

Progress Report 2: A Study on Convergence Results of Stochastic Gradient Methods

B09902055 Weiping Li, B09902073 Chun-Neng Chu

2023.5.16

Outline

1. Past Progress
 - a. Difference between Algorithms
 - b. Constraint Tradeoff
2. Goals
3. Standard SGD Convergence Proof Format
4. Insights from ADAGRAD-Norm (Xiaoyu Li et al.)
 - a. Non-arithmetic Lemmas
 - b. Necessity of Changing Update Sequence - Lemma 3
 - c. Necessity of Smoothness Constant - Lemma 3 & 8
 - d. Comparison with ADAGRAD-Norm (Rachel Ward et al.)
5. Future Plans
6. QA
7. Appendix
8. References

Difference between Algorithms: Update Sequence

Algorithm 1 ADAGRAD-Norm (Xiaoyu Li et al.)

- a. 1: Input: Initialize $x_0 \in R^d, b_0 > 0, \eta > 0$
2: **for** $t = 1, 2, \dots$ **do**
3: Generate $\xi_{t-1}, G_{t-1} = G(x_{t-1}, \xi_{t-1})$
4: $x_t \leftarrow x_{t-1} - \frac{\eta}{b_{t-1}} G_{t-1}$
5: $b_t^2 \leftarrow b_{t-1}^2 + \|G_{t-1}\|^2$
6: **end for**
-

Algorithm 2 ADAGRAD-Norm (Rachel Ward et al.)

- b. 1: Input: Initialize $x_0 \in R^d, b_0 > 0, \eta > 0$
2: **for** $t = 1, 2, \dots$ **do**
3: Generate $\xi_{t-1}, G_{t-1} = G(x_{t-1}, \xi_{t-1})$
4: $b_t^2 \leftarrow b_{t-1}^2 + \|G_{t-1}\|^2$
5: $x_t \leftarrow x_{t-1} - \frac{\eta}{b_t} G_{t-1}$
6: **end for**
-

Constraint Tradeoff

The two algorithms, with their individual constraints, can be proven to have the same complexity for convergence with respect to iteration T : $O(\frac{1}{\sqrt{T}})$ However, their constraints differ.

Constraints		
Constraints↓ Algorithm →	Xiaoyu Li et al	Rachel Ward et al.
M-smooth	✓	✓
$E \left[\ \nabla f(x_t) - G(x_t, \xi_t)\ ^2 \right] \leq \sigma^2$	✓	✓
$E_{\xi}[G(x, \xi)] = \nabla f(x)$	✓	✓
$f > -\infty$	✓	✓
know smoothness constant M	✓	
L-Lipschitz		✓

Goals

1. Find the reason for the different update sequence between the two algorithms.
2. Understand why Xiaoyu Li et al. require prior knowledge of specific smoothness constant in the proof.
3. Understand why Rachel Ward et al. require Lipschitz constraint in proof.
4. Can we acquire the same convergence rate with Rachel Ward et al.'s algorithm, but using the constraints of Xiaoyu Li et al.? Vice versa?

Standard SGD Convergence Proof Format

Intuition: We want to find a complexity bound for $\|\nabla f(x_t)\|^a$ in the form of $O(\frac{1}{T^\alpha})$. Hence, we use the following thought process.

① $\xleftarrow{\text{Markov's}}$ ② \leftarrow ③ $\xleftarrow{\text{Trick}}$ ④ $\xleftarrow{\text{Trick}}$ ⑤

$$\textcircled{1} \quad P(\min_{1 \leq t \leq T} \|\nabla f(x_t)\|^2 = O(\frac{1}{T^\alpha})) \geq 1 - \delta$$

$$\textcircled{2} \quad E[\min_{1 \leq t \leq T} \|\nabla f(x_t)\|^a]^{\frac{2}{a}} = O(\frac{1}{T^\alpha})$$

$$\textcircled{3} \quad E[(\sum_{t=1}^T \|\nabla f(x_t)\|^2)^{\frac{a}{2}}] = O(\frac{1}{T^{\alpha - \frac{a}{2}}})$$

$$\textcircled{4} \quad \sum_{t=1}^T \eta_t^* E[\|\nabla f(x_{t-1})\|^2] \leq \\ f(x_0) - f^* + \sum_{t=1}^T \frac{\eta_t^2 \|G(x_{t-1}, \xi_{t-1})\|^2 M}{2}$$

$$\textcircled{5} \quad |f(x_t) - f(x_{t-1}) - \langle \nabla f(x_{t-1}), x_t - x_{t-1} \rangle| \leq \frac{M}{2} \|x_t - x_{t-1}\|^2$$

Note: η_t is the learning rate at the t^{th} iteration; η_t^* may be η_t or the estimation of η_t ; f^* is the optimal target function value.

