

Stochastic Approximation Beyond Gradient

Gersende Fort
CNRS

Institut de Mathématiques de Toulouse, France



European Meeting of Statisticians 2023, Warsaw, Poland

In collaboration with

- Aymeric Dieuleveut, Ecole Polytechnique, CMAP, France
- Eric Moulines, Ecole Polytechnique, CMAP, France
- Hoi To Wai, Chinese Univ. of Hong-Kong, Hong-Kong

Talk based on the paper:

- *Stochastic Approximation Beyond Gradient for Signal Processing and Machine Learning*

by A. Dieuleveut, G. Fort, E. Moulines and H.-T. Wai
HAL-03979922 and arXiv:2302.11147

Partly funded by

Fondation Simone et Cino Del Duca, Project OpSiMorE



Outline

- 1 Stochastic Approximation, beyond gradient
- 2 Two examples
- 3 Non asymptotic convergence bounds in expectation
- 4 Variance Reduction by the SPIDER control variate

I. Stochastic Approximation

Stochastic Approximation: solve the root-finding problem

Solve

$$\omega \in \mathbb{R}^d \quad \text{s.t.} \quad h(\omega) = 0$$

when **only stochastic estimates** of the *mean field* h are available.

by an iterative algorithm

$$\omega_{k+1} = \omega_k + \gamma_{k+1} H(\omega_k, X_{k+1})$$

where

- γ_{k+1} is a positive step size
- $H(\omega_k, X_{k+1})$ is a stochastic oracle for $h(\omega_k)$.

Time discretization of the ODE

$$\frac{d\omega}{dt} = h(\omega(t))$$

yields $\tau_{k+1} = \tau_k + \gamma_{k+1} h(\tau_k)$.

Beyond the gradient case

The gradient case:

- Solve $\operatorname{argmin}_{\omega} f(\omega)$ “by” solving $\nabla f = 0$ when $-\nabla f(\omega) = \mathbb{E}[H(\omega, X)]$

Expected Risk Minimization	for batch data	$f(\omega) = (1/n) \sum_{i=1}^n \ell(\omega, Z_i)$
	for streaming data	$f(\omega) = \mathbb{E}[\ell(\omega, X)]$

- Available oracles given ω_k , a random variable X_{k+1} and the stochastic gradient term $H(\omega_k, X_{k+1})$.

Expected Risk Minimization	for batch data	$H(\omega, X_{k+1}) = -\mathbb{D}_{10} \ell(\omega, Z_{X_{k+1}})$	$X_{k+1} \in \{1, \dots, n\}$
	for streaming data	$H(\omega, X_{k+1}) = -\mathbb{D}_{10} \ell(\omega, X_{k+1})$	

Two extensions:

- The function h is not necessarily a gradient
- The oracle can be **biased**

$$\mathbb{E} \left[H(\omega_k, X_{k+1}) \middle| \mathcal{F}_k \right] \neq h(\omega_k)$$

$$\mathcal{F}_k := \sigma(\omega_0, X_1, \dots, X_k)$$

II. Two examples of *SA beyond gradient*

1st ex.: Compressed gradient

- **Compression operator** $\mathcal{C}(x, U)$, if the cost of storing/transmitting $\mathcal{C}(x, U)$ is less than the cost of storing/transmitting x

- First family:

$$\omega_{k+1} = \omega_k + \gamma_{k+1} \mathcal{C}(H(\omega_k, X_{k+1}), U_{k+1})$$

Ex. The Gauss-Southwell coordinate descent estimator $\mathcal{C}(x, u) = x_u e_u \quad u \in \{1, \dots, d\}$

- Second family:

$$\omega_{k+1} = \omega_k + \gamma_{k+1} H(\mathcal{C}(\omega_k, U_{k+1}), X_{k+1})$$

Ex. Stoch Gdt for deep learning, the Straight-Through Estimator quantizes the model ω_k before computing the oracle.

- Third family:

$$\omega_{k+1} = \mathcal{C}(\omega_k + \gamma_{k+1} H(\omega_k, X_{k+1}), U_{k+1})$$

Ex. Low precision Stoch Gdt: the model is quantized after computing the oracle.

