

Nonasymptotic Analysis of Stochastic Gradient Descent with the Richardson–Romberg Extrapolation

Alexey Naumov

International Laboratory of Stochastic Algorithms and High-Dimensional Inference
HSE University



October 10, 2024

Joint work with



Denis Belomestny
(Duisburg-Essen University & HSE)



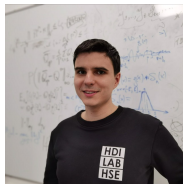
Alain Durmus
(ENS Paris-Saclay)



Marina Sheshukova
(HSE)



Eric Moulines
(Ecole Polytechnique)



Sergey Samsonov
(HSE)

Stochastic minimization problem

- Consider the strongly convex minimization problem, which admits a unique solution θ^*

$$f(\theta) \quad \min_{\theta \in \mathbb{R}^d} . \quad (1)$$

- The access to $f(\theta)$ is available only through the (unbiased) noisy observations $r F(\theta, \xi)$, where ξ is a random variable on (Z, \mathcal{Z}) .
- We solve the problem (1) using the SGD with constant step size γ , starting from initial distribution ν :

$$\theta_{k+1}^{(\gamma)} = \theta_k^{(\gamma)} - \gamma r F(\theta_k^{(\gamma)}, \xi_{k+1}), \quad \theta_0 \sim \nu, \quad (2)$$

where ξ_k is a sequence of i.i.d. random variables.

- Define noise function as

$$\varepsilon_k(\theta) = r F(\theta, \xi_k) - r f(\theta), \quad (3)$$

and noise covariance matrix as $\Sigma_\varepsilon^* = \mathbb{E}[r F(\theta^*, \xi) \otimes r F(\theta^*, \xi)]$

Polyak-Ruppert averaged estimator

Consider the Polyak-Ruppert averaged estimator

$$\bar{\theta}_n^{(\gamma)} = \frac{1}{n} \sum_{k=n+1}^{2n} \theta_k^{(\gamma)}. \quad (4)$$

- Under appropriate assumptions on f and γ_k ,

$$\rho_{\bar{\theta}_{n_0, n}}(\bar{\theta}_{n_0, n} - \theta^*) \stackrel{d}{\sim} \mathcal{N}(0, H^{-1} \Sigma_{\varepsilon}^* H^{-1}), \quad n \rightarrow \infty, \quad (5)$$

where $H = r^2 f''(\theta^*)$; e.g. Fort [2015]).

- We examine the mean-squared error bounds in the following form:

$$\mathbb{E}^{1/2}[\|\bar{\theta}_{n_0, n} - \theta^*\|^2] \leq \frac{\sqrt{\text{Tr}(H^{-1} \Sigma_{\varepsilon}^* H^{-1})}}{n^{1/2}} + \frac{C(f, d)}{n^{1/2+\delta}} + \dots \quad (6)$$

- The goal of the work is to obtain the result in the form (6) with the best possible constant δ .

Related works: Polyak-Ruppert averaged estimator

- | Previous studies considered decreasing step size in the dynamics (2):
 - | In [Moulines and Bach \[2011\]](#) for strongly convex functions it was shown that

$$\mathbb{E}^{1/2}[k\bar{\theta}_n - \theta \mid k^2] \leq \frac{\sqrt{\text{Tr}(H^{-1}\Sigma_\varepsilon^*H^{-1})}}{\rho_{\bar{n}}} + O(n^{-7/12})$$

- | [Gadat and Panloup \[2023\]](#) improved the result of [Moulines and Bach \[2011\]](#) and for a certain class of functions f , including strongly convex functions it was shown that

$$\mathbb{E}^{1/2}[k\bar{\theta}_n - \theta \mid k^2] \leq \frac{\sqrt{\text{Tr}(H^{-1}\Sigma_\varepsilon^*H^{-1})}}{\rho_{\bar{n}}} + O(n^{-5/8})$$