Xiaoyu Li et al. - Non-arithmetic Segments: Lemma 3

When surveying Xiaoyu Li et al.'s article, we noticed that only two segments in the proof require non-arithmetic lemmas (3, 8).

Lemma 3: Assume f is M -smooth and

$E[G(x_{t-1}, \xi_{t-1})] = \nabla f(x_{t-1})$. Then, the iterates of SGD with stepsizes $\eta_t \in R^d$ satisfy the following inequality

$$E \left[\sum_{t=1}^T \langle \nabla f(x_{t-1}), \eta_t \nabla f(x_{t-1}) \rangle \right] \leq f(x_{t-1}) - f^* \\ + \frac{M}{2} E \left[\sum_{t=1}^T \|\eta_t G(x_{t-1}, \xi_{t-1})\|^2 \right]$$

Xiaoyu Li et al. - Non-arithmetic Segments: Lemma 8

Lemma 8: Assume f is M -smooth, $E[G(x_{t-1}, \xi_{t-1})] = \nabla f(x_{t-1})$ and the stochastic gradient satisfies

$E[\exp(\|\nabla f(x) - g(x, \xi)\|^2/\sigma^2)] \leq \exp(1), \forall x$. Then,

$$E \left[\sum_{t=1}^T \eta_t^2 \|G(x_{t-1}, \xi_{t-1})\|^2 \right] \leq K + \frac{4\eta^2}{b_0^2} (1 + \ln T) \sigma^2 \\ + \frac{4\eta}{b_0} E \left[\sum_{t=1}^T \eta_t \|\nabla f(x_{t-1})\|^2 \right],$$

where

$$K = 2\eta^2 \ln \left(\sqrt{b_0^2 + 2T\sigma^2} + \sqrt{2} E \left[\sqrt{\sum_{t=1}^T \|\nabla f(x_{t-1})\|^2} \right] \right) - 2\eta^2 \ln(b_0)$$

Xiaoyu Li et al. - Changing Update Sequence for Lemma 3

In the proof of Lemma 3, there is an intermediary step that requires the following:

$$E_{\xi_{t-1}} [\langle \nabla f(x_{t-1}), \eta_t (\nabla f(x_{t-1}) - G(x_{t-1}, \xi_{t-1})) \rangle] = \\ \langle \nabla f(x_{t-1}), \eta_t \nabla f(x_{t-1}) - \eta_t E_t [G(x_{t-1}, \xi_{t-1})] \rangle = 0$$

This equation requires that η_t is independent to ξ_{t-1} . The two terms are independent due to the fact that at the t^{th} iteration, η_t is decided by ξ_0 to ξ_{t-2} . Hence, η_t can be taken out of the expectation.

Xiaoyu Li et al. - Smoothness Constant from Lemma 3 & 8

In the the next three slides, we demonstrate why concrete knowledge on the value of smoothness constant M is necessary.

$$\begin{aligned} & E \left[\sum_{t=1}^T \eta_t^2 \|G(x_{t-1}, \xi_{t-1})\|^2 \right] \\ &= E \left[\sum_{t=1}^T \eta_{t+1}^2 \|G(x_{t-1}, \xi_{t-1})\|^2 + \sum_{t=1}^T \|G(x_{t-1}, \xi_{t-1})\|^2 (\eta_t^2 - \eta_{t+1}^2) \right] \\ &\leq 2\eta^2 \ln \left(\sqrt{b_0^2 + 2T\sigma^2} + \sqrt{2} E \left[\sqrt{\sum_{t=1}^T \|\nabla f(x_{t-1})\|^2} \right] \right) - 2\eta^2 \ln(b_0) \\ &\quad + \frac{4\eta^2}{b_0^2} (1 + \ln T) \sigma^2 + \frac{4\eta}{b_0} E \left[\sum_{t=1}^T \eta_t \|\nabla f(x_{t-1})\|^2 \right] \end{aligned}$$

Omitting the majority of tricks used by Xiaoyu Li et al., we can claim that the term $\sqrt{2}E \left[\sqrt{\sum_{t=1}^T \|\nabla f(x_{t-1})\|^2} \right]$ can be dropped,

and $\sum_{t=1}^T \eta_{t+1}^2 \|G(x_{t-1}, \xi_{t-1})\|^2$ is bounded by $O(\ln(\sqrt{T}))$.

Intuitively speaking, the term $\sum_{t=1}^T \|G(x_{t-1}, \xi_{t-1})\|^2 (\eta_t^2 - \eta_{t+1}^2)$ is the penalty caused by "borrowing" the red term, which is from the next iteration, to bound the stochastic gradient norm squared.

