



# Gaussian Approximation and Multiplier Bootstrap for Stochastic Gradient Descent

Marina Sheshukova<sup>1</sup>, Sergey Samsonov<sup>1</sup> Denis Belomestny<sup>2,1</sup>, Eric Moulines<sup>3,4</sup>, Qi-Man Shao<sup>5</sup>, Zhuo-Song Zhang<sup>5</sup> Alexey Naumov<sup>1,6</sup>,

<sup>1</sup>HSE University    <sup>2</sup>Duisburg–Essen University    <sup>3</sup>CMAF, UMR 7641, École Polytechnique    <sup>4</sup>Mohamed Bin Zayed University of AI

<sup>5</sup>Southern University of Science and Technology <sup>6</sup>Steklov Mathematical Institute of the Russian Academy of Sciences

## Stochastic Gradient Descent (SGD)

- We aim to estimate  $\theta^* \in \arg \min_{\theta \in \mathbb{R}^d} f(\theta)$  with access only to the noisy gradients  $\nabla F(\theta, \xi)$  such that  $\nabla f(\theta) = \mathbb{E}_{\xi \sim \mathbb{P}_\xi} [\nabla F(\theta, \xi)]$ . Here  $\xi$  is a noise variable with the distribution  $\mathbb{P}_\xi$ . We assume that  $\theta^*$  is the unique minimizer.

- SGD with Polyak–Ruppert (PR) averaging:

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

$$\bar{\theta}_n = n^{-1} \sum_{k=0}^{n-1} \theta_k. \quad (2)$$

- Polyak–Juditsky central limit theorem (see [2]) implies asymptotic normality

$$\sqrt{n}(\bar{\theta}_n - \theta^*) \xrightarrow{d} \mathcal{N}(0, \Sigma_\infty), \quad n \rightarrow \infty, \quad (3)$$

where the covariance matrix  $\Sigma_\infty$  is minimax-optimal.

## Key questions

- What is the rate of convergence in (3)?
- How can (3) be leveraged to construct confidence sets for  $\theta^*$ , given that  $\Sigma_\infty$  is unknown in practice?

To quantify convergence rates in (3), we employ convex distance, which is defined for random vectors  $X, Y \in \mathbb{R}^d$  as

$$\mathbf{d}_C(X, Y) = \sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)|, \quad \text{where } \mathcal{C}(\mathbb{R}^d) \text{ denotes the collection of convex subsets of } \mathbb{R}^d.$$

## Confidence sets based on Multiplier bootstrap procedure

- Let  $\mathcal{W}^{n-1} = \{w_\ell\}_{1 \leq \ell \leq n-1}$  be i.i.d. random variables with distribution  $\mathbb{P}_w$ , each with mean  $\mathbb{E}[w_1] = 1$  and variance  $\text{Var}[w_1] = 1$ . Assume  $\mathcal{W}^{n-1}$  is independent of  $\Xi^{n-1} = \{\xi_\ell\}_{1 \leq \ell \leq n-1}$ . Procedure is based on perturbing the trajectory (1) (see [1])

$$\begin{aligned} \theta_k^b &= \theta_{k-1}^b - \alpha_k w_k \{ \nabla f(\theta_{k-1}^b) + g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k) \}, \quad k \geq 1, \quad \theta_0^b = \theta_0, \\ \bar{\theta}_n^b &= n^{-1} \sum_{k=0}^{n-1} \theta_k^b, \quad n \geq 1. \end{aligned} \quad (4)$$

Note that, when generating different weights  $w_k$ , we can draw samples from the conditional distribution of  $\bar{\theta}_n^b$  given the data  $\Xi^{n-1}$ .

- The core principle: "bootstrap world" probabilities  $\mathbb{P}(\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) \in B | \Xi^{n-1})$  are close to  $\mathbb{P}(\sqrt{n}(\bar{\theta}_n - \theta^*) \in B)$  for  $B \in \mathcal{C}(\mathbb{R}^d)$ .

## Assumptions

**A1.** The function  $f$  is two times continuously differentiable and  $L_1$ -smooth on  $\mathbb{R}^d$ . Moreover, we assume that  $f$  is  $\mu$ -strongly convex on  $\mathbb{R}^d$ .

**A2.** For each  $k \geq 1$ ,  $\zeta_k$  admits the decomposition  $\zeta_k = \eta(\xi_k) + g(\theta_{k-1}, \xi_k)$ , where

- 1  $\{\xi_k\}_{k=1}^{n-1}$  is a sequence of i.i.d. random variables on  $(\mathcal{Z}, \mathcal{Z})$  with distribution  $\mathbb{P}_\xi$ ,  $\eta : \mathcal{Z} \rightarrow \mathbb{R}^d$  is a function such that  $\mathbb{E}[\eta(\xi_1)] = 0$  and  $\mathbb{E}[\eta(\xi_1)\eta(\xi_1)^\top] = \Sigma_\xi$ . Moreover,  $\lambda_{\min}(\Sigma_\xi) > 0$ .

- 2 The function  $g : \mathbb{R}^d \times \mathcal{Z} \rightarrow \mathbb{R}^d$  satisfies  $\mathbb{E}[g(\theta, \xi_1)] = 0$  for any  $\theta \in \mathbb{R}^d$ . Moreover, there exists  $L_2 > 0$  such that for any  $\theta, \theta' \in \mathbb{R}^d$ , it holds that

$$\|g(\theta, \xi) - g(\theta', \xi)\| \leq L_2 \|\theta - \theta'\| \quad \text{and} \quad g(\theta^*, z) = 0 \quad \text{for all } z \in \mathcal{Z}. \quad (5)$$

- 3 There exist  $C_{1,\xi}, C_{2,\xi} > 0$  such that  $\mathbb{P}_\xi$ -almost surely that  $\|\eta(\xi)\| \leq C_{1,\xi}$  and  $\sup_\theta \|g(\theta, \xi)\| \leq C_{2,\xi}$ .