Analysis of SGD: Assumptions on function f

Assumption A1

The function f is μ -strongly convex on \mathbb{R}^d , that is, it is continuously differentiable and there exists a constant $\mu > 0$, such that for any $\theta, \theta^0 \in \mathbb{R}^d$, it holds that

$$\frac{\mu}{2} \|\theta - \theta^0\|^2 \leq f(\theta) - f(\theta^0) \leq L \|\theta - \theta^0\|. \quad (7)$$

Assumption A2

The function f is 4 times continuously differentiable and L_2 -smooth on \mathbb{R}^d , i.e., it is continuously differentiable and there is a constant $L_2 > 0$, such that for any $\theta, \theta^0 \in \mathbb{R}^d$,

$$\| \nabla f(\theta) - \nabla f(\theta^0) \| \leq L_2 \|\theta - \theta^0\|. \quad (8)$$

Moreover, f has uniformly bounded 3-rd and 4-th derivatives, such that

$$\| \nabla^i f(\theta) \| \leq L_i, \text{ for } i \in \{3, 4\}. \quad (9)$$

Analysis of SGD: Assumptions on the noisy gradient $r F$

Assumption A3(p)

$\{\xi_k\}_{k \in \mathbb{N}}$ is a sequence of independent and identically distributed (i.i.d.) random variables with distribution P_ξ , such that ξ_i and θ_0 are independent and for any $\theta \in \mathbb{R}^d$ it holds that

$$\mathbb{E}_{\xi \sim P_\xi}[r F(\theta, \xi)] = r f(\theta).$$

Moreover, there exists τ_p , such that $\mathbb{E}^{1/p}[k r F(\theta^*, \xi) k^p] \leq \tau_p$, and for any $q = 2, \dots, p$ it holds with some $L_1 > 0$ that for any $\theta_1, \theta_2 \in \mathbb{R}^d$,

$$L_1^q \|\theta_1 - \theta_2\|^q \leq \mathbb{E}_{\xi \sim P_\xi}[k r F(\theta_1, \xi) - r F(\theta_2, \xi) k^q]. \quad (10)$$

Note that, A3(p) generalizes the well-known L_1 -co-coercivity assumption, see [Dieuleveut et al. \[2020b\]](#). A sufficient condition is to assume that $F(\theta, \xi)$ is P_ξ -a.s. convex with respect to $\theta \in \mathbb{R}^d$.

Kantorovich - Wasserstein distance

Definition

The function $c : Z \times Z \rightarrow \mathbb{R}_+$ is called distance-like if it is symmetric, lower semi-continuous and $c(x, y) = 0$ if and only if $x = y$.

Definition

For two probability measures ξ and ξ^θ we denote by $\mathcal{C}(\xi, \xi^\theta)$ the set of couplings of two probability measures, that is, for any $C \in \mathcal{C}(\xi, \xi^\theta)$ and any $A \subseteq Z$ it holds $C(Z \times A) = \xi^\theta(A)$ and $C(A \times Z) = \xi(A)$. We define

$$\mathbf{W}_c(\xi, \xi^\theta) = \inf_{C \in \mathcal{C}(\xi, \xi^\theta)} \int_{Z \times Z} c(z, z^\theta) C(dz, dz^\theta). \quad (11)$$

Analysis of SGD: Bias

- Under assumptions A1-A3(2) the sequence $f\theta_k^{(\gamma)} g_{k \geq N}$ is a homogeneous Markov chain with the Markov kernel

$$Q_\gamma(\theta, A) = \int_{\mathbb{R}^d} \mathbb{1}_A(\theta - \gamma F(\theta, z)) P_\xi(dz), \quad \theta \in \mathbb{R}^d, A \subseteq \mathbb{B}(\mathbb{R}^d); \quad (12)$$

see Dieuleveut et al. [2020a].

