



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

Marina Sheshukova¹, Denis Belomestny^{2,1}, Alain Durmus³, Eric Moulines^{3,4}, Alexey Naumov^{1,5}, Sergey Samsonov¹

¹HSE University ²Duisburg–Essen University ³CMAF, UMR 7641, École Polytechnique

⁴Mohamed Bin Zayed University of AI

⁵Steklov Mathematical Institute of the Russian Academy of Sciences

Problem setting and known results

- Stochastic gradient methods aim to minimize a function f with access only to the noisy gradients:

$$\min_{\theta \in \mathbb{R}^d} f(\theta), \quad \nabla f(\theta) = \mathbb{E}_{\xi \sim \mathbb{P}_\xi} [\nabla F(\theta, \xi)]. \quad (1)$$

Here ξ is a noise variable with the distribution \mathbb{P}_ξ and θ^* is the unique minimizer. The standard SGD algorithm writes as

$$\theta_{k+1} = \theta_k - \gamma_{k+1} \nabla F(\theta_k, \xi_{k+1}), \quad \theta_0 \in \mathbb{R}^d.$$

- A well-known technique for training stabilization is the Polyak-Ruppert averaging:

$$\bar{\theta}_{n_0, n} = n^{-1} \sum_{k=n_0+1}^{n+n_0} \theta_k,$$

where n_0 is the burn-in period. It is known that the averaged iterate is asymptotically normal:

$$\sqrt{n}(\bar{\theta}_{n_0, n} - \theta^*) \xrightarrow{d} N(0, \Sigma_\infty), \quad n \rightarrow \infty,$$

where the covariance matrix Σ_∞ is minimax-optimal, see [1]. Important research direction for the first-order optimization methods is to obtain the finite-sample bounds of the form

$$\mathbb{E}^{1/2} [\|\bar{\theta}_{n_0, n} - \theta^*\|^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\infty}}{n^{1/2}} + \frac{C(f, d)}{n^{1/2+\delta}} + \mathcal{R}(\|\theta_0 - \theta^*\|, n), \quad (2)$$

where $\mathcal{R}(\|\theta_0 - \theta^*\|, n)$ is a (transient) term corresponding to the initial condition. Our aim is to obtain a version of the above bound with the best-known constant $\delta > 0$. The best known counterpart of (2) is the bound of [2], which obtains the result with $\delta = 1/4$ with the Root-SGD algorithm.

Assumptions

We aim to solve the problem (1) using SGD with a constant step size $\gamma > 0$, starting from initial distribution ν :

$$\theta_{k+1}^{(\gamma)} = \theta_k^{(\gamma)} - \gamma \nabla F(\theta_k^{(\gamma)}, \xi_{k+1}), \quad \theta_0^{(\gamma)} = \theta_0 \sim \nu. \quad (3)$$

We focus on the convergence to θ^* of the Polyak-Ruppert averaged estimator defined for any $n \geq 1$, by

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

A1. The function f is continuously differentiable and μ -strongly convex on \mathbb{R}^d .

A2. The function f is 4 times continuously differentiable and L_2 -smooth on \mathbb{R}^d , that is, there is a constant $L_2 \geq 0$, such that for any $\theta, \theta' \in \mathbb{R}^d$, $\|\nabla f(\theta) - \nabla f(\theta')\| \leq L \|\theta - \theta'\|$. Moreover, f has bounded 3-rd and 4-th derivatives.

A3(p). Random variables $\{\xi_k\}_{k \in \mathbb{N}}$ are independent and identically distributed (i.i.d.) random variables with distribution \mathbb{P}_ξ , such that ξ_i and θ_0 are independent and for any $\theta \in \mathbb{R}^d$ it holds that

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

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

$$L^{q-1} \|\theta_1 - \theta_2\|^{q-2} \langle \nabla f(\theta_1) - \nabla f(\theta_2), \theta_1 - \theta_2 \rangle \geq \mathbb{E}_{\xi \sim \mathbb{P}_\xi} [\|\nabla F(\theta_1, \xi) - \nabla F(\theta_2, \xi)\|^q].$$

Our contributions

- We show that a version of SGD algorithm with constant step size, Polyak-Ruppert averaging, and Richardson-Romberg extrapolation lead to the root-MSE bound (2) with $\delta = 1/4$ when applied to strongly convex minimization problems. This result requires that the number of samples n is known a priori to optimize the step size γ .
- We obtain high-order moment bounds on the error, that is, we obtain for $p \geq 2$ the bounds of the form

$$\mathbb{E}^{1/p} [\|\bar{\theta}_n^{(RR)} - \theta^*\|^p] \leq \frac{c_0 p^{1/2} \sqrt{\text{Tr} \Sigma_\infty}}{n^{1/2}} + \frac{C(f, d, p)}{n^{3/4}} + \mathcal{R}(\|\theta_0 - \theta^*\|, n, p),$$

where c_0 is a universal constant, and $\bar{\theta}_n^{(RR)}$ is a counterpart of $\bar{\theta}_{n_0, n}$ when using Richardson-Romberg extrapolation.