With some tricks, it can be bounded by

$$\frac{4\eta^2}{b_0^2} (1 + \ln T) \sigma^2 + \frac{4\eta}{b_0} E \left[\sum_{t=1}^T \eta_t \|\nabla f(x_{t-1})\|^2 \right].$$

When lemma 3 is applied in an attempt to bound the term

$E \left[\sum_{t=1}^T \eta_t \|\nabla f(x_{t-1})\|^2 \right]$, scaling the same term on the RHS of lemma 8 by $\frac{M}{2}$, we find the inequality on the next page.

$$\begin{aligned}
& \left(1 - \frac{2\eta M}{\sqrt{b_0^2}}\right) E \left[\sum_{t=1}^T \eta_t \|\nabla f(x_{t-1})\|^2 \right] \leq f(x_0) - f^* \\
& + M \left(\eta^2 \ln \left(\sqrt{b_0^2 + 2T\sigma^2} + \sqrt{2} E \left[\sqrt{\sum_{t=1}^T \|\nabla f(x_{t-1})\|^2} \right] \right) - \frac{\eta^2 \ln(b_0^2)}{2} \right) \\
& + \frac{2\eta M}{b_0^2} (1 + \ln T) \sigma^2.
\end{aligned}$$

Consider the case where $\left(1 - \frac{2\eta M}{\sqrt{b_0^2}}\right) \leq 0$, we can see that the inequality always holds, and thus we gain no information of the term $E \left[\sum_{t=1}^T \eta_t \|\nabla f(x_{t-1})\|^2 \right]$. Hence we need to know the constant M to initialize η and b_0 .

What Happens when Xiaoyu Li et al.'s Constraints Applied on Rachel Ward et al.

1. Since for the t^{th} iteration, Rachel Ward et al. updates the learning rate (η_t) before the weights, we cannot take η_t out of the expectation in the intermediary step that requires
$$E_{\xi_{t-1}} [\langle \nabla f(x_{t-1}), \eta_t (\nabla f(x_{t-1}) - G(x_{t-1}, \xi_{t-1})) \rangle] = \langle \nabla f(x_{t-1}), \eta_t \nabla f(x_{t-1}) - \eta_t E_t [G(x_{t-1}, \xi_{t-1})] \rangle = 0$$
2. In Rachel Ward et al.'s article, η_t^* is an estimation of η_t instead of exactly η_t . Since theorem 4 in Xiaoyu Li et al. requires descending η_t^* and it is hard to confirm whether the estimation is indeed descending, the proof cannot be generalized to Rachel Ward et al. trivially.

Future Plans

1. Find the reason for the different update sequence between the two algorithms.
2. Understand why Xiaoyu Li et al. require prior knowledge of specific smoothness constant in the proof.
3. Understand why Rachel Ward et al. require Lipschitz constraint in proof.
4. Can we acquire the same convergence rate with Rachel Ward et al.'s algorithm, but using the constraints of Xiaoyu Li et al.? Vice versa?

Appendix: Proof of Lemma 3

From the definition of M -smooth, we have

$$|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq \frac{M}{2} \|y - x\|^2 \quad \text{Thus,}$$

$$\begin{aligned} f(x_t) &\leq f(x_{t-1}) + \langle \nabla f(x_{t-1}), x_t - x_{t-1} \rangle + \frac{M}{2} \|x_t - x_{t-1}\|^2 \\ &= f(x_{t-1}) + \langle \nabla f(x_{t-1}), \eta_t (\nabla f(x_{t-1}) - G(x_{t-1}, \xi_{t-1})) \rangle \\ &\quad - \langle \nabla f(x_{t-1}), \eta_t \nabla f(x_{t-1}) \rangle + \frac{M}{2} \|\eta_t G(x_{t-1}, \xi_{t-1})\|^2. \end{aligned}$$

Taking the conditional expectation with respect to ξ_0, \dots, ξ_{t-2} , we have that $E_{t-1}[\langle \nabla f(x_{t-1}), \eta_t (\nabla f(x_{t-1}) - G(x_{t-1}, \xi_{t-1})) \rangle] = \langle \nabla f(x_{t-1}), \eta_t \nabla f(x_{t-1}) - \eta_t E_t[G(x_{t-1}, \xi_{t-1})] \rangle = 0$.

Appendix: Proof of Lemma 3

Hence, from the law of total expectation, we have

$$E[\langle \nabla f(x_{t-1}), \eta_t \nabla f(x_{t-1}) \rangle] \leq$$
$$E\left[f(x_{t-1}) - f(x_t) + \frac{M}{2} \|\eta_t g(x_{t-1}, \xi_{t-1})\|^2\right].$$

Summing over $t = 1$ to T and lower bounding $f(x_T)$ with f^* , we have the stated bound.