- Introduce distance-like function

$$c(\theta, \theta^0) = k\theta - \theta^0 k(k\theta - \theta^0 k + k\theta^0 - \theta^0 k + c_0\gamma^{1/2}), \quad (13)$$

where we set $c_0 = 2^{3/2}\tau_2/\mu^{1/2}$

Lemma 1

Assume A1-A3(2). Then, for any $\gamma \in (0; \frac{1}{2L}]$ the Markov kernel Q_γ admits a unique invariant probability measure π_γ . Moreover, for all $\theta \in \mathbb{R}^d$ and $k \in \mathbb{N}$

$$\mathbf{W}_c(\nu Q_\gamma^k, \pi_\gamma) \leq 4(1/2)^{k/m(\gamma)} \mathbf{W}_c(\nu, \pi_\gamma), \text{ where } m(\gamma) = d2 \frac{\log 4}{\gamma\mu} e \quad (14)$$

- However, unless the function f is quadratic,

$$\mathbb{E}_{\pi_\gamma}[\theta] \neq \theta^*.$$

Analysis of SGD: Bias

- | We consider the following condition.

Assumption C1(p)

There exist constants $D_{\text{last},p}, C_{\text{step},p} \geq 2$ depending only on p , such that for any step size $\gamma \geq (0, 1/(L C_{\text{step},p})]$, and any initial distribution ν it holds that

$$\mathbb{E}_{\nu}^{2/p} [k\theta_k^{(\gamma)} - \theta^* k^p] \leq (1 - \gamma\mu)^k \mathbb{E}_{\nu}^{2/p} [k\theta_0 - \theta^* k^p] + D_{\text{last},p} \gamma \tau_p^2 / \mu. \quad (15)$$

Moreover, for the stationary distribution π_{γ} it holds that

$$\mathbb{E}_{\pi_{\gamma}}^{2/p} [k\theta_0^{(\gamma)} - \theta^* k^p] \leq D_{\text{last},p} \gamma \tau_p^2 / \mu. \quad (16)$$

- | It is important to recognize that C1(p) is not independent from the previous assumptions A1-A3(p). In particular, [Dieuleveut et al., 2020a, Lemma 13] implies that, under A1-A3(p) with $p \geq 2$, the bound (16) holds for $\gamma \geq (0, 1/(L C_{\text{step},p})]$ with some constants $D_{\text{last},p}$ and $C_{\text{step},p}$, which depends only upon p .

Analysis of SGD: Bias

Proposition 2, Theorem 4 in [Dieuleveut et al. \[2020a\]](#)

A1-A3(6), **C1(6)**. Then, for any $\gamma \geq (0, 1/(L C_{\text{step},6})]$, the following bias expansion holds:

$$\mathbb{E}_{\pi_\gamma}[\theta] := \int_{\mathbb{R}^d} u \pi_\gamma(du) = \theta^* + \gamma \Delta_1 + O(\gamma^{3/2}), \quad (17)$$

$$\mathbb{E}_{\pi_\gamma}[(\theta - \theta^*)^2] := \int_{\mathbb{R}^d} (u - \theta^*)^2 \pi_\gamma(du) = \gamma \Delta_2 + O(\gamma^{3/2}), \quad (18)$$

where $\Delta_1 \in \mathbb{R}^d$, $\Delta_2 \in \mathbb{R}^{d \times d}$ are constants independent of the step size γ . Moreover, for any starting point $\theta_0 \in \mathbb{R}^d$, it holds that

$$\mathbb{E}[\bar{\theta}_n] = \theta^* + \gamma \Delta_1 + O(\gamma^{3/2}) + \Delta_1(k\theta_0 - \theta^*, \gamma, n), \quad (19)$$

where $k\Delta_1(k\theta_0 - \theta^*, \gamma, n)k \leq \frac{e^{-\gamma\mu(n+1)/2}}{n\gamma\mu} (k\theta_0 - \theta^*k + \frac{\rho}{\beta} \frac{\sqrt{\gamma T_2}}{\mu})$.

