

Finite-Sum Optimization: Adaptivity to Smoothness and Loopless Variance Reduction

Bastien Batardière, Joon Kwon

February 5, 2024

Context

- Large scale optimization
 - Number of samples $n \gg 1$
 - Number of parameters $d \gg 1$
- Examples
 - ▶ scRNA seq data where $n \approx d \approx 20000$
 - ▶ Any log-likelihood with lots of samples.

Problem

- ▶ Goal:

$$\arg \min_{\theta \in \mathbb{R}^d} f(\theta)$$

where f has a finite-sum structure:

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

- ▶ Example of loss functions:

1. $f_i(\theta) = -\log P_\theta(X_i = x_i)$
2. Variational ELBOs can be designed this way:

$$J(\theta, \phi | Y) = \frac{1}{n} \sum_{i=1}^n J_i(\theta, \phi_i | Y_i)$$

with Y the data, ϕ the variational parameters, θ the model parameters.

Outline

Assumptions

Stochastic Gradient Descent

SAGA

AdaGrad

AdaSAGA

Convergence results

Numerical experiments

Assumptions

1. Each f_i is convex and differentiable
2. Each f_i is L -smooth:

$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|, \quad x, y \in \mathbb{R}^d, \quad 1 \leq i \leq n.$$

Examples:

- ▶ Linear regression $f_i(\theta) = \|x_i^\top \theta - y_i\|^2, (x_i, y_i)_{1 \leq i \leq n}$
- ▶ Logistic regression

$$f_i(\theta) = \log(\mathbb{P}_\theta(X_i = x_i))y_i + \log(1 - \mathbb{P}_\theta(X_i = x_i))(1 - y_i)$$

Goal: reach an $\epsilon > 0$ solution: find θ such that

$$f(\theta) - f(\theta^*) \leq \epsilon$$

Coefficient L is not always known.

Stochastic Gradient Descent [Robbins, 1951]

- ▶ Stochastic GD (SGD): replace $\nabla f(\theta_t)$ of classic Gradient Descent with an unbiased low cost gradient:

$$i_t \sim \text{unif}(1 \cdots n)$$
$$\theta_{t+1} = \theta_t - \eta_t \nabla f_{i_t}(\theta_t)$$

- ▶ Unbiased gradient estimator:

$$\mathbb{E}[\nabla f_{i_t}(\theta_t)] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\theta_t) = \nabla f(\theta_t)$$

- ▶ If $\eta_t = \frac{1}{\sqrt{t}}$, complexity in $O(\frac{1}{\epsilon^2})$ iteration

Stochastic Averaging GrAdient (SAGA)

[Defazio et al., 2014]

- ▶ Store the previously computed gradient in a matrix $\phi^{(t)} \in \mathbb{R}^{n \times d}$ to estimate the gradient of the whole dataset.

We aim to have

$$\frac{1}{n} \sum_{i=1}^n \phi_i^{(t)} \approx \nabla f(\theta_t)$$

- ▶ Given $\phi^{(1)} \in \mathbb{R}^{n \times d}, \theta_1$, perform

$$i_t \sim \text{unif}(1 \cdots n) \tag{1}$$

$$\mathbf{g}_t = \nabla f_{i_t}(\theta_t) - \phi_{i_t}^{(t)} + \frac{1}{n} \sum_{i=1}^n \phi_i^{(t)} \tag{2}$$

$$\theta_{t+1} = \theta_t - \eta \mathbf{g}_t \tag{3}$$

$$\text{Update } \phi_i^{(t+1)} = \begin{cases} \phi_i^{(t)} & \text{if } i \neq i_t \\ \nabla f_{i_t}(\theta_t) & \text{if } i = i_t \end{cases} \tag{4}$$

Update of the matrix ϕ , when $n = 3, d = 1$

$$(\phi_0 \quad \phi_0 \quad \phi_0)$$

Update of the matrix ϕ , when $n = 3, d = 1$

$$\begin{array}{l} (\phi_0 \quad \phi_0 \quad \phi_0) \\ t = 1, i(t) = 2 \implies \nabla f_2(\theta_1) \text{ is computed} \\ \Downarrow \end{array}$$

Update of the matrix ϕ , when $n = 3, d = 1$

$$\begin{array}{c} (\phi_0 \quad \phi_0 \quad \phi_0) \\ t = 1, i(t) = 2 \implies \nabla f_2(\theta_1) \text{ is computed} \\ \Downarrow \\ (\phi_0 \quad \nabla f_2(\theta_1) \quad \phi_0) \end{array}$$

