

NONASYMPTOTIC ANALYSIS OF STOCHASTIC GRADIENT DESCENT WITH THE RICHARDSON–ROMBERG EXTRAPOLATION

Marina Sheshukova^{1,5} **Denis Belomestny**^{2,1} **Alain Durmus**³ **Eric Moulines**^{3,4}
Alexey Naumov^{1,6} **Sergey Samsonov**¹

¹HSE University ²Duisburg-Essen University ³CMAP, UMR 7641, Ecole polytechnique

⁴Mohamed Bin Zayed University of AI ⁵Skolkovo Institute of Science and Technology

⁶Steklov Mathematical Institute of Russian Academy of Sciences

{msheshukova, anaumov, svsamsonov}@hse.ru

{alain.durmus, eric.moulines}@polytechnique.edu

denis.belomestny@uni-due.de

ABSTRACT

We address the problem of solving strongly convex and smooth minimization problems using stochastic gradient descent (SGD) algorithm with a constant step size. Previous works suggested to combine the Polyak-Ruppert averaging procedure with the Richardson-Romberg extrapolation to reduce the asymptotic bias of SGD at the expense of a mild increase of the variance. We significantly extend previous results by providing an expansion of the mean-squared error of the resulting estimator with respect to the number of iterations n . We show that the mean-squared error can be decomposed into the sum of two terms: a leading one of order $\mathcal{O}(n^{-1/2})$ with explicit dependence on a minimax-optimal asymptotic covariance matrix, and a second-order term of order $\mathcal{O}(n^{-3/4})$, where the power $3/4$ is best known. We also extend this result to the high-order moment bounds. Our analysis relies on the properties of the SGD iterates viewed as a time-homogeneous Markov chain. In particular, we establish that this chain is geometrically ergodic with respect to a suitably defined weighted Wasserstein semimetric.

1 INTRODUCTION

Stochastic gradient methods are a fundamental approach for solving a wide range of optimization problems, with a broad range of applications including generative modeling (Goodfellow et al., 2014; 2016), empirical risk minimization (Van der Vaart, 2000), and reinforcement learning (Sutton & Barto, 2018; Schulman et al., 2015; Mnih et al., 2015). These methods solve the stochastic minimization problem

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

where ξ is a random variable with distribution \mathbb{P}_ξ , and the gradient ∇f of the function f can be accessed only through (unbiased) noisy estimates ∇F . Throughout this paper, we consider strongly convex minimization problems admitting a unique solution θ^* . Arguably the simplest and one of the most widely used approaches to solve (1) is the stochastic gradient descent (SGD). This algorithm constructs the sequence of updates

$$\theta_{k+1} = \theta_k - \gamma_{k+1} \nabla F(\theta_k, \xi_{k+1}), \quad \theta_0 \in \mathbb{R}^d, \quad (2)$$

where $\{\gamma_k\}_{k \in \mathbb{N}}$ are step sizes, either diminishing or constant, and $\{\xi_k\}_{k \in \mathbb{N}}$ is an i.i.d. sequence with distribution \mathbb{P}_ξ . The algorithm (2) can be viewed as a special instance of the Robbins-Monro procedure (Robbins & Monro, 1951). While the SGD algorithm remains one of the core algorithms in statistical inference, its performance can be enhanced by means of additional techniques that use e.g., momentum (Qian, 1999), averaging (Ruppert, 1988; Polyak & Juditsky, 1992), or variance reduction (Defazio et al., 2014; Nguyen et al., 2017). In particular, the celebrated Polyak-Ruppert

algorithm proceeds with a trajectory-wise averaging of the estimates

$$\bar{\theta}_{n_0, n} = \frac{1}{n} \sum_{k=n_0+1}^{n+n_0} \theta_k \quad (3)$$

for some $n_0 > 0$. It is known (Polyak & Juditsky, 1992; Fort, 2015), that under appropriate assumptions on f and γ_k , the sequence of estimates $\{\bar{\theta}_{n_0, n}\}_{n \in \mathbb{N}}$ is asymptotically normal, that is,

$$\sqrt{n}(\bar{\theta}_{n_0, n} - \theta^*) \xrightarrow{d} \mathcal{N}(0, \Sigma_\infty), \quad n \rightarrow \infty \quad (4)$$

where \xrightarrow{d} denotes the convergence in distribution and $\mathcal{N}(0, \Sigma_\infty)$ denotes the zero-mean Gaussian distribution with covariance matrix Σ_∞ , which is asymptotically optimal, see Fort (2015) for a discussion. On the other hand, quantitative counterparts of (4) rely on the root mean squared error (root-MSE) bounds of the form

$$\mathbb{E}^{1/2}[\|\bar{\theta}_{n_0, n} - \theta^*\|^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\infty}}{n^{1/2}} + \frac{C(f, d)}{n^{1/2+\delta}} + \mathcal{R}(\|\theta_0 - \theta^*\|, n). \quad (5)$$

Here $\mathcal{R}(\|\theta_0 - \theta^*\|, n)$ is a remainder term which reflects the dependence upon initial condition, $C(f, d)$ is some instance-dependent constant and $\delta > 0$. There are many studies establishing (5) for Polyak-Ruppert averaged SGD under various model assumptions, including Moulines & Bach (2011); Gadat & Panloup (2023). In particular, Li et al. (2022) derived the bound (5) with the rate $\delta = 1/4$, which is best known among the first-order methods. However, their results apply to a modified two-timescale algorithm with multiple restarts. In our work, we show that the same non-asymptotic upper bound is achieved by a simple modification of the estimate $\bar{\theta}_{n_0, n}$ based on Richardson-Romberg extrapolation. The main contributions of the current paper are as follows:

- 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 (5) with $\delta = 1/4$ when applied to strongly convex minimization problems. We obtain this result by leveraging the analysis of iterates of the constant step-size SGD as a Markov chain. It is important to note that this result is obtained for a fixed step size γ of order $1/\sqrt{n}$ with n being a total number of iterations. This result requires that the number of samples, n , is known a priori to optimize the step size γ .
- We generalize the above result and obtain high-order moment bounds on the error. Under a similar step size $\gamma = 1/\sqrt{n}$, we obtain for $p \geq 2$ bounds of the form

$$\mathbb{E}^{1/p}[\|\bar{\theta}_n^{(RR)} - \theta^*\|^p] \leq \frac{c_0 p^{1/2} \sqrt{\text{Tr} \Sigma_\infty}}{n^{1/2}} + \frac{C(f, d, p)}{n^{3/4}} + \mathcal{R}(\|\theta_0 - \theta^*\|, n, p), \quad (6)$$

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, see related definitions in Section 4. Our proof is based on a novel version of the Rosenthal inequality, which might be of independent interest.

The rest of the paper is organized as follows. First, we provide a literature review on the non-asymptotic analysis of the first order optimization methods, with an emphasis on the constant step-size algorithms and Richardson-Romberg procedure in Section 2. Then, we provide analysis of the constant step size SGD viewed as a Markov chain together with the properties of the Polyak-Ruppert averaged estimator (3) in Section 3. In Section 4, we discuss the Richardson-Romberg extrapolation applied to the Polyak-Ruppert averaged SGD and derive the respective 2-nd and p -th moment error bounds.

Notations and definitions. For $\theta_1, \dots, \theta_k$ being the iterates of stochastic first-order method, we denote $\mathcal{F}_k = \sigma(\theta_0, \theta_1, \dots, \theta_k)$. Let (Z, d_Z) be a complete separable metric space equipped with its Borel σ -algebra \mathcal{Z} . We call a function $c : Z \times Z \rightarrow \mathbb{R}_+$ a *distance-like* function, if it is symmetric, lower semi-continuous and $c(x, y) = 0$ if and only if $x = y$, and there exists $q \in \mathbb{N}$ such that $(d(x, y) \wedge 1)^q \leq c(x, y)$. We denote by $\mathcal{C}(\xi, \xi')$ the set of couplings of probability measures ξ and ξ' , that is, a set of probability measures on $(Z^2, \mathcal{Z}^{\otimes 2})$, such that for any $\Pi \in \mathcal{C}(\xi, \xi')$ and any $A \in \mathcal{Z}$ it holds $\Pi(Z \times A) = \xi'(A)$ and $\Pi(A \times Z) = \xi(A)$. We define the Wasserstein semimetric associated to the distance-like function $c(\cdot, \cdot)$, as

$$\mathbf{W}_c(\xi, \xi') = \inf_{\Pi \in \mathcal{C}(\xi, \xi')} \int_{Z \times Z} c(z, z') \Pi(dz, dz'). \quad (7)$$