Appendix: Proof of Lemma 8

Lemma 10 states: If $x > 0, \eta > 0$, then $\ln(\frac{1}{x}) \geq \eta \left(1 - x^{\frac{1}{\eta}}\right)$.

$$\begin{aligned} E \left[\sum_{t=1}^T \eta_{t+1}^2 \|G(x_{t-1}, \xi_t)\|^2 \right] &= E \left[\sum_{t=1}^T \frac{\eta^2 \|G(x_{t-1}, \xi_{t-1})\|^2}{\left(b_0^2 + \sum_{i=1}^t \|g(x_{i-1}, \xi_{i-1})\|^2\right)} \right] \\ &\leq 2\eta^2 E \left[\ln \left(\sqrt{b_0^2 + \sum_{t=1}^T \|g(x_{t-1}, \xi_{t-1})\|^2} \right) \right] \\ &\leq 2\eta^2 \ln \left(\sqrt{b_0^2 + 2T\sigma^2} + \sqrt{2} E \left[\sqrt{\sum_{t=1}^T \|\nabla f(x_{t-1})\|^2} \right] \right) \end{aligned}$$

Appendix: Proof of Lemma 8

Where in first inequality we used Lemma 10 and in the third one we used Jensen's inequality. Putting things together, we have

$$\begin{aligned} & E \left[\sum_{t=1}^T \eta_t^2 \|G(x_{t-1}, \xi_{t-1})\|^2 \right] \\ &= E \left[\sum_{t=1}^T \eta_{t+1}^2 \|G(x_{t-1}, \xi_{t-1})\|^2 + \sum_{t=1}^T \|G(x_{t-1}, \xi_{t-1})\|^2 (\eta_t^2 - \eta_{t+1}^2) \right] \\ &\leq 2\eta^2 \ln \left(\sqrt{b_0^2 + 2T\sigma^2} + \sqrt{2} E \left[\sqrt{\sum_{t=1}^T \|\nabla f(x_{t-1})\|^2} \right] \right) \\ &\quad + \frac{4\eta^2}{b_0^2} (1 + \ln T) \sigma^2 + \frac{4\eta}{b_0^{2\frac{1}{2}}} E \left[\sum_{t=1}^T \eta_t \|\nabla f(x_{t-1})\|^2 \right] \end{aligned}$$

Appendix

1. L-Lipschitz gradient implies L-smooth:

- (a) f is L-smooth if f is continuously differentiable and
$$\forall x, y \in \text{dom } f, \langle \nabla f(y) - \nabla f(x), y - x \rangle \leq L \|y - x\|^2$$
- (b) ∇f is L-Lipschitz iff
$$\forall x, y \in \text{dom } f, \|\nabla f(y) - \nabla f(x)\| \leq L \|y - x\|$$
- (c) Cauchy-Schwarz inequality states that $\forall u, v \in$ an inner product space, $|\langle u, v \rangle| \leq \|u\| \|v\|$.

By (a), (b), (c) can write

$$\langle \nabla f(y) - \nabla f(x), y - x \rangle \leq \|\nabla f(y) - \nabla f(x)\| \|y - x\| \leq L \|y - x\|^2$$

2. Convergence rate - expectation form to probability form:

Markov's inequality: $P(X \geq a) \leq \frac{E[X]}{a}$, if X is a non-negative random variable and $a > 0$

$$\rightarrow P(\min_{0 \leq t \leq N} \|\nabla f(x_t)\|^2 \geq \frac{E[(\min_{0 \leq t \leq N} \|\nabla f(x_t)\|^2)]}{\delta}) \leq \delta$$

$$\rightarrow P(\min_{0 \leq t \leq N} \|\nabla f(x_t)\|^2 \leq \frac{E[(\min_{0 \leq t \leq N} \|\nabla f(x_t)\|^2)]}{\delta}) \geq 1 - \delta$$

$$\rightarrow H(\delta, N, D^*) = \frac{E[(\min_{0 \leq t \leq N} \|\nabla f(x_t)\|^2)]}{\delta}$$

References

1. Quoc Tran-Dinh. Sublinear Convergence Rates of Extragradient-Type Methods: A Survey on Classical and Recent Developments.
2. Xiaoyu Li, Francesco Orabona. On the Convergence of Stochastic Gradient Descent. with Adaptive Stepsizes
3. Rachel Ward, Xiaoxia Wu, Léon Bottou. AdaGrad stepsizes: Sharp convergence over nonconvex landscapes.
4. Yen-Huan Li. Optimization algorithms Lecture 3.
5. Chih-Jen Lin. Optimization Methods for Deep Learning: Convergence of stochastic gradient methods.