



Toward Deeper Understanding of Momentum in Nonconvex Stochastic Optimization

Tianyi Liu*, Zhehui Chen*, Enlu Zhou*, Tuo Zhao*
 *School of Industrial and System Engineering at Georgia Institute of Technology



Background

Consider an empirical risk minimization problem,

$$\min_{\theta} \mathcal{F}(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i, \theta)),$$

- $\{(x_i, y_i)\}_{i=1}^n$: n observations.
- ℓ : loss function.
- f : decision function associated with θ .

A Popular Algorithm: **Momentum SGD (Heavy Ball)**

$$\begin{aligned} \theta_{k+1} &= \theta_k - \eta \nabla \ell(y_i, f(x_i, \theta_k)) \xrightarrow{\text{SGD}} \\ &\quad + \mu(\theta_k - \theta_{k-1}) \xrightarrow{\text{Momentum}} \end{aligned}$$

- η : step size.
- $\mu \in [0, 1]$: momentum parameter.

Questions:

- Does Momentum **accelerate** the algorithm?
- Does Momentum **help escape from saddle**?

Streaming PCA

• A simple, but nontrivial example:

$$\min_v -v^\top \mathbb{E}_{x \sim \mathcal{D}}[xx^\top]v$$

subject to $v \in \mathcal{M}$,

where $\mathcal{M} = \{v \in \mathbb{R}^p : v^\top v = 1\}$ is the stiefel manifold.

• Streaming data:

At the i -th iteration, independently sample $x_i \sim \mathcal{D}$.

• Assumption 1 (Eigen-gap):

$\mathbb{E}[x] = 0$ and $\Sigma = \mathbb{E}[xx^\top]$ is positive definite with eigenvalues

$$\lambda_1 > \lambda_2 \geq \dots \geq \lambda_p > 0$$

associated with unit eigenvectors $v^{(1)}, v^{(2)}, \dots, v^{(p)}$.

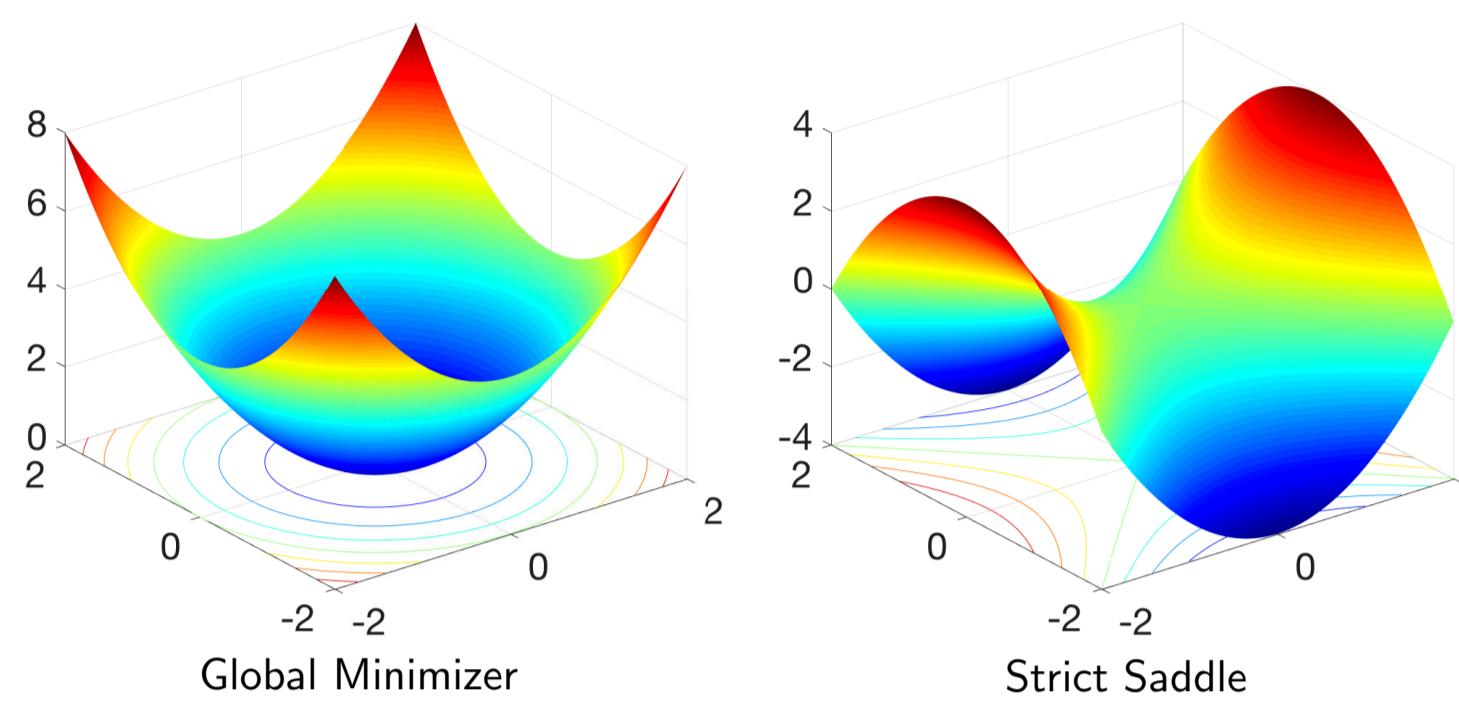
• Assumption 2 (Bounded Input):