Bound for Polyak-Ruppert averaged estimator

Theorem 3

Assume [A1-A3\(6\)](#), [C1\(6\)](#). Then for any $\gamma \geq (0, 1/(L C_{\text{step},6})]$ and any $n \geq N$, the sequence of Polyak-Ruppert estimates [\(4\)](#) satisfies

$$\mathbb{E}_\nu^{1/2}[kH(\bar{\theta}_n^{(\gamma)} - \theta^*)k^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\varepsilon^*}}{n^{1/2}} + \frac{1}{\gamma^{1/2}n} + \gamma + \frac{\gamma^{1/2}}{n^{1/2}} + R_1(n, \gamma, k\theta_0 - \theta^*k), \quad (20)$$

where

$$R_1(n, \gamma, k\theta_0 - \theta^*k) \leq \frac{e^{\gamma\mu(n+1)/2}}{\gamma n} (\mathbb{E}_\nu^{1/2}[k\theta_0 - \theta^*k^2] + \mathbb{E}_\nu^{1/2}[k\theta_0 - \theta^*k^4])$$

Theorem 3: sketch of proof

| Note that

$$\theta_{k+1} - \theta^* = \theta_k - \theta^* - \gamma(r f(\theta_k) + \varepsilon_{k+1}(\theta_k)),$$

where $\varepsilon_{k+1}(\theta_k)$ is a martingale-difference sequence w.r.t. F_k .

| Set

$$\eta(\theta) = r f(\theta) - H(\theta - \theta^*),$$

| We get

$$\theta_{k+1} - \theta^* = (I - \gamma H)(\theta_k - \theta^*) - \gamma \varepsilon_{k+1}(\theta_k) - \gamma \eta(\theta_k),$$

| Rearranging the terms, we obtain

$$H(\theta_k - \theta^*) = \frac{\theta_k - \theta_{k+1}}{\gamma} - \varepsilon_{k+1}(\theta_k) - \eta(\theta_k). \quad (21)$$

Theorem 3: sketch of proof

- | Summing the recurrence (21), we obtain that

$$\begin{aligned} H(\bar{\theta}_n - \theta^*) &= \frac{\theta_{n+1} - \theta^*}{\gamma n} - \frac{\theta_{2n} - \theta^*}{\gamma n} - \frac{1}{n} \sum_{k=n+1}^{2n} \{\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\} \\ &\quad + \frac{1}{n} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*) - \frac{1}{n} \sum_{k=n+1}^{2n} \eta(\theta_k). \end{aligned} \tag{22}$$

- | Applying the 3-rd order Taylor expansion with integral remainder, we get that

$$\eta(\theta_k) = \nabla f(\theta_k) - H(\theta_k - \theta^*) = \frac{1}{2} \left(\int_0^1 \nabla^3 f(t\theta^* + (1-t)\theta_k) dt \right) (\theta_k - \theta^*)^2.$$

- | Using A2, we thus obtain that

$$\|\eta(\theta_k)\| \leq \frac{1}{2} L_3 \|\theta_k - \theta^*\|^2.$$

- | Hence, applying Minkowski's inequality, we get

$$\begin{aligned} \mathbb{E}_\nu^{1/2} [\|H(\bar{\theta}_n - \theta^*)\|^2] &\leq \frac{\mathbb{E}_\nu^{1/2} [\|\theta_{n+1} - \theta^*\|^2]}{\gamma n} + \frac{\mathbb{E}_\nu^{1/2} [\|\theta_{2n} - \theta^*\|^2]}{\gamma n} \\ &+ \frac{1}{n} \mathbb{E}_\nu^{1/2} [\|\sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\|^2] + \frac{1}{n} \mathbb{E}_\nu^{1/2} [\|\sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*)\|^2] \\ &\quad + \frac{L_3}{2n} \sum_{k=n+1}^{2n} \mathbb{E}_\nu^{1/2} [\|\theta_k - \theta^*\|^4] \end{aligned}$$

Theorem 3: sketch of proof

| We have

$$\mathbb{E}_\nu \left[k \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta_k) \varepsilon_{k+1}(\theta^*) k^2 \right] = \sum_{k=n+1}^{2n} f \mathbb{E}_\nu \left[k \varepsilon_{k+1}(\theta_k) \varepsilon_{k+1}(\theta^*) k^2 \right]$$

| Since ξ_k are i.i.d, we have

$$\mathbb{E}_\nu^{1/2} \left[k \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*) k^2 \right] = \sqrt{n \operatorname{Tr} \Sigma_\varepsilon^*} \quad (23)$$

| Using [A3](#), we get

$$\mathbb{E}_\nu \left[k \varepsilon_{k+1}(\theta_k) \varepsilon_{k+1}(\theta^*) k^2 \right] \leq L^2 \mathbb{E}_\nu \left[k \theta_k \theta^* k^2 \right]. \quad (24)$$

| Applying [C1](#)[4], we complete the proof.