Note that $\mathbf{W}_c(\xi, \xi')$ is not necessarily a distance, as it may fail to satisfy the triangle inequality. In the particular case of $Z = \mathbb{R}^d$, and $c_p(x, y) = \|x - y\|^p$, $x, y \in \mathbb{R}^d$, $p \geq 1$, we use a short notation for $\mathbf{W}_p(\xi, \xi')$ defined by $\mathbf{W}_p^p(\xi, \xi') = \mathbf{W}_{c_p}(\xi, \xi')$. For $x, y \in \mathbb{R}^d$ denote by $x \otimes y$ the tensor product of x and y and by $x^{\otimes k}$ the k -th tensor power of x . In addition, for a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ we denote by $\nabla^k f(\theta)$ the k -th differential of f , that is $\nabla^k f(\theta)_{i_1, \dots, i_k} = \frac{\partial^k f}{\partial \theta_{i_1} \dots \partial \theta_{i_k}}$. For any tensor $M \in (\mathbb{R}^d)^{\otimes(k-1)}$, we define $\nabla^k f(\theta)M \in \mathbb{R}^d$ by the relation $(\nabla^k f(\theta)M)_l = \sum_{i_1, \dots, i_{k-1}} M_{i_1, \dots, i_{k-1}} \nabla^k f(\theta)_{i_1, \dots, i_{k-1}, l}$, where $l \in \{1, \dots, d\}$. Also for any tensor $M \in (\mathbb{R}^d)^{\otimes(k-1)}$ we define $\|M\| = \sup_{x^l \neq 0, l \in \{1, \dots, k\}} \frac{\sum_{i_1, \dots, i_k} M_{i_1, \dots, i_k} x_{i_1}^1 \dots x_{i_k}^k}{\|x^1\| \dots \|x^k\|}$. For two sequences $\{a_n\}_{n \in \mathbb{N}}$ and $\{b_n\}_{n \in \mathbb{N}}$ we write $a_n \lesssim b_n$, if there is an absolute constant $c > 0$, such that $a_n \leq cb_n$ for any $n \in \mathbb{N}$. Throughout this paper we use c_0 for an absolute constant, which values may vary from line to line.

2 LITERATURE REVIEW

Constant step-size SGD has been widely studied in literature. Its bias and MSE were studied for strongly convex problems in (Dieuleveut et al., 2020), both for the last iterate and under the Polyak-Ruppert averaging. Yu et al. (2021) consider bias and asymptotic normality of the SGD's last iterate for non-convex problems under Polyak-Lojasiewicz condition and its generalizations. Vlatakis-Gkaragkounis et al. (2024) consider constant step-size methods for solving variational inequalities, focusing on characterizing the bias and establishing asymptotic normality. Merad & Gaïffas (2023) study convergence of constant step-size SGD iterates in \mathbf{W}_p distance, $p \geq 1$, and provide concentration bounds, both for θ_n and $\bar{\theta}_{n_0, n}$.

Moulines & Bach (2011) derive the bound (5) with $\delta = 1/6$ for the case of strongly convex functions f . Gadat & Panloup (2023) obtained an MSE counterpart of (5) of the form. Li et al. (2022) suggested the Root-SGD algorithm combining the ideas of the two-timescale stochastic approximation and using multiple restarts, establishing a counterpart of (5) with $\delta = 1/4$. The recent series of papers (Huo et al., 2023; Zhang & Xie, 2024; Zhang et al., 2024) investigate stochastic approximation algorithms with both i.i.d. and Markovian data and constant step sizes. The authors consider both linear SA problems and Q -learning, quantify bias, and propose precise characterization of the bias together with a Richardson-Romberg extrapolation procedure. However, these results only consider 2-nd moment of the error and provide MSE bounds of order $\mathcal{O}(1/n) + \mathcal{O}(\gamma)$ with no explicit expression for the leading term.

Richardson-Romberg extrapolation. Richardson-Romberg (RR) extrapolation is a technique used to improve the accuracy of numerical approximations (Hildebrand, 1987), such as those from numerical differentiation or integration. It involves using approximations with different step sizes and then extrapolating to reduce the error, typically by removing the leading term in the error expansion. The one-step RR extrapolation was introduced to reduce the discretization error induced by an Euler scheme to simulate stochastic differential equation in Talay & Tubaro (1990), and later generalized for non-smooth functions in Bally & Talay (1996). This technique was extended using multistep discretizations in Pagès (2007). RR extrapolation have been applied to Stochastic Gradient Descent (SGD) methods in Durmus et al. (2016), Merad & Gaïffas (2023) and Huo et al. (2024a), to improve convergence and reduce error in optimization problems, particularly when dealing with noisy or high-variance gradient estimates. Recent papers (Allmeier & Gast, 2024; Zhang & Xie, 2024; Huo et al., 2024a; Kwon et al., 2024) consider applications of RR to different stochastic approximation problems with constant step-size, including Q -learning, and single- and two-timescale stochastic approximation.

3 FINITE-TIME ANALYSIS OF THE SGD DYNAMICS FOR STRONGLY CONVEX MINIMIZATION PROBLEMS

3.1 GEOMETRIC ERGODICITY OF SGD ITERATES

We consider the following assumption on the function f in the minimization problem (1).

A1. The function f is continuously differentiable and μ -strongly convex on \mathbb{R}^d , that is, there exists a constant $\mu > 0$, such that for any $\theta, \theta' \in \mathbb{R}^d$, it holds that

$$\frac{\mu}{2} \|\theta - \theta'\|^2 \leq f(\theta) - f(\theta') - \langle \nabla f(\theta'), \theta - \theta' \rangle. \quad (8)$$

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 $\theta, \theta' \in \mathbb{R}^d$,

$$\|\nabla f(\theta) - \nabla f(\theta')\| \leq L_2 \|\theta - \theta'\|. \quad (9)$$

Moreover, f has bounded 3-rd and 4-th derivatives, that is, there exist $L_3, L_4 \geq 0$ such that

$$\|\nabla^i f(\theta)\| \leq L_i \text{ for } i \in \{3, 4\}. \quad (10)$$

We aim to solve the problem (1) using SGD with a constant step size, starting from initial distribution ν . That is, for $k \geq 0$ and a step size $\gamma \geq 0$, we consider the following recurrent scheme

$$\theta_{k+1}^{(\gamma)} = \theta_k^{(\gamma)} - \gamma \nabla F(\theta_k^{(\gamma)}, \xi_{k+1}), \quad \theta_0^{(\gamma)} = \theta_0 \sim \nu, \quad (11)$$

where $\{\xi_k\}_{k \in \mathbb{N}}$ is a sequence satisfying the following condition.

A3 (p). $\{\xi_k\}_{k \in \mathbb{N}}$ is a sequence of 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 $L_1 > 0$ that for any $\theta_1, \theta_2 \in \mathbb{R}^d$,

$$L_1^{q-1} \|\theta_1 - \theta_2\|^{q-2} \langle \nabla f(\theta_1) - \nabla f(\theta_2), \theta_1 - \theta_2 \rangle \geq \mathbb{E}_{\xi \sim \mathbb{P}_\xi} [\|\nabla F(\theta_1, \xi) - \nabla F(\theta_2, \xi)\|^q]. \quad (12)$$

Assumption **A3(p)** generalizes the well-known L_1 -co-coercivity assumption, see [Dieuleveut et al. \(2020\)](#). A sufficient condition which allows for **A3(p)** is to assume that $F(\theta, \xi)$ is \mathbb{P}_ξ -a.s. convex with respect to $\theta \in \mathbb{R}^d$. For ease of notation, we set

$$L = \max(L_1, L_2, L_3, L_4), \quad (13)$$

and trace only dependence upon L in our subsequent bounds. In this paper we focus on the convergence to θ^* of the Polyak-Ruppert averaging estimator defined for any $n \geq 0$,

$$\bar{\theta}_n^{(\gamma)} = \frac{1}{n} \sum_{k=n+1}^{2n} \theta_k^{(\gamma)}. \quad (14)$$

Many previous studies instead consider $\bar{\vartheta}_n^{(\gamma)} = \frac{1}{n-n_0} \sum_{k=n_0+1}^n \theta_k^{(\gamma)}$ rather than $\bar{\theta}_n^{(\gamma)}$, where $n \geq n_0 + 1$ and n_0 denotes a burn-in period. However, when the sample size n is sufficiently large, the choice of the optimal burn-in size n_0 affects the leading terms in the MSE bound of $\bar{\theta}_n^{(\gamma)} - \theta^*$ only by a constant factor. Therefore, we focus on (14), or equivalently, use $2n$ observations and set $n_0 = n$.