2nd ex.: Stochastic Expectation Maximization in the curved exponential family

- The goal

$$\operatorname{argmin}_{\theta} f(\theta) := -\log \int_{\mathcal{D}} p(x; \theta) \nu(\mathrm{d}x) \qquad p(x; \theta) = \xi(x) \exp(\langle S(x), \phi(\theta) \rangle - \psi(\theta))$$

- The EM algorithm

$$\begin{aligned} \theta_k &\xrightarrow[\text{compute an expectation}]{\text{E-step}} \bar{S}(\theta_k) \xrightarrow[\text{optimize}]{\text{M-step}} \theta_{k+1} := \mathsf{T}(\bar{S}(\theta_k)) \\ s_k &:= \bar{S}(\theta_k) \xrightarrow[\text{optimize}]{\text{M-step}} \mathsf{T}(s_k) \xrightarrow[\text{compute an expectation}]{\text{E-step}} s_{k+1} := \bar{S}(\mathsf{T}(s_k)) \end{aligned}$$

where

$$\bar{S}(\theta) := \int_{\mathcal{D}} S(x) \pi(x; \theta) \nu(\mathrm{d}x) \qquad \mathsf{T}(s) := \operatorname{argmin}_{\theta} \psi(\theta) - \langle s, \phi(\theta) \rangle$$

- Limiting points of EM

EM finds θ_{\star} solving the root-finding pbm

$$\mathsf{T}(\bar{S}(\theta)) = \theta$$

or

EM finds s_{\star} solving the root-finding pbm

$$\bar{S}(\mathsf{T}(s)) = s$$

and then set

$$\theta_{\star} = \mathsf{T}(s_{\star}).$$

Which oracle of \bar{S} , in order to solve

$$\bar{S}(\mathbf{T}(\omega)) - \omega = 0 \qquad \bar{S}(\mathbf{T}(\omega)) := \int S(x) \pi(x; \mathbf{T}(\omega)) \nu(\mathrm{d}x)$$

- Stochastic Approximation EM (SAEM) when \bar{S} intractable

$$\omega_{k+1} = \omega_k + \gamma_{k+1} \left(\frac{1}{M} \sum_{m=1}^M S(X_{k+1,m}) - \omega_k \right)$$

where $X_{k+1,m}$ obtained from self-normalized importance sampling, MCMC, ...

- Mini-batch EM: when $S(x) := (1/n) \sum_{i=1}^n S_i(x)$ and large n

$$\omega_{k+1} = \omega_k + \gamma_{k+1} \left(\frac{1}{b} \sum_{i \in X_{k+1}} \bar{S}_i(\mathbf{T}(\omega_k)) - \omega_k \right).$$

III. Non-asymptotic convergence bounds in expectation

The assumptions

$$\omega_{k+1} = \omega_k + \gamma_{k+1} H(\omega_k, X_{k+1})$$

On the oracles, for some $W \geq 0$:

- **L^2 -moment** $\mathbb{E} [\|H(\omega_k, X_{k+1})\|^2] < \infty$
- **Growth of the mean field** $\exists c_0, c_1, \forall \omega \quad \|h(\omega)\|^2 \leq c_0 + c_1 W(\omega)$
- **Bias** $\exists \tau_0, \tau_1, \forall k, \mathbb{E} [H(\omega_k, X_{k+1}) | \mathcal{F}_k] - h(\omega_k) \|^2 \leq \tau_0 + \tau_1 W(\omega_k) \text{ a.s.}$
- **Variance**
 $\exists \sigma_0, \sigma_1, \forall k, \mathbb{E} [\|H(\omega_k, X_{k+1}) - \mathbb{E} [H(\omega_k, X_{k+1}) | \mathcal{F}_k]\|^2 | \mathcal{F}_k] \leq \sigma_0^2 + \sigma_1^2 W(\omega_k) \text{ a.s.}$

Ex. Gradient: $W(\omega) = \|\nabla f(\omega)\|^2$

Stoch EM: $W(\omega) = \|\bar{S}(\mathbf{T}(\omega)) - \omega\|^2$

A smooth Lyapunov function: V

- **Lower bounded** $\inf_{\omega} V(\omega) > -\infty$
- **Smooth fct** $V \text{ is } C^1 \text{ and } \exists L_V, \quad \|\nabla V(\omega) - \nabla V(\omega')\| \leq L_V \|\omega - \omega'\|$
- **Lyapunov V and control W** $\exists \rho \geq 0, \forall \omega \quad \langle \nabla V(\omega), h(\omega) \rangle \leq -\rho W(\omega)$