$$\|x\|_2 \leq C_p,$$

where C_p is a constant (possibly depends on p).

Landscape

- Stationary point without constraints: $\nabla \mathcal{F}(v) = 0$.



- Stationary point for Streaming PCA:
Manifold Gradient:

$$\nabla_{\mathcal{M}} \mathcal{F}(v) = (I - vv^\top) \Sigma v = 0.$$

\Rightarrow **Stationary Points:**

$$\underbrace{\pm v^{(1)}}_{\text{Global Minimizers}}, \underbrace{\pm v^{(2)}, \dots, \pm v^{(p)}}_{\text{Strict Saddle and Maximizers}}.$$

Intuition

- MSGD for Streaming PCA:

$$v_{k+1} = v_k + \mu(v_k - v_{k-1}) + \eta \underbrace{(I - v_k v_k^\top) X_k X_k^\top v_k}_{\nabla_{\mathcal{M}} F_k(v_k): \text{Stoc. Approx. of } \nabla_{\mathcal{M}} \mathcal{F}(v_k)}.$$

- ODE Approximation:

$$\text{Discrete: } \frac{v_{k+1} - v_k}{\eta} = \mu \frac{v_k - v_{k-1}}{\eta} + \nabla_{\mathcal{M}} F_k(v_k).$$

weakly $\Downarrow \eta \rightarrow 0$

$$\text{Continuous: } dV = \mu dV + \nabla_{\mathcal{M}} \mathcal{F}(V).$$

Discrete/Stochastic \rightarrow Continuous/Deterministic.
Similar to the Law of Large Number, not reliable!

- SDE Approximation ($v_k^\eta - v^{(i)} = O(\sqrt{\eta})$):

Consider the normalized error

$$u_k^\eta = \frac{v_k^\eta - v^{(i)}}{\sqrt{\eta}}.$$

$$\text{Discrete: } \frac{u_{k+1} - u_k}{\sqrt{\eta}} = \mu \frac{u_k - u_{k-1}}{\sqrt{\eta}} + \nabla_{\mathcal{M}} F_k(u_k).$$

weakly $\Downarrow \eta \rightarrow 0$

$$\text{Continuous: } dU = \mu dU + \nabla_{\mathcal{M}} \mathcal{F}(U) dt + \Sigma dB_t.$$

Randomness Returns.
Similar to the Central Limit Theorem!

Diffusion Approximation

Theorem 1 (Global Convergence).

Chosen $\mu \in [0, 1)$, $v_0 \in \mathcal{M}$, denote $V^\eta(t) = v_{\lfloor t/\eta \rfloor}$.

$$V^\eta(\cdot) \xrightarrow[\text{weakly converge}]{\eta \rightarrow 0} V(\cdot) \text{ satisfying}$$

$$\text{ODE } dV = \frac{1}{1-\mu} [\Sigma V - V^\top \Sigma V] dt, V(0) = v_0.$$

Theorem 2 (Local Behavior).

Denote $U^\eta(t) = u_{\lfloor t/\eta \rfloor}^\eta$. Based on $v_k^\eta - v^{(i)} = O(\sqrt{\eta})$,

$$\text{for } i \neq j, U_i^\eta(t) \xrightarrow[\text{weakly converge}]{\eta \rightarrow 0} U_i(t) \text{ satisfying}$$

$$\text{SDE } dU_i = \underbrace{\left(\frac{\lambda_i - \lambda_j}{1-\mu}\right) U_i dt}_{\text{Drift}} + \underbrace{\left(\frac{\alpha_{i,j}}{1-\mu}\right) dB_t}_{\text{Noise}}.$$

Proof Technique: Fixed-State-Chain [1].