Properties of $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ viewed as a Markov chain. Under assumptions **A1**, **A2** and **A3(2)**, the sequence $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ defined by the relation (11) is a time-homogeneous Markov chain with the Markov kernel

$$Q_\gamma(\theta, A) = \int_{\mathbb{R}^d} \mathbb{1}_A(\theta - \gamma \nabla F(\theta, z)) \mathbb{P}_\xi(dz), \quad \theta \in \mathbb{R}^d, A \in \mathcal{B}(\mathbb{R}^d), \quad (15)$$

where $\mathcal{B}(\mathbb{R}^d)$ is a Borel σ -field of \mathbb{R}^d . In [Dieuleveut et al. \(2020\)](#) it has been established that, under the stated assumptions, Q_γ admits a unique invariant distribution π_γ , if γ is small enough. Previous studies, such as [Dieuleveut et al. \(2020\)](#) or [Merad & Gaïffas \(2023\)](#), studied the convergence of the distributions of $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ to π_γ in the 2-Wasserstein distance \mathbf{W}_2 , associated with the Euclidean distance in \mathbb{R}^d . Our main results require to switch to the non-standard distance-like function, which is defined under **A1** and **A3(2)** as follows:

$$c(\theta, \theta') = \|\theta - \theta'\| \left(\|\theta - \theta^*\| + \|\theta' - \theta^*\| + \frac{2\sqrt{2}\tau_2\sqrt{\gamma}}{\sqrt{\mu}} \right), \quad \theta, \theta^* \in \mathbb{R}^d. \quad (16)$$

Here the constants τ_2 and μ are given in **A3(2)** and **A1**, respectively. This distance-like function is designed to analyze $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ under **A1** and **A3(2)**. In particular, it depends on the step size γ and θ^* . Our first main result establishes *geometric ergodicity* of the Markov kernel Q_γ with respect to the distance-like function c from (16).

Proposition 1. Assume A1, A2, and A3(2). Then for any $\gamma \in (0; 1/(2L)]$, the Markov kernel Q_γ defined in (15) admits a unique invariant distribution π_γ . Moreover, Q_γ is geometrically ergodic with respect to the cost function c , that is, for any initial distribution ν on \mathbb{R}^d and $k \in \mathbb{N}$,

$$\mathbf{W}_c(\nu Q_\gamma^k, \pi_\gamma) \leq 4(1/2)^{k/m(\gamma)} \mathbf{W}_c(\nu, \pi_\gamma), \quad (17)$$

where $m(\gamma) = \lceil 2 \log 4/(\gamma\mu) \rceil$.

Discussion. The proof of Proposition 1 is provided in Appendix A.1. Properties of the invariant distribution π_γ were previously studied in literature, see e.g. Dieuleveut et al. (2020). In particular, it is known (Dieuleveut et al., 2020, Lemma 13), that the 2-nd moment of π_γ scales linearly with γ :

$$\int_{\mathbb{R}^d} \|\theta - \theta^*\|^2 \pi_\gamma(d\theta) \lesssim \frac{\gamma\tau_2}{\mu}. \quad (18)$$

This bound yields, using Lyapunov’s inequality, that

$$\int_{\mathbb{R}^d \times \mathbb{R}^d} \|\theta - \theta'\| \pi_\gamma(d\theta) \pi_\gamma(d\theta') \lesssim \sqrt{\frac{\gamma\tau_2}{\mu}}.$$

At the same time, expectation of the cost function $c(\theta, \theta')$ scales linearly with the step size γ :

$$\int_{\mathbb{R}^d} c(\theta, \theta') \pi_\gamma(d\theta) \pi_\gamma(d\theta') \lesssim \frac{\gamma\tau_2}{\mu}. \quad (19)$$

The property (19) is crucial to obtain tighter (with respect to the step size γ) error bounds for the Richardson-Romberg estimator, as well as in the Rosenthal inequality for additive functional of $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ derived in Proposition 8. Precisely, the additional $\sqrt{\gamma}$ factor obtained in (19) would allow us to obtain sharper bounds on the remainder terms in Theorem 6. Next, we analyze the error $\theta_\infty^{(\gamma)} - \theta^*$ where $\theta_\infty^{(\gamma)}$ is distributed according to the stationary distribution π_γ . To this end, we consider the following condition.

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 $\gamma \in (0, 1/(L C_{\text{step},p})]$, and any initial distribution ν it holds that

$$\mathbb{E}_\nu^{2/p} [\|\theta_k^{(\gamma)} - \theta^*\|^p] \leq (1 - \gamma\mu)^k \mathbb{E}_\nu^{2/p} [\|\theta_0 - \theta^*\|^p] + D_{\text{last},p} \gamma \tau_p^2 / \mu. \quad (20)$$

Moreover, for the stationary distribution π_γ it holds that

$$\mathbb{E}_{\pi_\gamma}^{2/p} [\|\theta_\infty^{(\gamma)} - \theta^*\|^p] \leq D_{\text{last},p} \gamma \tau_p^2 / \mu. \quad (21)$$

It is important to note that C1 is not independent from the preceding assumptions A1 - A3(p). In particular, Dieuleveut et al. (2020, Lemma 13) establishes that, under A1, A2, and A3(p), the bound (21) holds for $\gamma \in (0, 1/(L C_{\text{step},p})]$ with some constants $D_{\text{last},p}$ and $C_{\text{step},p}$ depending only upon p . Unfortunately, it is difficult to obtain precise dependence of $C_{\text{step},p}$ and $D_{\text{last},p}$ on p , and to obtain (21) with precise numerical constants. Existing studies (Gadat & Panloup, 2023; Merad & Gaïffas, 2023) either use different set of assumptions or do not explicitly characterize their dependence on p . That is why we prefer to state C1(p) as a separate assumption. In the subsequent bounds we use C1(p) together with A1, A2, and A3(p), tracking the dependence of our bounds upon $C_{\text{step},p}$ and $D_{\text{last},p}$. We leave the derivation of C1(p) with precise constants $D_{\text{last},p}, C_{\text{step},p}$ as an interesting direction for future research.

Under assumption C1, we control the fluctuations of $\theta_k^{(\gamma)}$ around θ^* . However, unless f is quadratic, it is known that $\int_{\mathbb{R}^d} \theta \pi_\gamma(d\theta) \neq \theta^*$. In the next proposition, we quantify this bias under weaker assumptions than those in Dieuleveut et al. (2020, Theorem 4).

Proposition 2. Assume A1, A2, A3(6), and C1(6). Then there exist such $\Delta_1 \in \mathbb{R}^d, \Delta_2 \in \mathbb{R}^{d \times d}$, not depending upon γ , that for any $\gamma \in (0, 1/(L C_{\text{step},6})]$, it holds

$$\bar{\theta}_\gamma := \int_{\mathbb{R}^d} \theta \pi_\gamma(d\theta) = \theta^* + \gamma \Delta_1 + B_1 \gamma^{3/2}, \quad (22)$$

$$\bar{\Sigma}_\gamma := \int_{\mathbb{R}^d} (\theta - \theta^*)^{\otimes 2} \pi_\gamma(d\theta) = \gamma \Delta_2 + B_2 \gamma^{3/2}. \quad (23)$$

Here $B_1 \in \mathbb{R}^d$ and $B_2 \in \mathbb{R}^{d \times d}$ satisfy $\|B_1\| \leq \frac{L}{\mu} C_1 + \frac{LD_{\text{last},6}^3 \tau_6^3}{2\mu^{5/2}}$, $\|B_2\| \leq C_1$, where C_1 defined in (51) is a constant independent of γ . Moreover, for any initial distribution ν on \mathbb{R}^d , it holds that

$$\mathbb{E}_\nu[\bar{\theta}_n^{(\gamma)}] = \theta^* + \gamma \Delta_1 + B_1 \gamma^{3/2} + \mathcal{R}_1(\theta_0 - \theta^*, \gamma, n), \quad (24)$$

where

$$\|\mathcal{R}_1(\theta_0 - \theta^*, \gamma, n)\| \lesssim \frac{e^{-\gamma\mu(n+1)/2}}{n\gamma\mu} \left(\mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^2] + \frac{\sqrt{\gamma\tau_2}}{\sqrt{\mu}} \right). \quad (25)$$

The proof is provided in Appendix A. Results of this type are known in the literature for stochastic approximation algorithms, see e.g. Huo et al. (2024b) and Allmeier & Gast (2024). The additive term Δ_1 vanishes when the function f is quadratic, see (Moulines & Bach, 2011).

3.2 ANALYSIS OF THE POLYAK-RUPPERT AVERAGED ESTIMATOR $\bar{\theta}_n^{(\gamma)}$.

In this section, we analyze the finite-sample properties of the estimator $\bar{\theta}_n^{(\gamma)}$ from (14). The analysis is based on techniques previously used in Moulines & Bach (2011), as well as in the analysis of the Polyak-Ruppert averaged LSA (Linear Stochastic Approximation) algorithms, see Mou et al. (2020); Durmus et al. (2024). Below we derive the key representation for the error $\bar{\theta}_n^{(\gamma)} - \theta^*$, following Moulines & Bach (2011). Define the k -th step noise level at the point $\theta \in \mathbb{R}^d$ by:

$$\varepsilon_k(\theta) = \nabla F(\theta, \xi_k) - \nabla f(\theta), \quad (26)$$

and $\varepsilon_{k+1}(\theta_k^{(\gamma)})$ is a martingale-difference sequence w.r.t. $(\mathcal{F}_k)_{k \in \mathbb{N}}$. Then the recurrence (11) takes form

$$\theta_{k+1}^{(\gamma)} - \theta^* = \theta_k^{(\gamma)} - \theta^* - \gamma(\nabla f(\theta_k^{(\gamma)}) + \varepsilon_{k+1}(\theta_k^{(\gamma)})). \quad (27)$$