Ex. Gradient: $V(\omega) = f(\omega)$

Stoch EM: $V(\omega) = f(\mathbf{T}(\omega))$

If **biased** oracles i.e. $\tau_0 + \tau_1 > 0$, additional conditions

$$c_V := \sup_{\omega} \frac{\|\nabla V(\omega)\|^2}{W(\omega)} < \infty, \quad \sqrt{c_V} (\sqrt{\tau_0}/2 + \sqrt{\tau_1}) < \rho.$$

Main theorem

Theorem 1, Dieuleveut-F.-Moulines-Wai (2023)

In addition to the previous assumptions, assume that $\gamma_k \in (0, \gamma_{\max})$. Then, for any $T \geq 1$

$$\begin{aligned} \sum_{k=0}^{T-1} \frac{\gamma_{k+1} \mu_{k+1}}{\sum_{\ell=0}^{T-1} \gamma_{\ell+1} \mu_{\ell+1}} \mathbb{E}[W(\omega_k)] & \leq 2 \frac{\mathbb{E}[V(\omega_0)] - \min V}{\sum_{\ell=0}^{T-1} \gamma_{\ell+1} \mu_{\ell+1}} && \text{initial cond.} \\ & + L_V \eta_0 \frac{\sum_{k=0}^{T-1} \gamma_{k+1}^2}{\sum_{\ell=0}^{T-1} \gamma_{\ell+1} \mu_{\ell+1}} \\ & + 2b_0 \frac{\sum_{k=0}^{T-1} \gamma_{k+1}}{\sum_{\ell=0}^{T-1} \gamma_{\ell+1} \mu_{\ell+1}} && b_0 = 0 \text{ iff unbiased oracles} \end{aligned}$$

- Meaningful results under the assumptions

$$\min_{\omega: d(\omega, \{h=0\}) > \epsilon} W(\omega) > 0 \quad \forall \epsilon > 0$$

A Robbins-Siegmund type inequality

Lemma 9, Dieuleveut-F.-Moulines-Wai (2023)

$$\begin{aligned} \mathbb{E} \left[V(\omega_{k+1}) \middle| \mathcal{F}_k \right] &\leq V(\omega_k) - \underbrace{\gamma_{k+1} \mu_{k+1}}_{\substack{\text{positive for} \\ \gamma_{k+1} \text{ small enough}}} W(\omega_k) \\ &\quad + \underbrace{\gamma_{k+1} b_0}_{\substack{\geq 0 \text{ and zero} \\ \text{iff unbiased oracles}}} \\ &\quad + \gamma_{k+1}^2 \tilde{b} \end{aligned}$$

- From the assumptions on the Lyapunov function V

$$V(\omega_{k+1}) \leq V(\omega_k) + \langle \nabla V(\omega_k), \omega_{k+1} - \omega_k \rangle + \frac{L_V}{2} \|\omega_{k+1} - \omega_k\|^2.$$

- Use

$$\omega_{k+1} - \omega_k = \gamma_{k+1} h(\omega_k) + \gamma_{k+1} \left(H(\omega_k, X_{k+1}) - h(\omega_k) \right).$$

- Negative term: $\langle \nabla V(\omega_k), \gamma_{k+1} h(\omega_k) \rangle \leq -\rho \gamma_{k+1} W(\omega_k)$
- Apply the conditional expectations $\mathbb{E}[\cdot | \mathcal{F}_k]$ and use the assumptions on the bias and variance of the oracles.

Corollary 1: which parameter ?

On the left hand side:

$$\sum_{k=0}^{T-1} \frac{\gamma_{k+1} \mu_{k+1}}{\sum_{\ell=0}^{T-1} \gamma_{\ell+1} \mu_{\ell+1}} \mathbb{E} [W(\omega_k)]$$

- If **W is convex**,

$$W \left(\sum_{k=0}^{T-1} \frac{\gamma_{k+1} \mu_{k+1}}{\sum_{\ell=0}^{T-1} \gamma_{\ell+1} \mu_{\ell+1}} \omega_k \right) \leq \sum_{k=0}^{T-1} \frac{\gamma_{k+1} \mu_{k+1}}{\sum_{\ell=0}^{T-1} \gamma_{\ell+1} \mu_{\ell+1}} W(\omega_k)$$

adopt a convex combination of the iterates

- Otherwise,

$$\sum_{k=0}^{T-1} \frac{\gamma_{k+1} \mu_{k+1}}{\sum_{\ell=0}^{T-1} \gamma_{\ell+1} \mu_{\ell+1}} W(\omega_k) = \mathbb{E} [W(\omega_R)] \quad \mathbb{P}(R = k) \propto \gamma_{k+1} \mu_{k+1}.$$