Update of the matrix ϕ , when $n = 3, d = 1$

$$\begin{array}{l} (\phi_0 \quad \phi_0 \quad \phi_0) \\ t = 1, i(t) = 2 \implies \nabla f_2(\theta_1) \text{ is computed} \\ \Downarrow \\ (\phi_0 \quad \nabla f_2(\theta_1) \quad \phi_0) \\ t = 2, i(t) = 3 \implies \nabla f_3(\theta_2) \text{ is computed} \\ \Downarrow \end{array}$$

Update of the matrix ϕ , when $n = 3, d = 1$

$$\begin{array}{c} (\phi_0 \quad \phi_0 \quad \phi_0) \\ t = 1, i(t) = 2 \implies \nabla f_2(\theta_1) \text{ is computed} \end{array}$$

\Downarrow

$$\begin{array}{c} (\phi_0 \quad \nabla f_2(\theta_1) \quad \phi_0) \\ t = 2, i(t) = 3 \implies \nabla f_3(\theta_2) \text{ is computed} \end{array}$$

\Downarrow

$$(\phi_0 \quad \nabla f_2(\theta_1) \quad \nabla f_3(\theta_2))$$

Update of the matrix ϕ , when $n = 3, d = 1$

$$\begin{array}{l} (\phi_0 \quad \phi_0 \quad \phi_0) \\ t = 1, i(t) = 2 \implies \nabla f_2(\theta_1) \text{ is computed} \\ \Downarrow \\ (\phi_0 \quad \nabla f_2(\theta_1) \quad \phi_0) \\ t = 2, i(t) = 3 \implies \nabla f_3(\theta_2) \text{ is computed} \\ \Downarrow \\ (\phi_0 \quad \nabla f_2(\theta_1) \quad \nabla f_3(\theta_2)) \\ t = 3, i(t) = 3 \implies \nabla f_3(\theta_3) \\ \Downarrow \end{array}$$

Update of the matrix ϕ , when $n = 3, d = 1$

$$\begin{array}{c} (\phi_0 \quad \phi_0 \quad \phi_0) \\ t = 1, i(t) = 2 \implies \nabla f_2(\theta_1) \text{ is computed} \end{array}$$

\Downarrow

$$\begin{array}{c} (\phi_0 \quad \nabla f_2(\theta_1) \quad \phi_0) \\ t = 2, i(t) = 3 \implies \nabla f_3(\theta_2) \text{ is computed} \end{array}$$

\Downarrow

$$\begin{array}{c} (\phi_0 \quad \nabla f_2(\theta_1) \quad \nabla f_3(\theta_2)) \\ t = 3, i(t) = 3 \implies \nabla f_3(\theta_3) \end{array}$$

\Downarrow

$$(\phi_0 \quad \nabla f_2(\theta_1) \quad \nabla f_3(\theta_3))$$

Adagrad [Duchi et al., 2011]

- ▶ Given $\theta_1, G_0 = 0_{\mathbb{R}^d}$, perform

$$i_t \sim \text{unif}(1 \cdots n)$$

$$g_t = \nabla f_{i_t}(\theta_t)$$

$$G_t = G_{t-1} + g_t \odot g_t$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t}} \odot g_t$$

Goal: adapting the learning rate according to the data.

AdaSAGA

Goal: use both variance reduction from SAGA and adaptation from Adagrad.

- ▶ Given $\phi^{(1)} \in \mathbb{R}^{n \times d}$, θ_1 , $G_0 = 0_{\mathbb{R}^d}$, perform

$$i_t \sim \text{unif}(1 \cdots n)$$

$$g_t = \nabla f_{i_t}(\theta_t) - \phi_{i_t}^{(t)} + \frac{1}{n} \sum_{i=1}^n \phi_i^{(t)}$$

$$G_t = G_{t-1} + g_t \odot g_t$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t}} \odot g_t$$

$$\text{Update } \phi_i^{(t+1)} = \begin{cases} \phi_i^{(t)} & \text{if } i \neq i_t \\ \nabla f_{i_t}(\theta_t) & \text{if } i = i_t \end{cases}$$

Convergence results

Theorem

Let $T \geq 1, \eta > 0, x^{(1)} \in \mathbb{R}^d$ and $(x^{(t)})_{1 \leq t \leq T}$ be a sequence of iterates defined by AdaSAGA. Then, under convexity and L -smoothness, it holds that:

