

Non-reversibility can accelerate the mixing process of Markov chains. The following experiment compares MCGD with reversible and non-reversible Markov chains over the same graph with 20 nodes. The objective is a least square problem with data distributed on the graph.

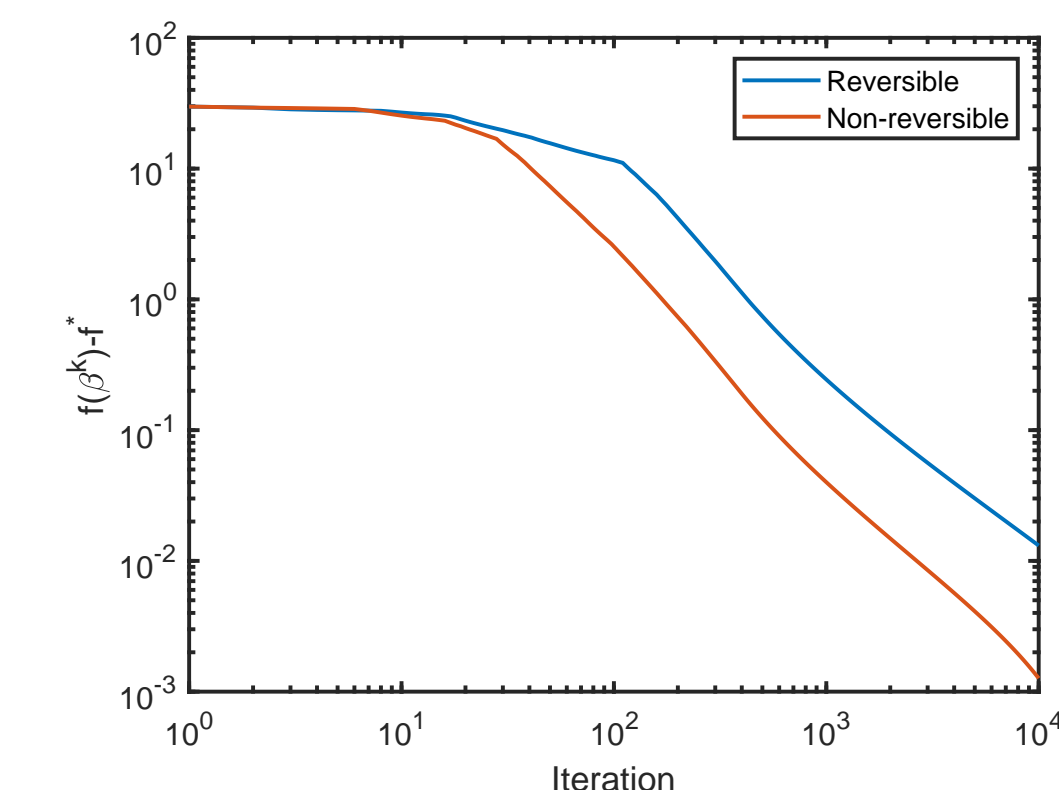


Figure 1: The second largest eigenvalues of reversible and non-reversible Markov chains are 0.75 and 0.66 respectively.

REFERENCES

- JC Duchi, A Agarwal, M Johansson, MI Jordan. *Ergodic mirror descent*. SIAM Journal on Optimization, 2012.
- B Johansson, M Rabi, M Johansson. *A simple peer-to-peer algorithm for distributed optimization in sensor networks*. In Decision and Control, IEEE, 2007.
- B Johansson, M Rabi, M Johansson. *A randomized incremental subgradient method for distributed optimization in networked systems*. SIAM Journal on Optimization, 2009.
- SS Ram, A Nedic, VV Veeravalli. *Incremental stochastic subgradient algorithms for convex optimization*. SIAM Journal on Optimization, 2009.