stop at a random time / choose randomly one of the iterates

Corollary 2: Constant step size

$$\gamma := \frac{\gamma_{\max}}{2} \wedge O\left(\frac{1}{\sqrt{T}}\right)$$

then

$$\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E}[W(\omega_k)] \leq \underbrace{B}_{\text{null}} + \frac{A_1}{\sqrt{T}} \wedge \frac{A_2}{T}$$

iff unbiased oracles

- If **unbiased** oracles,
the RHS goes to zero when $T \rightarrow +\infty$ by choosing $\gamma_k = \gamma_{\max}/\sqrt{T}$
the convergence rate of SA is $O(1/\sqrt{T})$.
- If **biased oracles**: the RHS can not be made small when the step sizes are constant.

Corollary 3: ϵ -approximate stationarity

In non-convex optimization: in general, it is intractable to find a global minimum or to test if a point is a local minimum.

- **Stationarity as a convergence criterion:** For a precision ϵ , find a random stopping time R s.t.

$$\mathbb{E}[W(\omega_R)] \leq \epsilon.$$

- When the oracles are **unbiased**: choose R *uniform* on $\{1, \dots, T\}$ for T larger than

$$T(\epsilon) := \frac{A_3}{\epsilon^2} \vee \frac{A_4}{\epsilon}.$$

high-precision regime:	$T(\epsilon) = O(1/\epsilon^2)$	step size $\gamma_\epsilon = O(\epsilon)$
low-precision regime:	$T(\epsilon) = O(1/\epsilon)$	step size $\gamma_\epsilon = \gamma_{\max}/2$

- Section III-A, Dieuleveut-F.-Moulines-Wai (2023) **explicit constants**, not detailed here

IV. Variance reduction by SPIDER

Control variates

- The oracles are not unique:

$$\mathbb{E}[H(\omega, X)] = h(\omega) \implies \mathbb{E}[H(\omega, X) + U] = h(\omega) \quad \text{where} \quad \mathbb{E}[U] = 0.$$

- Choose U **correlated with the natural oracle** $H(\omega, X)$ s.t.

$$\text{Var}(H(\omega, X) + U) \ll \text{Var}(H(\omega, X))$$

The SPIDER control variate when h is a finite sum

Adapted from the gradient case: **Stochastic Path-Integrated Differential Estimator**

For problems of the form

$$\omega : \quad h(\omega) = 0 \quad \text{when} \quad h(\omega) = \frac{1}{n} \sum_{i=1}^n h_i(\omega) \quad \text{and } n \text{ large}$$

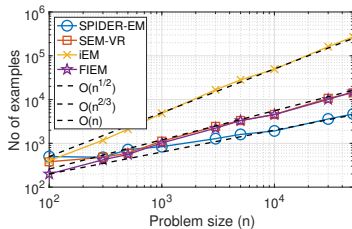
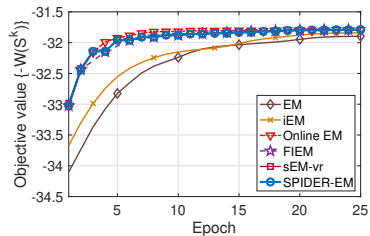
- At iteration $\#k$, a natural oracle for $h(\omega_k)$ is

$$H(\omega_k, X_{k+1}) := \frac{1}{b} \sum_{i \in X_{k+1}} h_i(\omega_k) \quad X_{k+1} \text{ mini-batch from } \{1, \dots, n\}, \text{ of size } b$$

- The **SPIDER oracle** is

$$H_{k+1}^{\text{SP}} := \frac{1}{b} \sum_{i \in X_{k+1}} h_i(\omega_k) + \underbrace{H_k^{\text{SP}}}_{\text{oracle for } h(\omega_{k-1})} - \underbrace{\frac{1}{b} \sum_{i \in X_{k+1}} h_i(\omega_{k-1})}_{\text{oracle for } h(\omega_{k-1})}$$

Efficiency ... via plots (here)



Conclusion

- A unifying framework for SA, that covers gradient SA, non-gradient SA, possibly with biased oracles is introduced.
- Explicit controls of convergence in expectation are provided.
- From which are deduced: stopping rules strategies, constant step sizes strategies, rates of convergence.

Stochastic Approximation Beyond Gradient for Signal Processing and Machine Learning

by A. Dieuleveut, G. Fort, E. Moulines and H.-T. Wai

HAL-03979922 and arXiv:2302.11147