

Convergence Analysis of a Momentum Algorithm with Adaptive Step Size for Nonconvex Optimization

Anas Barakat, Pascal Bianchi

LTCI, Télécom Paris, Institut polytechnique de Paris

**11th OPT Workshop on Optimization for Machine Learning
December 14th, 2019**



A Momentum Algorithm with Adaptive Step Size

- ▶ ADAM famous **BUT** convergence issues (Reddi et al., 2018).
- ▶ Several variants : Yogi, AdaBound, AdaShift, Nadam, QHAdam, RAdam ...
- ▶ **Goal** : convergence rates for adaptive algorithms (ADAM in particular) for **nonconvex** optimization.

Algorithm

$$\begin{cases} x_{n+1} = x_n - a_{n+1} p_{n+1} \\ p_{n+1} = p_n + b (\nabla f(x_n) - p_n) \end{cases}$$

where $a_n \in \mathbb{R}_+^d$ $b \geq 0$, $x_0, p_0 \in \mathbb{R}^d$.

Contributions

Main Idea

Clipping the effective step size a_{n+1} :


$$0 < \delta \leq a_{n+1} \leq a_{sup}(L) \quad (1)$$

Results

- ▶ $O(1/n)$ convergence rate for ADAM in deterministic and stochastic settings.
(control of $\min_{0 \leq k \leq n-1} \|\nabla f(x_k)\|^2$).
- ▶ Convergence rate analysis of the objective function using the Kurdyka-Łojasiewicz (KL) property.

Thank you for your attention

Feel free to come to my poster




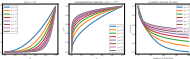
ANAS

Convergence Analysis of a Momentum Algorithm with Adaptive Step Size for Nonconvex Optimization

Anas Barakat and Pascal Bianchi

LTCI, Télécom Paris, Institut Polytechnique de Paris, France
anas.barakat@telecom-paristech.fr