Remark 3. $O-U$ processes $\{U_i\}$ characterize the local algorithmic behavior:

- $\lambda_i < \lambda_j$, process converges.
- $\lambda_i > \lambda_j$, process diverges.

Reference: [1] Kushner and Yin. Stochastic approximation and recursive algorithms and applications, Stochastic modelling and applied probability, vol. 35.

Implication

• **Phase I (Around Saddle):** Momentum enlarges variance, encourages exploitation and **helps escape!**

• **Phase II (Traverse between stationary points):** Drift dominating variance leads to deterministic trajectory.

• **Phase III (Around Optima):** Momentum enlarges variance and hurts convergence. **Step size annealing!**

| Phase | Step Size η | Asymptotic Complexity |
|-------|---|---|
| I | $\frac{(\lambda_1 - \lambda_2)\epsilon}{\phi}$ | $\frac{(1-\mu)\phi}{\epsilon(\lambda_1 - \lambda_2)^2} \log\left(\frac{(1-\mu)\delta^2}{\epsilon}\right)$ |
| II | $\frac{(\lambda_1 - \lambda_2)\epsilon}{\phi}$ | $\frac{(1-\mu)\phi}{\epsilon(\lambda_1 - \lambda_2)^2} \log\left(\frac{1-\delta^2}{\delta^2}\right)$ |
| III | $\frac{(1-\mu)(\lambda_1 - \lambda_2)\epsilon}{\phi}$ | $\frac{\phi}{\epsilon(\lambda_1 - \lambda_2)^2} \log\left(\frac{\delta^2}{\epsilon}\right)$ |

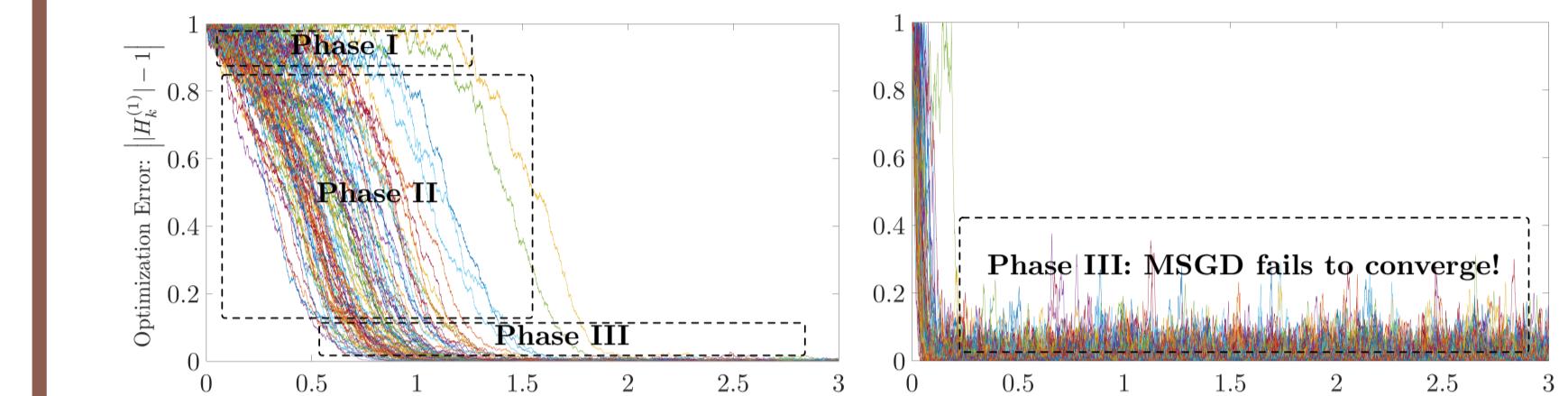
ϕ is a constant related to $\alpha_{i,j}$.

Momentum **accelerates** the algorithm in Phase I and II.

Experiments

- Streaming PCA (Synthetic Data):

$$\Sigma = \text{diag}\{4, 3, 2, 1\}, \mu = 0.9, \eta = 5 \times 10^{-5}.$$



◊ Momentum helps escape from saddle.

◊ MSGD converges after step size annealing.