**A3.** There exist  $L_3, \beta > 0$  such that for all  $\theta$  with  $\|\theta - \theta^*\| \leq \beta$ , it holds

$$\|\nabla^2 f(\theta) - \nabla^2 f(\theta^*)\| \leq L_3 \|\theta - \theta^*\|. \quad (6)$$

**A4.** The stochastic gradient  $F(\theta, \xi) := \nabla f(\theta) + g(\theta, \xi) + \eta(\xi)$  is almost surely  $L_4$ -co-coercive, that is, for any  $\theta, \theta' \in \mathbb{R}^d$ , it holds  $\mathbb{P}_\xi$ -almost surely that

$$L_4 \langle F(\theta, \xi) - F(\theta', \xi), \theta - \theta' \rangle \geq \|F(\theta, \xi) - F(\theta', \xi)\|^2. \quad (7)$$

## Assumptions

**A5.** There exist constants  $0 < W_{\min} < W_{\max} < +\infty$ , such that  $W_{\min} \leq w_1 \leq W_{\max}$  a.s.

**A6.** Let  $\alpha_k = c_0 \{k_0 + k\}^{-\gamma}$ , where  $\gamma \in (1/2, 1)$ , an  $c_0$  satisfies  $c_0 W_{\max} \max(2L_4, \mu) \leq 1$  and  $k_0 \geq (\frac{2\gamma}{\mu c_0 W_{\min}})^{1/(1-\gamma)}$ .

**A7.** Number of observations  $n$  is large enough.

## Non-asymptotic multiplier bootstrap validity

**Theorem 1.** Assume **A1** - **A7**. Then with  $\mathbb{P}$  - probability at least  $1 - 1/n$ , it holds

$$\sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}(\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) \in B | \Xi^{n-1}) - \mathbb{P}(\sqrt{n}(\bar{\theta}_n - \theta^*) \in B)| \leq \frac{C_1 \log n}{n^{\gamma-1/2}}, \quad (8)$$

where  $C_1$  are some problem-specific constants.

- The proof is based on the gaussian approximation in real and bootstrap world.
- Key observation: instead of  $\mathcal{N}(0, \Sigma_\infty)$  we use  $\mathcal{N}(0, \Sigma_n)$ , where  $\Sigma_n$  is the covariance of the linearized recursion which correspond to (1) with additive noise  $\eta(\xi_k)$  (see [3] for details).

## Rate of convergence in the Polyak–Juditsky central limit theorem

**A8(p)** Conditions (i) and (ii) from **A2** holds. Moreover, there exists  $\sigma_p > 0$  such that  $\mathbb{E}^{1/p}[\|\eta(\xi_1)\|^p] \leq \sigma_p$ .

**A9** Suppose that  $\alpha_k = c_0/(k_0 + k)^\gamma$ , where  $\gamma \in (1/2, 1)$ ,  $k_0 \geq 1$ , and  $c_0$  satisfies  $2c_0 L_1 \leq 1$ .

**Theorem 2.** Assume **A1**, **A3**, **A8(4)**, **A9**. Then, with  $Y \sim \mathcal{N}(0, I_d)$  it holds that

$$\mathbf{d}_C(\sqrt{n}(\bar{\theta}_n - \theta^*), \Sigma_\infty^{1/2} Y) \leq \frac{C_2}{n^{\gamma-1/2}} + \frac{C_\infty}{n^{1-\gamma}}, \quad (9)$$

where  $C_2$  and  $C_\infty$  are some problem-specific constants.

## Lower bounds

**Theorem 3.** There exists the problem satisfying conditions **A1**, **A3**, **A8(4)**, **A9**, such that with  $Y \sim \mathcal{N}(0, I_d)$  for  $n$  large enough it holds that

$$\mathbf{d}_C(\sqrt{n}(\bar{\theta}_n - \theta^*), \Sigma_\infty^{1/2} Y) \geq \frac{C_3}{n^{1-\gamma}}. \quad (10)$$

The bound (10) implies that  $\sqrt{n}(\bar{\theta}_n - \theta^*)$  cannot be approximated by  $\mathcal{N}(0, \Sigma_\infty)$  faster than  $1/n^{1-\gamma}$ , and the rate in Equation (9) is tight for  $\gamma \in [3/4, 1)$ . This highlights the necessity of using  $\Sigma_n$  in the bootstrap result (Equation (8)).

## Takeaways

- We prove the first non-asymptotic validity result for multiplier bootstrap in averaged SGD, with rate  $n^{-(\gamma-1/2)}$  for  $\gamma \in (1/2, 1)$ ;
- For strongly convex problems, we show a tight Gaussian approximation rate  $n^{-1/4}$  in (3) with  $\gamma = 3/4$ .

## Acknowledgement

This work was supported by the Ministry of Economic Development of the Russian Federation (code 25-139-66879-1-0003).

## References

- [1] Yixin Fang, Jinfeng Xu, and Lei Yang. Online bootstrap confidence intervals for the stochastic gradient descent estimator. *Journal of Machine Learning Research*, 19(78):1–21, 2018.
- [2] Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992.
- [3] Marina Sheshukova, Sergey Samsonov, Denis Belomestny, Eric Moulines, Qi-Man Shao, Zhuo-Song Zhang, and Alexey Naumov. Gaussian approximation and multiplier bootstrap for stochastic gradient descent. *arXiv preprint arXiv:2502.06719*, 2025.