We set

$$\eta(\theta) = \nabla f(\theta) - H^*(\theta - \theta^*), \text{ where } H^* = \nabla^2 f(\theta^*) \in \mathbb{R}^{d \times d}. \quad (28)$$

We obtain from (27) with simple algebra that

$$H^*(\theta_k^{(\gamma)} - \theta^*) = \frac{\theta_k^{(\gamma)} - \theta_{k+1}^{(\gamma)}}{\gamma} - \varepsilon_{k+1}(\theta_k^{(\gamma)}) - \eta(\theta_k^{(\gamma)}). \quad (29)$$

Taking average of (29) for $k = n+1$ to $2n$, we arrive at the final representation:

$$H^*(\bar{\theta}_n^{(\gamma)} - \theta^*) = \frac{\theta_{n+1}^{(\gamma)} - \theta^*}{\gamma n} - \frac{\theta_{2n}^{(\gamma)} - \theta^*}{\gamma n} - \frac{1}{n} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta_k^{(\gamma)}) - \frac{1}{n} \sum_{k=n+1}^{2n} \eta(\theta_k^{(\gamma)}). \quad (30)$$

Equation (30) is the key ingredient in the proof of the next result where the variance of noise $\varepsilon_k(\theta^*)$ measured at the optimal point θ^* naturally appears, that is,

$$\Sigma_\varepsilon^* = \mathbb{E}_{\xi \sim \mathbb{P}_\xi}[\nabla F(\theta^*, \xi)^{\otimes 2}]. \quad (31)$$

Note that Σ_ε^* does not depend on the step size γ of (11) and is related to the "optimal" covariance matrix of the Polyak-Ruppert averaged iterates $\bar{\theta}_n^{(\gamma)}$, see Fort (2015). In our first main result below, we establish non-asymptotic properties of the averaged Polyak-Ruppert estimator (14).

Theorem 3. Assume A1, A2, A3(6), and C1(6). Then for any $\gamma \in (0, 1/(L C_{\text{step},6})]$, $n \in \mathbb{N}$, and initial distribution ν on \mathbb{R}^d , the sequence of Polyak-Ruppert estimates (14) satisfies

$$\mathbb{E}_\nu^{1/2}[\|H^*(\bar{\theta}_n^{(\gamma)} - \theta^*)\|^2] \leq \frac{\sqrt{\text{Tr } \Sigma_\varepsilon^*}}{\sqrt{n}} + \frac{C_2}{\gamma^{1/2}n} + C_3\gamma + \frac{C_4\gamma^{1/2}}{n^{1/2}} + \mathcal{R}_2(n, \gamma, \|\theta_0 - \theta^*\|), \quad (32)$$

where the constants C_2 to C_4 are defined in Appendix B (see equation (65)), and

$$\begin{aligned} \mathcal{R}_2(n, \gamma, \|\theta_0 - \theta^*\|) &= \frac{c_0(1 - \gamma\mu)^{(n+1)/2} L}{\gamma\mu n} \mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^2] \\ &\quad + \frac{c_0 L(1 - \gamma\mu)^{n+1}}{2n\gamma\mu} \mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^4], \end{aligned}$$

where c_0 is an absolute (numerical) constant.

The version of Theorem 3 with explicit constants together with the proof is provided in Appendix B, see Theorem 15. Note that the result of Theorem 3 is valid for arbitrary $\gamma \in (0, 1/(L C_{\text{step},6})]$. At the same time, this bound can be optimized over step size of the form $\gamma = n^{-\beta}$, $\beta \in (0, 1)$.

Corollary 4. *Under the assumptions of Theorem 3, provided that $n \geq (L C_{\text{step},6})^{3/2}$, it holds setting $\gamma = n^{-2/3}$ that*

$$\mathbb{E}_\nu^{1/2}[\|\mathbf{H}^*(\bar{\theta}_n^{(\gamma)} - \theta^*)\|^2] \leq \frac{\sqrt{\text{Tr } \Sigma_\varepsilon^*}}{n^{1/2}} + \frac{C(L, \mu)}{n^{2/3}} + \mathcal{R}_2(n, 1/n^{2/3}, \|\theta_0 - \theta^*\|), \quad (33)$$

where the expression for $C(L, \mu)$ can be traced from Appendix B, eq. (65).

Corollary 4 implies that, if n is known in advance and $\gamma = n^{-2/3}$, then $\bar{\theta}_n^{(\gamma)}$ satisfies (5) with $\delta = 1/6$. A closer inspection of the sum (30) reveals that $\mathbb{E}_{\pi_\gamma}[\eta(\theta_k^\gamma)] \asymp \gamma$, and we can not expect to provide a better bound for the term $\frac{1}{n} \sum_{k=n+1}^{2n} \eta(\theta_k^\gamma)$ compared to the one coming from the Minkowski's inequality. Thus, this is the *bias* of the stationary distribution, which prevents us from gaining the optimal second-order term w.r.t. the sample size n from Corollary 4.

Note that in case of deterministic problems $\varepsilon_k(\theta) = 0$ for any k and θ , and $\mathbf{C1}(p)$ is satisfied for any $p \geq 2$ with $D_{\text{last},p} = 0$. In such a setting, $\Sigma_\varepsilon^* = 0$, and the remainder terms are proportional to $D_{\text{last},p}$ with $p = 2, 4$, or 6, and therefore also vanishes. Therefore, Theorem 3 provides exponential convergence bounds, which are embedded in the remainder term. The decay rate of the second-order (w.r.t. n) term in (33) is well studied in the literature. In particular, Moulines & Bach (2011) obtained the bound of the same order $\tilde{O}(n^{-2/3})$ for the second-order term of the root-MSE bound of SGD algorithm with Polyak-Ruppert averaging. A similar rate under more general assumptions was reported in (Gadat & Panloup, 2023, Theorem 2). However, all these results are known to be suboptimal for first-order methods. The recent work by Li et al. (2022) shows that the best known second-order error term in the bounds (33) is of order $\mathcal{O}(n^{-3/4})$ and can be achieved by the Root-SGD algorithm. In the next section we mirror this bound using the constant step-size SGD algorithm combined with the Richardson-Romberg extrapolation technique.

4 RICHARDSON-ROMBERG EXTRAPOLATION

Our analysis presented in Theorem 3 was based on the summation by parts formula (30) and Taylor expansion of the gradient $\nabla f(\theta)$ in the vicinity of θ^* , yielding the remainder quantity $\eta(\theta)$. It is important to notice that

$$\int_{\mathbb{R}^d} \eta(\theta) \pi_\gamma(d\theta) \neq 0, \quad (34)$$

which prevents us from using more aggressive (larger) step sizes γ in the optimized bound (33). In this section we show that Richardson-Romberg extrapolation technique is sufficient to significantly reduce the bias associated with $\eta(\theta)$ and improve the second-order term in the MSE bound (33). Instead of considering a single SGD trajectory $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$, and then relying on the tail-averaged estimator $\bar{\theta}_n^{(\gamma)}$, we construct two parallel chains based on the same sequence $\{\xi_k\}_{k \in \mathbb{N}}$:

$$\begin{aligned} \theta_{k+1}^{(\gamma)} &= \theta_k^{(\gamma)} - \gamma \nabla F(\theta_k^{(\gamma)}, \xi_{k+1}), & \bar{\theta}_n^{(\gamma)} &= \frac{1}{n} \sum_{k=n+1}^{2n} \theta_k^{(\gamma)}, \\ \theta_{k+1}^{(2\gamma)} &= \theta_k^{(2\gamma)} - 2\gamma \nabla F(\theta_k^{(2\gamma)}, \xi_{k+1}), & \bar{\theta}_n^{(2\gamma)} &= \frac{1}{n} \sum_{k=n+1}^{2n} \theta_k^{(2\gamma)}. \end{aligned} \quad (35)$$

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

$$\bar{\theta}_n^{(RR)} := 2\bar{\theta}_n^{(\gamma)} - \bar{\theta}_n^{(2\gamma)}. \quad (36)$$

Note that it is possible to use different sources of randomness $\{\xi_k\}_{k \in \mathbb{N}}$ and $\{\xi'_k\}_{k \in \mathbb{N}}$ when constructing the sequences $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ and $\{\theta_k^{(2\gamma)}\}_{k \in \mathbb{N}}$, respectively. At the same time, it is possible to show the benefits of using the same sequence of random variables $\{\xi_k\}_{k \in \mathbb{N}}$ in (35). Indeed, consider the decomposition (30) and further expand the term $\eta(\theta)$ defined in (28) as

$$\eta(\theta) = \psi(\theta) + G(\theta),$$

where we have defined, for $\theta \in \mathbb{R}^d$, the following vector-valued functions:

$$\begin{aligned}\psi(\theta) &= (1/2) \nabla^3 f(\theta^*)(\theta - \theta^*)^{\otimes 2}, \\ G(\theta) &= (1/2) \left(\int_0^1 t^2 \nabla^4 f(t\theta^* + (1-t)\theta) dt \right) (\theta - \theta^*)^{\otimes 3}.\end{aligned}\quad (37)$$