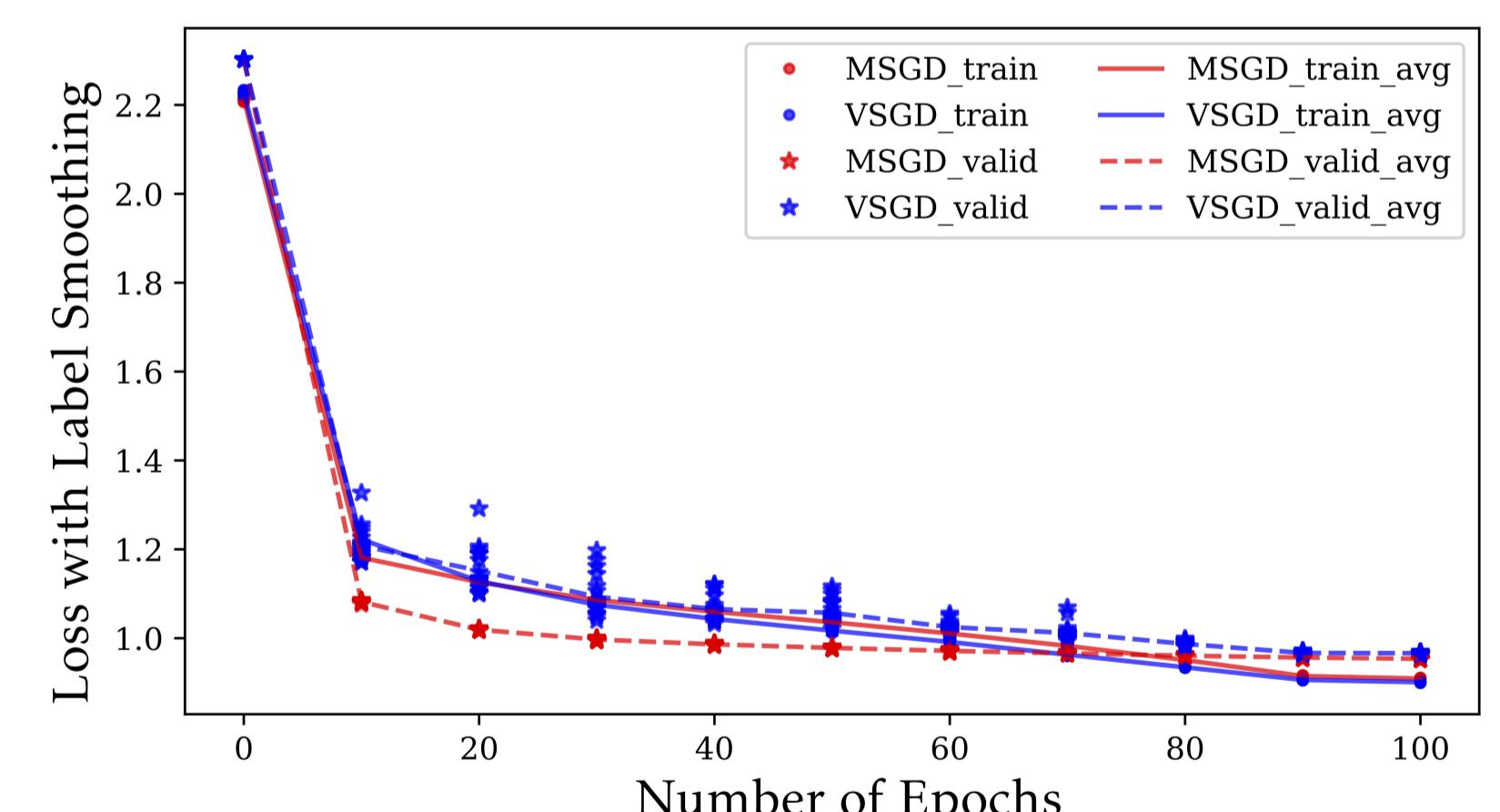
- Deep Neural Network with optimal setting:

Architecture: ResNet-9; Dataset: CIFAR-10;

Peak Step Size: $\eta_V = 2$; $\eta_M = 0.36$ with $\mu = 0.9$;

Step Size Linear Schedule:

0 – 20 epochs: $0 \nearrow \eta$; 21 – 100 epochs: $\eta \searrow 0$.



◊ The best setting of MSGD escapes from saddle much **faster** than that of VSGD and eventually **outperforms** that of VSGD.

Open Questions

1. Extension to more general problems;
2. Bridging asymptotic and nonasymptotic analyses.