SGD iterates as a Markov chain

- The sequence $\{\theta_k^{(\gamma)}\}$ defined by the relation (3) is a time-homogeneous Markov chain with the Markov kernel

$$Q_\gamma(\theta, A) = \int_{\mathbb{R}^d} 1_A(\theta - \gamma \nabla F(\theta, z)) \mathbb{P}_\xi(dz).$$

- Key ingredient of our result is to show that the Markov chain $\{\theta_k^{(\gamma)}\}$ converges to its invariant distribution π_γ in a Wasserstein distance \mathbf{W}_c , associated with the non-standard distance-like function:

$$c(\theta, \theta') = \|\theta - \theta'\| (\|\theta - \theta^*\| + \|\theta' - \theta^*\| + \frac{2\sqrt{2}\tau_2\sqrt{\gamma}}{\sqrt{\mu}}), \quad \theta, \theta' \in \mathbb{R}^d.$$

Theorem 1. Assume **A1**, **A2**, **A3**(4). Then for n large enough, any initial distribution ν , it holds setting $\gamma = n^{-2/3}$ that

$$\mathbb{E}_\nu^{1/2} [\|\mathbf{H}^*(\bar{\theta}_n^{(\gamma)} - \theta^*)\|^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\infty^*}}{n^{1/2}} + \frac{\mathbf{C}(L, \mu)}{n^{2/3}} + \mathcal{R}(n, 1/n^{2/3}, \|\theta_0 - \theta^*\|),$$

where the remainder term $\mathcal{R}(n, \gamma, \|\theta_0 - \theta^*\|) = (1 - \mu\gamma)^{n/2} \text{poly}(\gamma, n, \|\theta_0 - \theta^*\|)$.

Richardson-Romberg Extrapolation

Instead of considering a single SGD trajectory $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$, construct two parallel chains based on the same sequence $\{\xi_k\}_{k \in \mathbb{N}}$:

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

Then construct the Richardson-Romberg estimator, which allows to reduce the steady-state bias of SGD:

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

Theorem 2. Assume **A1**, **A2**, **A3**(6). Then for n large enough, any initial distribution ν , it holds setting $\gamma = n^{-1/2}$ that

$$\mathbb{E}_\nu^{1/2} [\|\mathbf{H}^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\infty^*}}{n^{1/2}} + \frac{\mathbf{C}(L, \mu)}{n^{3/4}} + \mathcal{R}(n, 1/\sqrt{n}, \|\theta_0 - \theta^*\|),$$

where the remainder term $\mathcal{R}(n, \gamma, \|\theta_0 - \theta^*\|) = (1 - \mu\gamma)^{n/2} \text{poly}(\gamma, n, \|\theta_0 - \theta^*\|)$. Here the power $n^{-3/4}$ is the best known among the first-order methods. This result can be generalized to p -th moment bounds with $p \geq 2$:

Theorem 3. Let $p \geq 2$ and assume **A1**, **A2**, **A3**(3p). Then for n large enough, any initial distribution ν , it holds setting $\gamma = n^{-1/2}$ that

$$\mathbb{E}_\nu^{1/p} [\|\mathbf{H}^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^p] \leq \frac{c_1 \sqrt{\text{Tr} \Sigma_\infty^*} p^{1/2}}{n^{1/2}} + \frac{\mathbf{C}(L, \mu, p)}{n^{3/4}} + \mathcal{R}(n, 1/\sqrt{n}, \|\theta_0 - \theta^*\|),$$

where $c_1 = 60e$ is an absolute constant. Note that the Richardson-Romberg procedure allows for more aggressive choice of step size (γ of order $n^{-1/2}$), compared to $n^{-2/3}$, which corresponds to $\bar{\theta}_n^{(\gamma)}$.

References

- Gersende Fort. Central limit theorems for stochastic approximation with controlled markov chain dynamics. *ESAIM: Probability and Statistics*, 19:60–80, 2015.
- Chris Junchi Li, Wenlong Mou, Martin Wainwright, and Michael Jordan. Root-sgd: Sharp nonasymptotics and asymptotic efficiency in a single algorithm. In Po-Ling Loh and Maxim Raginsky, editors, *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 909–981. PMLR, 02–05 Jul 2022.