
Stochastic Gradient Descent with Momentum Convergence Diagnostic for Nonconvex Optimization

Anonymous Author(s)

Affiliation

Address

email

1 Nonconvex case

We recall the SGD with Momentum update we are analyzing here:

$$\theta_{n+1} = \theta_n - \gamma_{n+1} \nabla \ell(\theta_n, \xi_{n+1}) + \beta(\theta_n - \theta_{n-1}) \quad (1)$$

where ℓ is the nonconvex loss function parametrized by $\theta \in \Theta \subset \mathbb{R}^p$ and ξ is some random noise.
 $\beta \in [0, 1)$ is the momentum parameter and γ_{n+1} the learning stepsize.

We also define $f(\theta) = \mathbb{E}[\ell(\theta, \xi)]$ the expected loss. Following Pflug convergence diagnostic test, we construct the following test statistics:

$$\nabla \ell(\theta_n, \xi_{n+1})^\top \nabla \ell(\theta_{n-1}, \xi_n) \quad (2)$$

and the goal will be to upperbound the expectation of this quantity in order to spot the two different phases through the iterates.

We make the following assumptions before analyzing (2).

H1. The loss function $\ell(\theta, \xi)$ is nonconvex w.r.t. the parameter θ .

We consider the very general setting where the loss function $\ell(\cdot; \xi)$ is (l, L) -smooth, see [Allen-Zhu, 2017, Zhou and Gu, 2019]

H2. There exist some constant $l \in \mathbb{R}$ and $L > 0$ such that for $(\theta, \vartheta) \in \Theta^2$:

$$\frac{l}{2} \|\theta - \vartheta\|^2 \leq \ell(\theta) - \ell(\vartheta) - \nabla \ell(\vartheta)^\top (\theta - \vartheta) \leq \frac{L}{2} \|\theta - \vartheta\|^2 \quad (3)$$

Note that if $l = -L$ we recover the conventional L -smoothness definition and if $l \geq 0$ (resp. $l > 0$) we have convexity (resp. strong convexity).

H3. There exists $K > 1$ such that

$$\mathbb{E} \left[(\theta_n - \theta_{n-1})^\top (\theta_{n-1} - \theta_{n-2}) \right] \geq -K \mathbb{E} \left[\|\theta_n - \theta_{n-1}\|^2 \right]$$

for large enough iteration index n .

Finally and classically (see [Ghadimi and Lan, 2013]) in nonconvex optimization, we make an assumption on the magnitude of the gradient:

H4. There exists a constant $G > 0$ such that

$$\|\nabla \ell(\theta, \xi)\| < G \quad \text{for any } \theta \text{ and } \xi$$

We recall an important convergence result for the SGD with Momentum update from [Yan et al., 2018]:

21 **Theorem 1.** [Yan et al., 2018] Under assumptions H 1, H 2, H 4 and the boundedness of the
 22 variance of the stochastic gradients, we have

$$\begin{aligned} \min_{k=0,\dots,n} \mathbb{E} \left[\|\nabla \ell(\theta_k)\|^2 \right] &\leq \frac{2(\ell(\theta_0) - \ell_*)(1-\beta)}{n+1} \max \left\{ \frac{2L}{1-\beta}, \frac{\sqrt{n+1}}{C} \right\} \\ &\quad + \frac{C}{\sqrt{n+1}} \frac{L\beta^2((1-\beta)s-1)^2(G^2 + \sigma^2) + L\sigma^2(1-\beta)^2}{(1-\beta)^3} \end{aligned} \quad (4)$$

23 We can easily check the following identity:

$$\nabla \ell(\theta_n, \xi_{n+1})^\top \nabla \ell(\theta_{n-1}, \xi_n) = \frac{1}{\gamma} \nabla \ell(\theta_n, \xi_{n+1})^\top (\theta_{n-1} - \theta_n) + \frac{\beta}{\gamma} \nabla \ell(\theta_n, \xi_{n+1})^\top (\theta_{n-1} - \theta_{n-2}) \quad (5)$$

24 Taking expectations on both sides and using Assumption H 2, we have:

$$\begin{aligned} \mathbb{E} \left[\nabla \ell(\theta_n, \xi_{n+1})^\top \nabla \ell(\theta_{n-1}, \xi_n) \right] &\leq \frac{1}{\gamma} \left[f(\theta_{n-1}) - f(\theta_n) - \frac{l}{2} \|\theta_{n-1} - \theta_n\|^2 \right] \\ &\quad + \frac{\beta}{\gamma} \left[f(\theta_n) - f(\theta_n + \theta_{n-2} - \theta_{n-1}) + \frac{L}{2} \|\theta_{n-1} - \theta_{n-2}\|^2 \right] \end{aligned} \quad (6)$$

25 which yields:

$$\begin{aligned} \mathbb{E} \left[\nabla \ell(\theta_n, \xi_{n+1})^\top \nabla \ell(\theta_{n-1}, \xi_n) \right] &\leq \frac{1}{\gamma} \left[f(\theta_{n-1}) - f(\theta^*) - \frac{l}{2} \|\theta_{n-1} - \theta_n\|^2 \right] \\ &\quad + \frac{\beta}{\gamma} \left[f(\theta_n) - f(\theta^*) + \frac{L}{2} \|\theta_{n-1} - \theta_{n-2}\|^2 \right] \end{aligned} \quad (7)$$

26 where θ^* is the global minimizer of the expected loss.

27 Denote $\Delta_n = \theta_n - \theta_{n-1}$ and observe that:

$$\|\Delta_n\|^2 = \gamma^2 \|\nabla \ell(\theta_{n-1}, \xi_n)\|^2 + 2\beta \Delta_n^\top \Delta_{n-1} - \beta^2 \|\Delta_{n-1}\|^2 \quad (8)$$

28 Using assumptions H 3 and using Theorem 1 we obtain:

$$\mathbb{E} \|\Delta_n\|^2 \leq \gamma^2 G^2 - (2\beta K + \beta^2) \mathbb{E} \left[\|\Delta_{n-1}\|^2 \right] \quad (9)$$

29 **References**

- 30 Z. Allen-Zhu. Natasha: Faster non-convex stochastic optimization via strongly non-convex parameter. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages
31 89–97. JMLR. org, 2017.
- 33 S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- 35 Y. Yan, T. Yang, Z. Li, Q. Lin, and Y. Yang. A unified analysis of stochastic momentum methods
36 for deep learning. *arXiv preprint arXiv:1808.10396*, 2018.
- 37 D. Zhou and Q. Gu. Lower bounds for smooth nonconvex finite-sum optimization. *arXiv preprint*
38 *arXiv:1901.11224*, 2019.