We further rewrite the decomposition (30) as

$$\begin{aligned}\mathbf{H}^*(\bar{\theta}_n^{(\gamma)} - \theta^*) &= \frac{\theta_{n+1}^{(\gamma)} - \theta^*}{\gamma n} - \frac{\theta_{2n}^{(\gamma)} - \theta^*}{\gamma n} - \frac{1}{n} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*) \\ &\quad - \frac{1}{n} \sum_{k=n+1}^{2n} \{\varepsilon_{k+1}(\theta_k^{(\gamma)}) - \varepsilon_{k+1}(\theta^*)\} - \frac{1}{n} \sum_{k=n+1}^{2n} \psi(\theta_k^{(\gamma)}) - \frac{1}{n} \sum_{k=n+1}^{2n} G(\theta_k^{(\gamma)}).\end{aligned}\quad (38)$$

Note that in the decomposition (38), the linear statistics $W = n^{-1} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*)$ does not depend upon γ . Moreover, when setting the step size $\gamma \simeq n^{-\beta}$ with an appropriate $\beta \in (0, 1)$, we can show that the moments of all other terms except for W in the r.h.s. of (38) are small (see Theorem 9 for more details). Hence, using the same sequence $\{\xi_k\}_{k \in \mathbb{N}}$ of noise variables in (35) yields an estimator $\bar{\theta}_n^{(RR)}$, such that its leading component of the variance still equals W . Hence, using the Richardson-Romberg procedure will increase only the second-order (w.r.t. n) components of the variance. At the same time, using different random sequences $\{\xi_k\}_{k \in \mathbb{N}}$ and $\{\xi'_k\}_{k \in \mathbb{N}}$ for $\bar{\theta}_n^{(\gamma)}$ and $\bar{\theta}_n^{(2\gamma)}$ increase the leading component of the MSE by a constant factor. Hence, it is preferable to use synchronous noise construction as introduced in (35).

Proposition 2 implies the following improved bound on the bias of $\bar{\theta}_n^{(RR)}$:

Proposition 5. Assume A1, A2, A3(6), and C1(6). Then, for any $\gamma \in (0, 1/(\mathbf{L} C_{\text{step},6})]$, and any initial distribution ν on \mathbb{R}^d , it holds that

$$\mathbb{E}_\nu[\bar{\theta}_n^{(RR)}] = \theta^* + B_3 \gamma^{3/2} + \mathcal{R}_3(\theta_0 - \theta^*, \gamma, n), \quad (39)$$

where $B_3 \in \mathbb{R}^d$ is a vector such that $\|B_3\| \leq \frac{L}{\mu} C_1 + \frac{L D_{\text{last},6}^{3/2} \tau_6^3}{2\mu^{5/2}}$, and

$$\|\mathcal{R}_3(\theta_0 - \theta^*, \gamma, n)\| \lesssim \frac{e^{-\gamma\mu(n+1)/2}}{n\gamma\mu} \left(\mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^2] + \frac{\sqrt{\gamma}\tau_2}{\sqrt{\mu}} \right).$$

The proof of Proposition 5 is provided in Appendix A. This result is a simple consequence of Proposition 5, since the linear in γ component of the bias $\gamma\Delta_1$ from (24) cancels out when computing $\bar{\theta}_n^{(RR)}$. We are now ready to formulate the main result for the Richardson-Romberg estimate $\bar{\theta}_n^{(RR)}$.

Theorem 6. Assume A1, A2, A3(6), and C1(6). Then for any $\gamma \in (0, 1/(\mathbf{L} C_{\text{step},6})]$, initial distribution ν and $n \in \mathbb{N}$, the Richardson-Romberg estimator $\bar{\theta}_n^{(RR)}$ defined in (36) satisfies

$$\begin{aligned}\mathbb{E}_\nu^{1/2}[\|\mathbf{H}^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^2] &\leq \frac{\sqrt{\text{Tr} \Sigma_\varepsilon^*}}{n^{1/2}} + \frac{C_{\text{RR},1} \gamma^{1/2}}{n^{1/2}} + \frac{C_{\text{RR},2}}{\gamma^{1/2} n} + C_{\text{RR},3} \gamma^{3/2} + \frac{C_{\text{RR},4} \gamma}{n^{1/2}} \\ &\quad + \mathcal{R}_4(n, \gamma, \|\theta_0 - \theta^*\|),\end{aligned}$$

where the constants $C_{\text{RR},1}$ to $C_{\text{RR},4}$ are defined in Appendix C (equation (73)), and

$$\begin{aligned}\mathcal{R}_4(n, \gamma, \|\theta_0 - \theta^*\|) &= \frac{c_0 \mathbf{L} (1 - \gamma\mu)^{(n+1)/2}}{n\gamma\mu} \\ &\quad \times \left(\mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^6] + \mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^4] + \mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^2] + \frac{D_{\text{last},4} \gamma \tau_4^2}{\mu} \right),\end{aligned}$$

with c_0 being an absolute constant.

Proof of Theorem 6 is provided in Appendix C. Similarly to Theorem 3, we can optimize the above bound setting γ depending upon n .

Corollary 7. Under the assumptions of Theorem 6, provided that $n \geq (\mathbf{L} C_{\text{step},6})^2$, it holds setting $\gamma = n^{-1/2}$ that

$$\mathbb{E}_\nu^{1/2}[\|\mathbf{H}^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\varepsilon^*}}{n^{1/2}} + \frac{C(\mathbf{L}, \mu)}{n^{3/4}} + \mathcal{R}_4(n, 1/\sqrt{n}, \|\theta_0 - \theta^*\|), \quad (40)$$

where the expression for $C(\mathbf{L}, \mu)$ can be traced from Appendix C, eq. (73).

Discussion. Note that the result of Corollary 7 is a counterpart of (5) with $\delta = 1/4$. This decay rate of the second order term is the same as for the Root-SGD algorithm of Li et al. (2022). However, the assumptions of Theorem 6 are stronger compared to the ones imposed by Li et al. (2022). In particular, in A2 we require that f is 4 times continuously differentiable and uniformly bounded. At the same time, Li et al. (2022) impose Lipschitz continuity of the Hessian of f . Our proof of Theorem 6 essentially relies on the 4-th order Taylor expansion, and it is not clear, if this assumption can be relaxed. We leave further investigations of this question for future research.

Now we generalize the previous result for the p -th moment bounds with $p \geq 2$. 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 8. *Let $p \geq 2$ and assume A1, A2, A3(2p), and C1(2p). Then for ψ defined in (37) and any $\gamma \in (0, 1/(L C_{\text{step}, 2p}))$, it holds that*

$$\mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=0}^{n-1} \{\psi(\theta_k^{(\gamma)}) - \pi_\gamma(\psi)\} \right\|^p \right] \lesssim \frac{L D_{\text{last}, 2p} p \tau_{2p}^2 \sqrt{n\gamma}}{\mu^{3/2}} + \frac{L D_{\text{last}, 2p} \tau_{2p}}{\mu^2}. \quad (41)$$

Discussion. Proof of Proposition 8 is provided in Appendix D.1. It is important to acknowledge that there are numerous Rosenthal-type inequalities for dependent sequences in the literature. Proposition 8 can be viewed as an analogue to the classical Rosenthal inequality for strongly mixing sequences, see (Rio, 2017, Theorem 6.3). However, it should be emphasized that the Markov chain $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ is geometrically ergodic under the assumptions A1-A3(p) only in sense of the weighted Wasserstein semi-metric $\mathbf{W}_c(\xi, \xi')$ with respect to a cost function c defined in (16). As a result, the sequence $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ does not necessarily satisfy strong mixing conditions. At the same time, $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ satisfies the τ -mixing condition, see Merlevède et al. (2011). However, the considered function $\psi(\theta)$ is quadratic, which means that the result of (Merlevède et al., 2011, Theorem 1) can not be directly applied. Bounds similar to (41) have been explored in (Durmus et al., 2023), but in Proposition 8 we obtain the bound with tighter dependence of the right-hand side upon γ . Below we provide the p -th moment bound together with corollary for the step size γ optimized w.r.t. n .

Theorem 9. *Let $p \geq 2$ and assume A1, A2, A3(3p), and C1(3p). Then for any step size $\gamma \in (0, 1/(L C_{\text{step}, 3p}))$, initial distribution ν , and $n \in \mathbb{N}$, the estimator $\bar{\theta}_n^{(RR)}$ defined in (36) satisfies*

$$\begin{aligned} \mathbb{E}_\nu^{1/p} [\|H^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^p] &\leq \frac{c_1 \sqrt{\text{Tr } \Sigma_\varepsilon^*} p^{1/2}}{n^{1/2}} + \frac{C_{\text{RR}, 5}}{n\gamma^{1/2}} + \frac{C_{\text{RR}, 6} \gamma^{1/2}}{n^{1/2}} + C_{\text{RR}, 7} \gamma^{3/2} \\ &\quad + \frac{c_2 p \tau_p}{n^{1-1/p}} + \frac{C_{\text{RR}, 8}}{n} + \mathcal{R}_5(n, \gamma, \|\theta_0 - \theta^*\|), \end{aligned} \quad (42)$$

where $c_1 = 60e$ and $c_2 = 60$ are absolute constants from the Pinelis version of Rosenthal inequality (Pinelis, 1994, Theorem 4.1), problem-specific constants $C_{\text{RR}, 5}$ to $C_{\text{RR}, 8}$ are defined in Appendix D (equation (101)), and

$$\mathcal{R}_5(n, \gamma, \|\theta_0 - \theta^*\|) = (1 - \gamma\mu)^{(n+1)/2} C_{f,p} (\mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^p] + \mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^{2p}] + \mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^{3p}]).$$

Here constant $C_{f,p}$ can be traced from (102).