| | | |
|---|--|---|
| Problem <ul style="list-style-type: none">$\min_{x \in \mathbb{R}^d} f(x)$$f$ non-convex differentiable.∇f is L-Lipschitz continuous.$\inf_{x \in \mathbb{R}^d} f(x) > -\infty$. | A descent lemma <p>$\forall x \in \mathbb{R}^d, \quad R_n := f(x_n) + \frac{1}{2n} \ x_n - \bar{x}\ ^2.$</p> <p>Lemma. Under previous assumptions, $\forall n \in \mathbb{N}, \forall x \in \mathbb{R}_n$,</p> $R_{n+1} \leq R_n - \frac{\alpha}{2n(n+1)} \ \bar{x}_{n+1} - \bar{x}_n\ ^2,$ <p>where $\bar{x}_{n+1} := 1 - \frac{\alpha_{n+1}}{2n}, \quad \frac{\alpha}{2n} = \frac{\alpha}{2n} \left(1 - \frac{\alpha}{2n} \right) \frac{1}{2n} = \frac{\alpha}{2n^2}.$</p> | KL inequality <ul style="list-style-type: none">satisfied by nonsmooth deep neural networks built from activations ReLU (Saxe et al., 2016), $t \mapsto t^2$ and log-exp ($\log(1 + e^x)$). <p>$\Phi_1 := \{x \in C^1([0, \eta]) \cap C^2([0, \eta]) : \Phi(0) = 0, \Phi \text{ concave and } \Phi' \geq 0\}.$</p> <p>Definition. (KL property [3, Appendix A]) A proper l.s.c. function $H : \mathbb{R}^d \rightarrow [-\infty, +\infty]$ has the KL property locally at $\bar{x} \in \text{dom } H$ if there exist $\alpha > 0, \rho \in \Phi_1$ and a neighborhood $\mathcal{U}(\bar{x})$ s.t. for all $x \in \mathcal{U}(\bar{x}) \cap \{H(x) \leq H(\bar{x}) + \alpha\}$:</p> $\ \nabla(\rho \circ H)(x) - H(x)\ \geq 1.$ <ul style="list-style-type: none">H becomes sharp under a reparameterization of its values through the so-called desingularizing function ρ.  |
| Summary <p>Main Idea : clipping the adaptive step size using a bound depending on $L(\nabla f)$.</p> <p>Contributions :</p> <ul style="list-style-type: none">sublinear rates in deterministic and stochastic context (no bounded gradients compared to [3], dimension free).convergence rates on the function value sequence under Kuratowski-Lojasiewicz (KL) property. | Deterministic setting <p>Theorem. Let previous assumptions hold. Assume $1 - \alpha < \beta \leq 1$ and :</p> <ul style="list-style-type: none">Let $t > 0$ s.t. $\alpha_{\text{step}} := \frac{\beta}{2} \left(1 - \frac{\alpha_{\text{step}}}{2n} - \log(1 - \frac{\alpha_{\text{step}}}{2n}) \right) \geq 0$.Let $t > 0$ s.t. $\forall n \in \mathbb{N}, \quad \beta \leq \alpha_{n+1} \leq \min(\alpha_{\text{step}}, \frac{\beta}{2n})$. <p>Then (R_n) is nonincreasing, $\lim \nabla f(x_n) \rightarrow 0$ as $n \rightarrow +\infty$ and</p> $\forall n \geq 1, \quad \min_{x \in \mathbb{R}_n} \ \nabla f(x_n)\ ^2 \leq \frac{4}{n^2} \left(\frac{H_n - \inf f}{\beta} + \beta n^2 \right).$ | KL rates (similar techniques to [3, 4]) <p>$\forall x \in \mathbb{R}^d \times \mathbb{R}^d, \quad H(x) := H(x, y) = f(x) + \frac{1}{2\beta} \ y\ ^2.$</p> <p>Theorem. Let $x_0 = (x_0, y_0)$ where $y_0 = \sqrt{2\beta} \nabla f(x_0)$, $f(x_0) = \inf f(x_0)$ where $\nabla f(x_0) = 0$. Suppose that f is concave, condition (1) holds and</p> <ul style="list-style-type: none">H is a KL function with KL exponent θ i.e. $\varphi(x) = \frac{\beta}{2} e^{\theta x}, \theta \in (0, 1]$. <p>(i) If $\theta = 1$, then $f(x_n)$ converges in a finite number of iterations.</p> <p>(ii) If $1/2 \leq \theta < 1$, then $\forall \eta \in (0, 1), C > 0$ s.t. $f(x_n) - f(x_0) \leq C \eta^C$.</p> <p>(iii) If $\theta < 1/2$, then $f(x_n) - f(x_0) = O(n^{-\frac{1}{1-\theta}})$.</p> |
| A momentum algorithm <p>$\begin{cases} x_{n+1} = x_n - \alpha_{n+1} \nabla f(x_n) \\ \bar{x}_{n+1} = \bar{x}_n + b(\nabla f(x_n) - \bar{x}_n) \end{cases}$</p> <ul style="list-style-type: none">constructive product.$\alpha_n \in \mathbb{R}^+$ may depend on the past gradients $g_n = \nabla f(x_n)$ and the iterates x_k for $k \leq n$.includes SGD, Heavy Ball, Antzoul and other adaptive algorithms [2]. | Stochastic setting <p>Theorem. Let previous assumptions hold. Assume $1 - \alpha < \beta \leq 1$ and :</p> <ul style="list-style-type: none">$\forall x \in \mathbb{R}^d, \quad \mathbb{E}[\ \nabla f(x, \zeta) - \nabla f(x)\ ^2] \leq \sigma^2.$Let $t > 0$ s.t. $\alpha_{\text{step}} := \frac{\beta}{2} \left(1 - \frac{\alpha_{\text{step}}}{2n} - \log(1 - \frac{\alpha_{\text{step}}}{2n}) \right) \geq 0$.Let $t > 0$ s.t. $\forall n \geq 1$, almost surely, $\beta \leq \alpha_{n+1} \leq \min(\alpha_{\text{step}}, \frac{\beta}{2n})$. $\mathbb{E}[\ \nabla f(x_n, \zeta)\ ^2] \leq \frac{4}{n^2} \left(\frac{H_n - \inf f}{\beta} + \frac{\sigma^2}{2n} \right)$ <p>where x_n is an iterate uniformly randomly chosen from $\{x_0, \dots, x_{n-1}\}$.</p> | References <p>[1] M. S. Bubeck, S. Bubeck, S. Bubeck, and S. Bubeck. Random algorithms for nonconvex optimization. In <i>Advances in Neural Information Processing Systems</i>, pages 9791–9801, 2019.</p> <p>[2] M. S. Bubeck, S. Bubeck, S. Bubeck, and S. Bubeck. Adaptive gradient algorithms for nonconvex optimization. In <i>ICML</i>, pages 1040–1050, 2019.</p> <p>[3] M. S. Bubeck, S. Bubeck, S. Bubeck, and S. Bubeck. Convergence rates of forward gradient schemes for nonconvex optimization. In <i>ICML</i>, pages 1040–1050, 2019.</p> <p>[4] M. S. Bubeck, S. Bubeck, S. Bubeck, and S. Bubeck. Convergence rates of forward gradient schemes for nonconvex optimization. In <i>ICML</i>, pages 1040–1050, 2019.</p> |
| MIM Assumption (verified for ANAS) <p>There exists $\alpha > 0$ s.t. $\alpha_{n+1} \leq \frac{\alpha}{2n}$.</p> | | |

For more details: article available on the Workshop page / arXiv.
AB, P. Bianchi. *Convergence Analysis of a Momentum Algorithm with Adaptive Step Size for Nonconvex Optimization*