Theorem 3: Discussion

- Setting optimal step size $\gamma \propto n^{-2/3}$, which yields an error bound of order:

$$\mathbb{E}^{1/2}[kH(\bar{\theta}_n - \theta^*)k^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\varepsilon^*}}{n^{1/2}} + O\left(\frac{1}{n^{2/3}}\right) + R_1(k\theta_0 - \theta^*k, n). \quad (25)$$

- There are different results in the literature that provide various decay rates of the second-order term in (25). However, all these results are known to be suboptimal for the first-order methods.
- In fact, the recent result of Li et al. [2022] shows that a second-order error can be achieved by modifying the SGD algorithm with averaging and control variates

$$\mathbb{E}^{1/2}[kH(\bar{\theta}_n - \theta^*)k^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\varepsilon^*}}{n^{1/2}} + O\left(\frac{1}{n^{3/4}}\right) + R_1^0(k\theta_0 - \theta^*k, n). \quad (26)$$

- Our goal is to obtain an analogue of the 2-moment bound (26), using a simpler algorithm.

Richardson-Romberg estimator

- | We construct two parallel chains based on the same sequence of noise variables $f_{\xi_k} g_{k,2N}$:

$$\begin{aligned}\theta_{k+1}^{(\gamma)} &= \theta_k^{(\gamma)} - \gamma r F(\theta_k^{(\gamma)}, \xi_{k+1}), & \bar{\theta}_n^{(\gamma)} &= \frac{1}{n} \sum_{k=n+1}^{2n} \theta_k^{(\gamma)}, \\ \theta_{k+1}^{(2\gamma)} &= \theta_k^{(2\gamma)} - 2\gamma r F(\theta_k^{(2\gamma)}, \xi_{k+1}), & \bar{\theta}_n^{(2\gamma)} &= \frac{1}{n} \sum_{k=n+1}^{2n} \theta_k^{(2\gamma)}.\end{aligned}\quad (27)$$

- | Based on $\bar{\theta}_n^{(\gamma)}$ and $\bar{\theta}_n^{(2\gamma)}$ defined above, we construct a Richardson-Romberg estimator as

$$\bar{\theta}_n^{(RR)} := 2\bar{\theta}_n^{(\gamma)} - \bar{\theta}_n^{(2\gamma)}.\quad (28)$$

- | Note that in the decomposition (22), the linear statistics $W = n^{-1} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*)$ does not depend upon γ . Hence, using the same sequence $f_{\xi_k} g_{k,2N}$ of noise variables in (27) yields an estimator $\bar{\theta}_n^{(RR)}$, such that its leading component of the variance still equals W .

Richardson-Romberg estimator: Bias

Proposition 4

Assume [A1-A3\(6\)](#), [C1\(6\)](#). Then, for any $\gamma \in (0, 1/(L C_{\text{step},6})]$, and any starting point $\theta_0 \in \mathbb{R}^d$, it holds that

$$\mathbb{E}_\nu[\bar{\theta}_n^{(RR)}] = \theta^* + O(\gamma^{3/2}) + \Delta_2(k\theta_0 - \theta^*, \gamma, n), \quad (29)$$

where $k\Delta_2(k\theta_0 - \theta^*, \gamma, n)k \leq 3 \frac{e^{-\gamma\mu(n+1)/2}}{n\gamma\mu} (\mathbb{E}_\nu^{1/2}[k\theta_0 - \theta^*k^2] + \frac{\rho_{2\gamma\tau_2}}{\beta\mu})$.