Corollary 10. *Under the assumptions of Theorem 9, provided that $n \geq (L C_{\text{step}, 3p})^2$, it holds setting $\gamma = n^{-1/2}$ that*

$$\mathbb{E}_\nu^{1/p} [\|H^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^p] \lesssim \frac{\sqrt{\text{Tr } \Sigma_\varepsilon^*} p^{1/2}}{n^{1/2}} + \frac{C(L, \mu, p)}{n^{3/4}} + \mathcal{R}_5(n, 1/\sqrt{n}, \|\theta_0 - \theta^*\|), \quad (43)$$

where the expression for $C(L, \mu, p)$ can be traced from Appendix D, eq. (101).

Discussion. Proof of Theorem 9 is provided in Appendix D. Note that the result above is a direct generalization of Theorem 6, which reveals the same scaling of the step size γ with respect to n . To the best of our knowledge, this is the first analysis of a first-order method, which provides a bound for the second-order term of order $\mathcal{O}(n^{-3/4})$ while keeping the precise leading term related to the minimax-optimal covariance matrix Σ_ε^* . Such results were previously known only for the setting of linear stochastic approximation (LSA), which corresponds to the case of quadratic function f in the initial minimization problem (1), see Mou et al. (2020); Durmus et al. (2024). Thus, Richardson-Romberg extrapolation applied to strongly convex minimization problems allows to restore the p -th moment error moment bounds as they were obtained in the LSA setting.

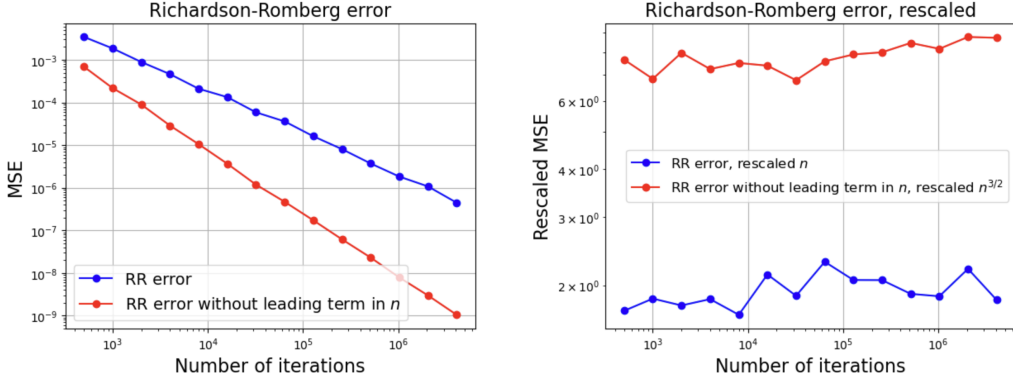


Figure 1: Left picture: Richardson-Romberg experimental error with and without the leading term $\frac{1}{n} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*)$. Right picture: same errors after rescaling by n and $n^{3/2}$, respectively.

5 NUMERICAL RESULTS

In this section we illustrate numerically the scale of the second-order terms in equation (40) in Corollary 7. We show that, for a particular minimization problem, setting $\gamma = n^{-1/2}$, we achieve the scaling of the second-order terms in root-MSE bounds of order $\mathcal{O}(n^{-3/4})$. We consider the problem

$$\min_{\theta \in \mathbb{R}} f(\theta), \quad f(\theta) = \theta^2 + \cos \theta,$$

with the stochastic gradient oracles $\nabla F(\theta, \xi)$ given by $\nabla F(\theta, \xi) = 2\theta - \sin \theta + \xi$, and $\xi \sim \mathcal{N}(0, 1)$. This example satisfies the assumptions A1, A2, A3(p) with any $p \geq 2$. We select different sample sizes n , choose $\gamma = 1/\sqrt{n}$, and construct the associated estimates $\bar{\theta}_n^{(\gamma)}$ and $\bar{\theta}_n^{(2\gamma)}$. Detailed description of the experimental setting is provided in Appendix E. Then for each n we compute the Richardson-Romberg estimates $\bar{\theta}_n^{(RR)}$ from (35) alongside with its versions without the leading term in n , i.e. $\bar{\theta}_n^{(RR)} + n^{-1} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*)$. We provide first the plot for $\|\bar{\theta}_n^{(RR)} - \theta^*\|^2$ and $\|\bar{\theta}_n^{(RR)} + n^{-1} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*) - \theta^*\|^2$, averaged over M parallel runs, in Figure 1. On the same figure we also provide the plots for rescaled errors

$$n\|\bar{\theta}_n^{(RR)} - \theta^*\|^2 \text{ and } n^{3/2}\|\bar{\theta}_n^{(RR)} + n^{-1} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*) - \theta^*\|^2,$$

also averaged over M parallel runs. The corresponding plot indicates that the proper scaling of the MSE of the remainder part is $n^{-3/2}$, that is, corresponding root-MSE scales as $\mathcal{O}(n^{-3/4})$, as predicted by Corollary 7.

6 CONCLUSION

In this paper, we have demonstrated that Polyak-Ruppert averaged SGD iterates with a constant step size achieve optimal root-MSE and p -th moment error bounds. More precisely, we have shown that these bounds admit both the sharp, optimal leading term, which aligns with the optimal covariance matrix Σ_∞ , and a second-order term of order $\mathcal{O}(n^{-3/4})$, which is best known among the first order methods. Directions for future research include, firstly, generalizing the proposed algorithm to the setting of dependent noise sequences $\{\xi_k\}_{k \in \mathbb{N}}$ in the stochastic gradients (1). Another natural question is to study the properties of $\bar{\theta}_n^{(RR)}$ under relaxed assumptions on f . In particular, it would be interesting to remove additional smoothness assumptions on f (bounded 3-rd and 4-th derivatives), and to relax the strong convexity condition A1. One more research direction is to study a relation between the parameter δ in (5) and the Berry-Esseen type bounds for the distribution of $\sqrt{n}(\bar{\theta}_{n_0, n} - \theta^*)$, following the technique of Shao & Zhang (2022).

7 ACKNOWLEDGEMENTS

The work of M. Sheshukova, A. Naumov, and S. Samsonov was prepared within the framework of the HSE University Basic Research Program. The work of E. Moulines has been partly funded by the European Union (ERC-2022-SYG-OCEAN-101071601). Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them.

REFERENCES

- Sebastian Allmeier and Nicolas Gast. Computing the Bias of Constant-step Stochastic Approximation with Markovian Noise. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 137873–137902. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/f949c1f490beb42124a267b7476cd353-Paper-Conference.pdf.
- V. Bally and D. Talay. The law of the euler scheme for stochastic differential equations. *Probability Theory and Related Fields*, 104(1):43–60, 1996. doi: 10.1007/BF01303802. URL <https://doi.org/10.1007/BF01303802>.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. *Advances in neural information processing systems*, 27, 2014.
- 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, 2020. doi: 10.1214/19-AOS1850. URL <https://doi.org/10.1214/19-AOS1850>.
- 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.
- Alain Durmus, Umut Simsekli, Eric Moulines, Roland Badeau, and Gaël Richard. Stochastic Gradient Richardson-Romberg Markov Chain Monte Carlo. *Advances in neural information processing systems*, 29, 2016.
- Alain Durmus, Eric Moulines, Alexey Naumov, Sergey Samsonov, and Marina Sheshukova. Rosenthal-type inequalities for linear statistics of Markov chains. *arXiv preprint arXiv:2303.05838*, 2023.
- Alain Durmus, Eric Moulines, Alexey Naumov, and Sergey Samsonov. Finite-time high-probability bounds for Polyak-Ruppert averaged iterates of linear stochastic approximation. *Mathematics of Operations Research*, 2024.
- 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>.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- Ian J. Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, Cambridge, MA, USA, 2016. <http://www.deeplearningbook.org>.
- Francis Begnaud Hildebrand. *Introduction to numerical analysis*. Courier Corporation, 1987.