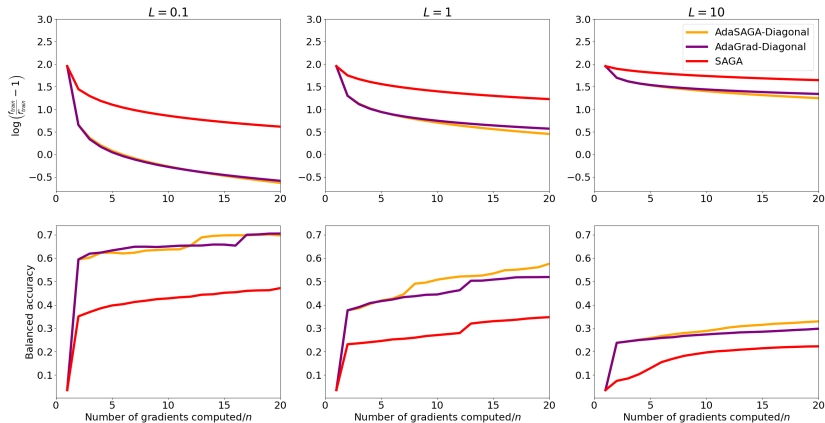
$$\mathbb{E}[f(\bar{x}_T)] - f(x^*) \leq \frac{\alpha \left(\eta + \frac{D^2}{2\eta} \right) \sqrt{4Ln\Delta_1} + 8L\alpha^2 \left(\eta + \frac{D^2}{2\eta} \right)^2}{T},$$

where $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x^{(t)}$, $\Delta_1 = f(x^{(1)}) - f(x^*)$ and $\alpha = 1$ (resp. $\alpha = \sqrt{d}$) for the Norm variant (resp. the Diagonal variant).

Complexity summary

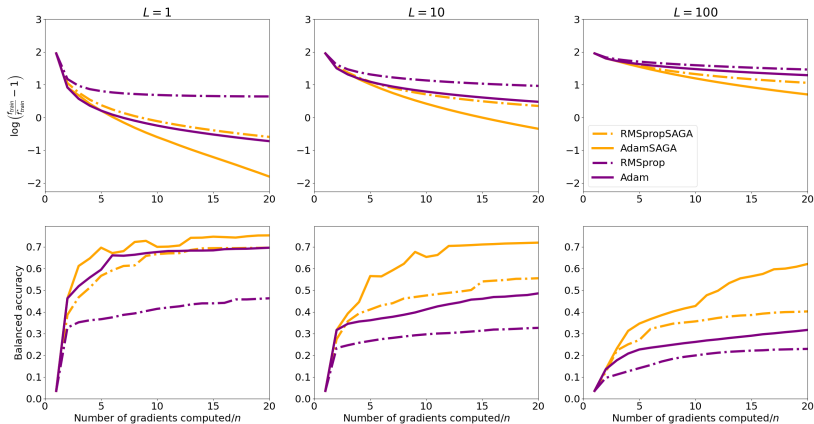
	Gradient complexity to reach an ϵ -solution
<i>GD</i>	$O\left(\frac{nL}{\epsilon}\right)$
<i>SGD</i>	$O\left(\frac{\sigma}{\epsilon^2}\right)$
<i>SAGA</i>	$O\left(\frac{n+L}{\epsilon}\right)$
<i>Adagrad</i>	$O\left(\frac{\sigma}{\epsilon^2}\right)$
<i>AdaSAGA</i>	$O\left(\frac{L+\sqrt{Ln}}{\epsilon}\right)$

Numerical experiments



Logistic regression on scMark(20000×8400). Comparison between AdaGrad, SAGA and AdaSAGA.

Numerical experiments



Comparison between Adam, RMSprop, AdamSAGA and RMSPropSAGA on SCMark dataset ($n = 20000$, $d = 8400$).


Conclusions


- ▶ Faster convergence than AdaGrad
- ▶ Adaptivity to smoothness
- ▶ More effective when replacing AdaGrad with Adam


Perspectives


- ▶ Proof when combine with Adam
- ▶ Combining with acceleration methods (e.g. Katyusha, Varag, Vrada)
- ▶ Non-convex objectives


Thank you for your attention


 Allen-Zhu, Z. (2018).
Katyusha: The first direct acceleration of stochastic gradient methods.

 Defazio, A., Bach, F., and Lacoste-Julien, S. (2014).
Saga: A fast incremental gradient method with support for non-strongly convex composite objectives.

 Dubois-Taine, B., Vaswani, S., Babanezhad, R., Schmidt, M., and Lacoste-Julien, S. (2021).
Svrg meets adagrad: Painless variance reduction.

 Duchi, J., Hazan, E., and Singer, Y. (2011).
Adaptive subgradient methods for online learning and stochastic optimization.
J. Mach. Learn. Res., 12(null):2121–2159.

 Kingma, D. P. and Ba, J. (2014).
Adam: A method for stochastic optimization.
ICML.

 Lan, G., Li, Z., and Zhou, Y. (2019).
A unified variance-reduced accelerated gradient method for