Bound for Richardson-Romberg estimator

Theorem 5

Assume A1-A4(6). Then for any $\gamma \geq (0; \min(\frac{1}{9L}, \frac{1}{C_3L})]$ and any $n \geq N$, the estimator defined in 28 satisfies

$$\begin{aligned} E_{\nu}^{1/2}[kH(\bar{\theta}_n^{(RR)} - \theta^*)k^2] &\leq \frac{\sqrt{\text{Tr} \Sigma_{\epsilon}^*}}{n^{1/2}} + \frac{\gamma^{1/2}}{n^{1/2}} + \frac{1}{\gamma^{1/2}n} + \frac{\gamma}{n^{1/2}} + \gamma^{3/2} \\ &\quad + R_2(n, \gamma, k\theta - \theta^*k), \end{aligned}$$

where

$$\begin{aligned} R_2(n, \gamma, k\theta_0 - \theta^*k) &\leq \frac{e^{-\gamma\mu(n+1)/2}}{\gamma n} (E_{\nu}^{1/2}[k\theta_0 - \theta^*k^2] \\ &\quad + E_{\nu}^{1/2}[k\theta_0 - \theta^*k^4] + E_{\nu}^{1/2}[k\theta_0 - \theta^*k^6] + \gamma) \end{aligned}$$

Theorem 5: sketch of proof

- Using the recursion 21, we obtain that

$$\begin{aligned}
 H(\bar{\theta}_n^{(RR)} \quad \theta^*) &= 2 \frac{\theta_{n+1}^\gamma \quad \theta^*}{\gamma n} \quad 2 \frac{\theta_{2n}^\gamma \quad \theta^*}{\gamma n} \quad \frac{\theta_{n+1}^{2\gamma} \quad \theta^*}{2\gamma n} + \frac{\theta_{2n}^{2\gamma} \quad \theta^*}{2\gamma n} \\
 &\quad \frac{2}{n} \sum_{k=n+1}^{2n} [\varepsilon_{k+1}(\theta_k^\gamma) \quad \varepsilon_{k+1}(\theta^*)] + \frac{1}{n} \sum_{k=n+1}^{2n} [\varepsilon_{k+1}(\theta_k^{2\gamma}) \quad \varepsilon_{k+1}(\theta^*)] \\
 &\quad + \frac{1}{n} \sum_{k=n+1}^{2n} [\varepsilon_{k+1}(\theta^*)] \quad \frac{1}{n} \sum_{k=n+1}^{2n} [2\eta(\theta_k^\gamma) \quad \eta(\theta_k^{2\gamma})]
 \end{aligned}$$

- Define the function $\psi(\theta) = (1/2)r^3 f(\theta) (\theta \quad \theta^*)^2$
- Applying the 4-rd order Taylor expansion with integral remainder, we get that

$$\eta(\theta) = \psi(\theta) + \frac{1}{6} \left(\int_0^1 r^4 f(t\theta^* + (1-t)\theta) dt \right) (\theta \quad \theta^*)^3, \quad (30)$$

- Using A2, we obtain

$$k \frac{1}{6} \left(\int_0^1 r^4 f(t\theta^* + (1-t)\theta) dt \right) (\theta \quad \theta^*)^3 \leq \frac{1}{6} L_4 k \theta \quad \theta^* k^3. \quad (31)$$

Theorem 5: sketch of proof

- The key technical element of the proof is to bound

$$\frac{1}{n} \mathbb{E}_{\nu}^{1/2} \left[k \sum_{k=n+1}^{2n} f_{\psi}(\theta_k) - \pi_{\gamma}(\psi) g k^2 \right]$$

- Using coupling technique, it can be shown that

$$\begin{aligned} \frac{1}{n} \mathbb{E}^{1/2} \left[k \sum_{k=n+1}^{2n} f_{\psi}(\theta_k) - \pi_{\gamma}(\psi) g k^2 \right] &\leq \frac{1}{n} \mathbb{E}_{\pi_{\gamma}}^{1/2} \left[k \sum_{k=n+1}^{2n} f_{\psi}(\theta_k) - \pi_{\gamma}(\psi) g k^2 \right] \\ &\quad + \frac{e^{-\gamma \mu(n+1)/2}}{\gamma n} \left(\mathbb{E}_{\nu}^{1/2} [k \theta_0 - \theta^* k^4] + \gamma \right) \end{aligned}$$