- Dongyan Huo, Yudong Chen, and Qiaomin Xie. Bias and extrapolation in Markovian linear stochastic approximation with constant stepsizes. In *Abstract Proceedings of the 2023 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, pp. 81–82, 2023.
- Dongyan Lucy Huo, Yudong Chen, and Qiaomin Xie. Effectiveness of constant stepsize in markovian lsa and statistical inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 20447–20455, 2024a.
- Dongyan Lucy Huo, Yixuan Zhang, Yudong Chen, and Qiaomin Xie. The collusion of memory and nonlinearity in stochastic approximation with constant stepsize. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 21699–21762. Curran Associates, Inc., 2024b.
- Jeongyeol Kwon, Luke Dotson, Yudong Chen, and Qiaomin Xie. Two-Timescale Linear Stochastic Approximation: Constant Stepsizes Go a Long Way. *arXiv preprint arXiv:2410.13067*, 2024.
- 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 (eds.), *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pp. 909–981. PMLR, 02–05 Jul 2022. URL <https://proceedings.mlr.press/v178/li22a.html>.
- Ibrahim Merad and Stéphane Gaïffas. Convergence and concentration properties of constant step-size sgd through markov chains. *arXiv preprint arXiv:2306.11497*, 2023.
- Florence Merlevède, Magda Peligrad, and Emmanuel Rio. A Bernstein type inequality and moderate deviations for weakly dependent sequences. *Probability Theory and Related Fields*, 151(3-4): 435–474, 2011.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Wenlong Mou, Chris Junchi Li, Martin J Wainwright, Peter L Bartlett, and Michael I Jordan. On linear stochastic approximation: Fine-grained Polyak-Ruppert and non-asymptotic concentration. In *Conference on Learning Theory*, pp. 2947–2997. PMLR, 2020.
- Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. *Advances in neural information processing systems*, 24, 2011.
- Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method for machine learning problems using stochastic recursive gradient. In *International Conference on Machine Learning*, pp. 2613–2621. PMLR, 2017.
- A. Osekowski. *Sharp Martingale and Semimartingale Inequalities*. Monografie Matematyczne 72. Birkhäuser Basel, 1 edition, 2012. ISBN 3034803699,9783034803694.
- Gilles Pagès. Multi-step Richardson-Romberg Extrapolation: Remarks on Variance Control and Complexity. *Monte Carlo Methods and Applications*, 13(1):37–70, 2007. doi: doi:10.1515/MCMA.2007.003. URL <https://doi.org/10.1515/MCMA.2007.003>.
- 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>.
- Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992.
- Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- Emmanuel Rio. *Asymptotic Theory of Weakly Dependent Random Processes*, volume 80. 2017. ISBN 978-3-662-54322-1. doi: 10.1007/978-3-662-54323-8.

- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- 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>.
- David Ruppert. Efficient estimations from a slowly convergent Robbins-Monro process. Technical report, Cornell University Operations Research and Industrial Engineering, 1988.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015.
- Qi-Man Shao and Zhuo-Song Zhang. Berry–Esseen bounds for multivariate nonlinear statistics with applications to M-estimators and stochastic gradient descent algorithms. *Bernoulli*, 28(3): 1548–1576, 2022.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018. URL <http://incompleteideas.net/book/the-book-2nd.html>.
- Denis Talay and Luciano Tubaro. Expansion of the global error for numerical schemes solving stochastic differential equations. *Stochastic Analysis and Applications*, 8(4):483–509, 1990. doi: 10.1080/07362999008809220. URL <https://doi.org/10.1080/07362999008809220>.
- Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- Emmanouil Vasileios Vlatakis-Gkaragkounis, Angeliki Giannou, Yudong Chen, and Qiaomin Xie. Stochastic methods in variational inequalities: Ergodicity, bias and refinements. In *International Conference on Artificial Intelligence and Statistics*, pp. 4123–4131. PMLR, 2024.
- Lu Yu, Krishnakumar Balasubramanian, Stanislav Volgushev, and Murat A Erdogdu. An analysis of constant step size SGD in the non-convex regime: Asymptotic normality and bias. *Advances in Neural Information Processing Systems*, 34:4234–4248, 2021.
- Yixuan Zhang and Qiaomin Xie. Constant stepsize q-learning: Distributional convergence, bias and extrapolation. *arXiv preprint arXiv:2401.13884*, 2024.
- Yixuan Zhang, Dongyan Huo, Yudong Chen, and Qiaomin Xie. Prelimit coupling and steady-state convergence of constant-stepsize nonsmooth contractive SA. *arXiv preprint arXiv:2303.05838*, 2024. URL <https://arxiv.org/abs/2404.06023>.

A PROOF OF PROPOSITION 1, PROPOSITION 2, AND PROPOSITION 5

Throughout this appendix we use c_0 for an absolute constant, which values may vary from line to line. In addition, when the upper index of $\theta_k^{(\gamma)}$ or $\theta_k^{(2\gamma)}$ is omitted, we assume the result applies to iterations of $\theta_k^{(\gamma)}$. The corresponding results for $\theta_k^{2\gamma}$ can be obtained by substituting 2γ instead of γ . We also provide some additional definitions related to the Markov kernels and kernel couplings, particularly useful when considering convergence in Wasserstein semimetric. Detailed exposition can be found in (Douc et al., 2018, Chapter 20).

Let $Q(z, A)$ be a Markov kernel on (Z, \mathcal{Z}) . We say that K is a Markov coupling of Q if for all $(z, z') \in Z^2$ and $A \in \mathcal{Z}$, $K((z, z'), A \times Z) = Q(z, A)$ and $K((z, z'), Z \times A) = Q(z', A)$. If K is a kernel coupling of Q , then for all $n \in \mathbb{N}$, K^n is a kernel coupling of Q^n and for any $\mathcal{C} \in \mathcal{C}(\xi, \xi')$, CK^n is a coupling of $(\xi Q^n, \xi' Q^n)$. Moreover, it holds that

$$\mathbf{W}_c(\xi Q^n, \xi' Q^n) \leq \int_{Z \times Z} K^n c(z, z') \mathcal{C}(dz dz'), \quad (44)$$

see (Douc et al., 2018, Corollary 20.1.4). For any probability measure \mathcal{C} on $(Z^2, \mathcal{Z}^{\otimes 2})$, we denote by $\mathbb{P}_{\mathcal{C}}^K$ and $\mathbb{E}_{\mathcal{C}}^K$ the probability and the expectation on the canonical space $((Z^2)^{\mathbb{N}}, (\mathcal{Z}^{\otimes 2})^{\otimes \mathbb{N}})$ such that the canonical process $\{(Z_n, Z'_n), n \in \mathbb{N}\}$ is a Markov chain with initial probability \mathcal{C} and Markov kernel K . We write $\mathbb{E}_{z, z'}^K$ instead of $\mathbb{E}_{\delta_{z, z'}}^K$.

To prove Proposition 1 we need the following auxiliary lemma about the last iterate of SGD algorithm. It can be found in (Dieuleveut et al., 2020, Lemma 10), but we provide its proof here for completeness.

Lemma 11. *Assume A1, A2, and A3(2). Then for any $\gamma \in (0; 1/(2L)]$ and any $k, r \in \mathbb{N}$ it holds that*

$$\mathbb{E}^{1/2}[\|\theta_{k+r} - \theta^*\|^2 | \mathcal{F}_k] \leq (1 - \gamma\mu)^{r/2} \|\theta_k - \theta^*\| + \frac{2^{1/2} \gamma^{1/2} \tau_2}{\mu^{1/2}} \quad (45)$$

Proof. Using the recurrence (11), we get

$$\begin{aligned} \mathbb{E}[\|\theta_{k+1} - \theta^*\|^2 | \mathcal{F}_k] &= \mathbb{E}[\|\theta_k - \theta^* - \gamma \nabla F(\theta_k, \xi_{k+1})\|^2 | \mathcal{F}_k] \\ &= \mathbb{E}[\|\theta_k - \theta^*\|^2 - 2\gamma \langle \nabla F(\theta_k, \xi_{k+1}), \theta_k - \theta^* \rangle + \gamma^2 \|F(\theta_k, \xi_{k+1})\|^2 | \mathcal{F}_k]. \end{aligned}$$

Applying A3(2), we get

$$\mathbb{E}[\|\theta_{k+1} - \theta^*\|^2 | \mathcal{F}_k] \leq \|\theta_k - \theta^*\|^2 - 2\gamma \langle \nabla f(\theta_k) - \nabla f(\theta^*), \theta_k - \theta^* \rangle + \gamma^2 \mathbb{E}[\|F(\theta_k, \xi_{k+1})\|^2 | \mathcal{F}_k].$$

Since $\nabla f(\theta^*) = 0$, using A3(2), we obtain

$$\begin{aligned} \mathbb{E}[\|F(\theta_k, \xi_{k+1})\|^2 | \mathcal{F}_k] &= \mathbb{E}[\|F(\theta_k, \xi_{k+1}) - F(\theta^*, \xi_{k+1}) + \varepsilon_{k+1}(\theta^*)\|^2 | \mathcal{F}_k] \\ &\leq 2L \langle \nabla f(\theta_k) - \nabla f(\theta^*), \theta_k - \theta^* \rangle + 2\tau_2^2. \end{aligned}$$

