
Convergence of Optimizers Without Bounded Gradient Assumption

Yufeng (Roger) Wang, Eric Xia
rogerwyf@uw.edu, ericxia@uw.edu

1 Proposal

The primary goal of most deep learning models is to minimize the loss function. Optimizers are crucial to this end by fine-tuning hyperparameters such that the model will actually converge on a minimum instead of overshooting or moving away from it, and do so in a computationally efficient manner. Whereas existing optimizers are intuitively straightforward in convex learning, in non-convex settings they (notably for Adam-type adaptive gradient methods) often require the strong assumption of a universally bounded gradient for guaranteeing deterministic behavior: there exists some $G > 0$ s.t.

$$\|\nabla f(x)\| \leq G$$

for all $x \in \mathbb{R}^d$, where $f(x)$ is the loss function and d is the dimensions the cost function has. While convenient, the imposition of this assumption can be hard to verify in practical settings, hence in recent years there has been a trend in research efforts towards analyzing existing optimization methods and proposing novel methods without this assumption.

We plan to survey and analyze some recent papers in this realm related to RMSProp [1], SGD [2], SGDM [3], ADAM [4], and self-tuning step sizes [5], comparing and contrasting their approaches, methodologies, and findings, and discussing their broader implications for the convergence problem and the general field of optimizers.

As a non-momentum specificity of Adam, RMSProp is known for performing well empirically, despite theoretical analyses of it suggesting divergence even for simple complex functions. [1]. With their paper published at ICRL 2021, Shi et al. examined a counter-example presented by Reddi et al. [6] in 2018 in their study of the convergence of Adam, discovering a gap in their analysis of convergence for large β_2 parameter. Shi et al. ran simulations based on this gap. Plotting whether convergence was achieved for different choices of β_2 for RMSprop, they found that in general, there was a continuous and nontrivial curve sloping upwards, from divergence to convergence, suggesting their conjecture that RMSprop converges for large enough β_2 . This paper is notable and interesting as it is the first to prove the convergence of RMSprop without the bounded gradient assumption.

In their 2021 paper "SGD for Structured Nonconvex Functions: Learning Rates, Minibatching and Interpolation", Gower et al. show that SGD converges at a rate of $\mathcal{O}(1/\sqrt{k})$ on quasar (strongly) convex functions and prove linear convergence to a neighborhood for functions satisfying the Polyak-Lojasiewicz condition, a generalization of strongly-convex functions, without any bounded gradient assumption, instead relying on the expected residual (ER) condition for support. [2] The significance of this study resides in the provided insights on the complexity of minibatching and determination of optimal batch sizes, as well as the demonstration that for models interpolating the training data, the ER condition could be discarded and still give SOTA results.

The third paper we chose is interesting for a similar reason: it is the first convergence guarantee for Multistage SGDM—Multistage refers to applying a constant stepsize which is then dropped by a constant factor to encourage fine-tuning of training, and the momentum weight is either kept unchanged or gradually increased. Although SGDM has an upper hand against SGD empirically, the theoretical understanding of momentum in the stochastic case is far from complete, given that existing analyses of SGDM either "provide worse convergence bounds than those of SGD, or assume Lipschitz or quadratic objectives, which fail to hold in practice." [3] More specifically, the researchers demonstrate that SGDM has the same convergence bound as SGD for both strongly convex and nonconvex functions.

Huang et. al [4] propose a faster, universal framework for adaptive gradients that includes most already existing adaptive gradient forms of optimizers. Under the current framework of adaptive gradients, the adaptive learning rate varies between individual algorithms, e.g. Adam [7] vs. AdaGrad-Norm [8]. Not only does the proposed framework have a universal adaptive learning rate, it also guarantees convergence without assuming boundedness of the gradient. Furthermore, they experimentally validate the outperformance of their algorithm against existing adaptive algorithms, demonstrating the significance of this proposition as it is both more generalizable theoretically and empirically better in performance.

Faw et. al's recent work in SGD analysis [5] holds a particular significance in their demonstration that neither the bounded gradient assumption nor the bounded variance assumption are necessary in AdaGrad-Norm [8], thus adaptive gradient methods converge in much broader situations than previously understood.

References

- [1] Naichen Shi, Dawei Li, Mingyi Hong, and Ruoyu Sun. RMSprop converges with proper hyper-parameter. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=3UDSdyIcBDA>.
- [2] Robert M. Gower, Othmane Sebbouh, and Nicolas Loizou. Sgd for structured nonconvex functions: Learning rates, minibatching and interpolation, 2020. URL <https://arxiv.org/abs/2006.10311>.
- [3] Yanli Liu, Yuan Gao, and Wotao Yin. An improved analysis of stochastic gradient descent with momentum. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 18261–18271. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/d3f5d4de09ea19461dab00590df91e4f-Paper.pdf>.
- [4] Feihu Huang, Junyi Li, and Heng Huang. Super-adam: Faster and universal framework of adaptive gradients, 2021. URL <https://arxiv.org/abs/2106.08208>.
- [5] Matthew Faw, Isidoros Tziotis, Constantine Caramanis, Aryan Mokhtari, Sanjay Shakkottai, and Rachel Ward. The power of adaptivity in sgd: Self-tuning step sizes with unbounded gradients and affine variance, 2022. URL <https://arxiv.org/abs/2202.05791>.
- [6] Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *CoRR*, abs/1904.09237, 2019. URL <http://arxiv.org/abs/1904.09237>.
- [7] Anas Barakat and Pascal Bianchi. Convergence and dynamical behavior of the adam algorithm for nonconvex stochastic optimization. *SIAM Journal on Optimization*, 31(1):244–274, 2021. doi: 10.1137/19M1263443. URL <https://doi.org/10.1137/19M1263443>.
- [8] Rachel Ward, Xiaoxia Wu, and Leon Bottou. Adagrad stepsizes: Sharp convergence over nonconvex landscapes, 2018.