- It can be shown that for any $\theta, \theta^0 \in \mathbb{R}^d$, it holds that

$$|k\psi(\theta) - \psi(\theta^0)k| \leq \frac{1}{2} L_3 c(\theta, \theta^0). \quad (32)$$

- Moreover, see [Douc et al. \[2018\]](#), for any start point $\theta_0 \in \mathbb{R}^d$, it holds

$$|jQ^k \psi(\theta_0) - \pi_{\gamma}(\psi)j| \leq L_3 (1/2)^{k/m(\gamma)} \mathbf{W}_c(\delta_{\theta_0}, \pi_{\gamma}).$$

Theorem 5: sketch of proof

- | For covariance term, we get

$$\mathbb{E}_{\pi_\gamma} [(\psi(\theta_0) - \pi_\gamma(\psi))^T (\psi(\theta_k) - \pi_\gamma(\psi))] = (1/2)^{k/m(\gamma)} \gamma^2$$

- | For variance term, we have

$$\mathbb{E}_{\pi_\gamma} [k(\psi(\theta_0) - \pi_\gamma(\psi))^2] = \mathbb{E}_{\pi_\gamma} [k(\psi(\theta_0))^2] = \gamma^2$$

- | Combining results, we obtain

$$\frac{1}{n} \mathbb{E}_{\pi_\gamma}^{1/2} [k \sum_{k=n+1}^{2n} f(\psi(\theta_k) - \pi_\gamma(\psi))^2] = \frac{\gamma}{n^{1/2}} + \frac{\gamma^{1/2}}{n^{1/2}}$$

- | It remains to note that from Proposition 2 $k(2\pi_\gamma(\psi) - \pi_{2\gamma}(\psi))k = \gamma^{3/2}$
- | To obtain result of Theorem 5 it remains to apply Minkowski's inequality to the decomposition 30.

Theorem 5: Discussion

- Setting optimal step size γ depending on the number of samples n , we arrive at $\gamma = n^{-1/2}$, which yields an error bound of order:

$$\mathbb{E}_\nu^{1/2}[kH(\bar{\theta}_n^{RR} - \theta^*)k^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\varepsilon^*}}{n^{1/2}} + O\left(\frac{1}{n^{3/4}}\right) + R_2(k\theta_0 - \theta^*k, n).$$

- Now we aim to generalize this result for the p -th moment bounds with $p \geq 2$.

Richardson-Romberg estimator, p -th moment

- The key technical element of our proof for the p -th moment bound is the following statement, which can be viewed as a version of Rosenthal's inequality [Rosenthal, 1970, Pinelis, 1994].

Proposition 6

Let $p \geq 2$ and assume **A1-A3**($2p$), and **C1**($2p$). Then for any $\gamma \geq (0, 1/(L C_{\text{step},2p})]$, it holds that

$$\mathbb{E}_{\pi_\gamma}^{1/p} \left[k \sum_{k=0}^{n-1} f(\psi(\theta_k^{(\gamma)})) - \pi_\gamma(\psi) k^p \right] \leq \frac{L D_{\text{last},2p} p \tau_{2p}^2}{\mu^{3/2}} \frac{1}{n^\gamma} + \frac{L D_{\text{last},2p} \tau_{2p}}{\mu^2},$$

where $\psi(\theta) = (1/2) r^{-3} f(\theta) (\theta - \theta^*)^2$.

Richardson-Romberg estimator, p -th moment

Repeating the proof of Theorem 5, and using Proposition 6, we obtain the following bound:

Theorem 7

Let $p \geq 2$ and assume **A1-A3**($3p$) and **C1**($3p$). Then for any $\gamma \geq (0, 1/(L C_{\text{step},3p}))$ and any $n \geq N$, the estimator defined in 28 satisfies

$$\begin{aligned} \mathbb{E}_{\nu}^{1/p}[kH(\bar{\theta}_n^{(RR)} - \theta^*)k^p] &\leq \frac{\sqrt{\text{Tr} \Sigma_{\varepsilon}^*}}{n^{1/2} p^{1/2}} + \frac{1}{n^{1-1/p}} + \frac{\gamma^{1/2}}{n^{1/2}} + \frac{1}{\gamma^{1/2} n} \\ &\quad + \frac{\gamma}{n^{1/2}} + \gamma^{3/2} + R_3(n, \gamma, k\theta - \theta^*k), \end{aligned}$$

where

$$\begin{aligned} R_3(n, \gamma, k\theta_0 - \theta^*k) &\leq \frac{e^{-\gamma \mu(n+1)/2}}{\gamma n} (\mathbb{E}_{\nu}^{1/p}[k\theta_0 - \theta^*k^p] \\ &\quad + \mathbb{E}_{\nu}^{1/p}[k\theta_0 - \theta^*k^{2p}] + \mathbb{E}_{\nu}^{1/p}[k\theta_0 - \theta^*k^{3p}] + \gamma) \end{aligned}$$

Theorem 7: Discussion

- | Setting optimal step size γ depending on the number of samples n , we arrive at $\gamma = n^{-1/2}$, which yields an error bound of order:

$$\mathbb{E}_\nu^{1/p}[\|kH(\bar{\theta}_n^{RR} - \theta^*)\|^p] \leq \frac{\sqrt{\text{Tr} \Sigma_\varepsilon^*} p^{1/2}}{n^{1/2}} + O\left(\frac{1}{n^{3/4}}\right) + R_3(k\theta_0 - \theta^*, k, n).$$

Thank you!

References I

- Aymeric Dieuleveut, Alain Durmus, and Francis Bach. Bridging the gap between constant step size stochastic gradient descent and Markov chains. *The Annals of Statistics*, 48(3):1348 – 1382, 2020a. doi: 10.1214/19-AOS1850. URL <https://doi.org/10.1214/19-AOS1850>.
- Aymeric Dieuleveut, Alain Durmus, and Bach Francis. Bridging the gap between constant step size stochastic gradient descent and markov chains. *Annals of Statistics*, 2020b.
- R. Douc, E. Moulines, P. Priouret, and P. Soulier. *Markov chains*. Springer Series in Operations Research and Financial Engineering. Springer, 2018. ISBN 978-3-319-97703-4.
- Gersende Fort. Central limit theorems for stochastic approximation with controlled markov chain dynamics. *ESAIM: Probability and Statistics*, 19:60–80, 2015.
- Sébastien Gadat and Fabien Panloup. Optimal non-asymptotic analysis of the Ruppert–Polyak averaging stochastic algorithm. *Stochastic Processes and their Applications*, 156:312–348, 2023. ISSN 0304-4149. doi: <https://doi.org/10.1016/j.spa.2022.11.012>. URL <https://www.sciencedirect.com/science/article/pii/S0304414922002447>.
- Chris Junchi Li, Wenlong Mou, Martin Wainwright, and Michael Jordan. Root-sgd: Sharp nonasymptotics and asymptotic efficiency in a single algorithm. In *Conference on Learning Theory*, pages 909–981. PMLR, 2022.
- Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. *Advances in neural information processing systems*, 24, 2011.
- Iosif Pinelis. Optimum Bounds for the Distributions of Martingales in Banach Spaces. *The Annals of Probability*, 22(4):1679 – 1706, 1994. doi: 10.1214/aop/1176988477. URL <https://doi.org/10.1214/aop/1176988477>.
- Haskell P. Rosenthal. On the subspaces of L^p ($p > 2$) spanned by sequences of independent random variables. *Israel J. Math.*, 8:273–303, 1970. ISSN 0021-2172. doi: 10.1007/BF02771562. URL <https://doi.org/10.1007/BF02771562>.