Using A1, A2, and the fact that $\gamma \leq 1/(2L)$, we get

$$\mathbb{E}[\|\theta_{k+1} - \theta^*\|^2 | \mathcal{F}_k] \leq (1 - 2\gamma\mu(1 - L\gamma)) \|\theta_k - \theta^*\|^2 + 2\gamma^2 \tau_2^2 \leq (1 - \gamma\mu) \|\theta_k - \theta^*\|^2 + 2\gamma^2 \tau_2^2.$$

Hence, applying tower property for conditional expectations, we obtain

$$\mathbb{E}[\|\theta_{k+r} - \theta^*\|^2 | \mathcal{F}_k] \leq (1 - \gamma\mu)^r \|\theta_k - \theta^*\|^2 + 2\gamma^2 \tau_2^2 \sum_{i=0}^{r-1} (1 - \gamma\mu)^i \leq (1 - \gamma\mu)^r \|\theta_k - \theta^*\|^2 + \frac{2\gamma\tau_2^2}{\mu}.$$

□

A.1 PROOF OF PROPOSITION 1

Consider the synchronous coupling construction defined by the recursions

$$\begin{aligned} \theta_{k+1}^{(\gamma)} &= \theta_k^{(\gamma)} - \gamma \nabla F(\theta_k^{(\gamma)}, \xi_{k+1}), & \theta_0^{(\gamma)} &= \theta \in \mathbb{R}^d, \\ \tilde{\theta}_{k+1}^{(\gamma)} &= \tilde{\theta}_k^{(\gamma)} - \gamma \nabla F(\tilde{\theta}_k^{(\gamma)}, \xi_{k+1}), & \tilde{\theta}_0^{(\gamma)} &= \tilde{\theta} \in \mathbb{R}^d. \end{aligned} \quad (46)$$

The pair $(\theta_k^{(\gamma)}, \tilde{\theta}_k^{(\gamma)})_{k \in \mathbb{N}}$ defines a Markov chain with the Markov kernel $K_\gamma(\cdot, \cdot)$, which is a coupling kernel of (Q_γ, Q_γ) . From now on we omit an upper index (γ) and write simply $(\theta_k, \tilde{\theta}_k)_{k \in \mathbb{N}}$. Applying now A3(2), we get for $\gamma \leq 1/L$ that

$$\begin{aligned} \mathbb{E}[\|\theta_{k+1} - \tilde{\theta}_{k+1}\|^2 | \mathcal{F}_k] &= \mathbb{E}[\|\theta_k - \tilde{\theta}_k - \gamma(\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}))\|^2 | \mathcal{F}_k] \\ &= \|\theta_k - \tilde{\theta}_k\|^2 + \gamma^2 \mathbb{E}[\|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^2 | \mathcal{F}_k] \\ &\quad - 2\gamma \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \\ &\leq (1 - \gamma\mu) \|\theta_k - \tilde{\theta}_k\|^2, \end{aligned} \quad (47)$$

where in the last inequality we additionally used $1 - 2\gamma\mu(1 - \gamma L/2) \leq 1 - \gamma\mu$. Similarly, for a cost function c defined in (16), we get using Hölder's and Minkowski's inequalities, that for any $r \in \mathbb{N}$

$$\begin{aligned} \mathbb{E}[c(\theta_{k+r}, \tilde{\theta}_{k+r}) | \mathcal{F}_k] &\leq \mathbb{E}^{1/2}[\|\theta_{k+r} - \tilde{\theta}_{k+r}\|^2 | \mathcal{F}_k] (\mathbb{E}^{1/2}[\|\theta_{k+r} - \theta^*\|^2 | \mathcal{F}_k] \\ &\quad + \mathbb{E}^{1/2}[\|\tilde{\theta}_{k+r} - \theta^*\|^2 | \mathcal{F}_k] + \frac{2^{3/2}\gamma^{1/2}\tau_2}{\mu^{1/2}}). \end{aligned}$$

Combining the above inequalities and applying Lemma 11, we obtain

$$\begin{aligned} \mathbb{E}[c(\theta_{k+r}, \theta'_{k+r}) | \mathcal{F}_k] &\leq (1 - \gamma\mu)^{r/2} \|\theta_k - \tilde{\theta}_k\| ((1 - \gamma\mu)^{r/2} (\|\theta_k - \theta^*\| + \|\tilde{\theta}_k - \theta^*\|) + \frac{2^{5/2}\gamma^{1/2}\tau_2}{\mu^{1/2}}) \\ &\leq 2(1 - \gamma\mu)^{r/2} c(\theta_k, \theta'_k). \end{aligned}$$

Note that $2(1 - \gamma\mu)^{r/2} \leq 2$ for any $r \leq m(\gamma) - 1$ and $2(1 - \gamma\mu)^{m(\gamma)/2} \leq 1/2$. Hence, applying the result of (Douc et al., 2018, Theorem 20.3.4), we obtain that the Markov kernel Q_γ admits a unique invariant distribution π_γ . Applying (44), we get

$$\mathbf{W}_c(\nu Q_\gamma^k, \pi_\gamma) \leq 2(1/2)^{\lfloor k/m(\gamma) \rfloor} \mathbf{W}_c(\nu, \pi_\gamma). \quad (48)$$

It remains to note that $(1/2)^{\lfloor k/m(\gamma) \rfloor} \leq 2(1/2)^{k/m(\gamma)}$, and the statement follows.

A.2 PROOF OF PROPOSITION 2

We first prove (22) and (23) and introduce some additional notations. Under A1 – A3(2), we define a matrix-valued function $\mathcal{C}(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ as

$$\mathcal{C}(\theta) = \mathbb{E}[\varepsilon_1(\theta)^{\otimes 2}]. \quad (49)$$

The result below is essentially based on an appropriate modification of the bounds presented in Dieuleveut et al. (2020, Lemma 18). A careful inspection of its proof reveals that we do not need additional assumptions on $\mathcal{C}(\theta)$, instead we use Lemma 14.

Lemma 12. Assume A1, A2, A3(6), and C1(6). Then, for any $\gamma \in (0, 1/(L C_{\text{step},6})]$, it holds

$$\bar{\theta}_\gamma - \theta^* = -(\gamma/2) \{H^*\}^{-1} \{\nabla^3 f(\theta^*)\} \mathbf{T} \mathcal{C}(\theta^*) + B_1 \gamma^{3/2}, \quad (50)$$

where $\bar{\theta}_\gamma$ is defined in (22), $\mathcal{C}(\theta)$ is defined in (49), and $B_1 \in \mathbb{R}^d$ satisfies $\|B_1\| \leq \frac{L}{\mu} C_1 + \frac{L D_{\text{last},6}^{3/2} \tau_6^3}{2\mu^{5/2}}$, with

$$C_1 = \frac{\sqrt{L} \tau_2^2}{\sqrt{C_{\text{step},6}} \mu} + \frac{1}{2} \left(\left(\frac{L^2 D_{\text{last},2}}{\mu^{3/2}} + \frac{L \sqrt{D_{\text{last},2}}}{\sqrt{\mu}} \right) \tau_2^2 + \frac{L D_{\text{last},6}^{3/2} \tau_6^3}{2\mu^{3/2}} + \frac{L^{1/2} D_{\text{last},4}^2 \tau_4^4}{4\mu^2 C_{\text{step},6}^{3/2}} \right) \quad (51)$$

Moreover,

$$\bar{\Sigma}_\gamma = \gamma \mathbf{T} \mathcal{C}(\theta^*) + B_2 \gamma^{3/2}, \quad (52)$$

where the operator $\mathbf{T} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$ is defined by the relation

$$\text{vec}(\mathbf{T}A) = (H^* \otimes I + I \otimes H^*)^{-1} \text{vec}(A)$$

for any matrix $A \in \mathbb{R}^{d \times d}$, and $B_2 \in \mathbb{R}^{d \times d}$ is a matrix, such that $\|B_2\| \leq C_1$.

Proof. Let $(\theta_k^{(\gamma)})_{k \in \mathbb{N}}$ be a recurrence defined in (11) with initial distribution $\theta_0 \sim \pi_\gamma$. Recall that θ_0 is independent from the noise variables $(\xi_k)_{k \geq 1}$. First, applying a third-order Taylor expansion of $\nabla f(\theta)$ around θ^* , we obtain

$$\nabla f(\theta) = H^*(\theta - \theta^*) + (1/2)\{\nabla^3 f(\theta^*)\}(\theta - \theta^*)^{\otimes 2} + G(\theta), \quad (53)$$

where $G(\theta)$ has a form

$$G(\theta) = \frac{1}{2} \left(\int_0^1 t^2 \nabla^4 f(t\theta^* + (1-t)\theta) dt \right) (\theta - \theta^*)^{\otimes 3}.$$

Thus, using A2,

$$\|G(\theta)\| \leq \frac{L_4}{2} \|\theta - \theta^*\|^3.$$

Integrating (53) with respect to π_γ , we get

$$H^*(\bar{\theta}_\gamma - \theta^*) + (1/2)\{\nabla^3 f(\theta^*)\} \left[\int_{\mathbb{R}^d} (\theta - \theta^*)^{\otimes 2} \pi_\gamma(d\theta) \right] = - \int_{\mathbb{R}^d} G(\theta) \pi_\gamma(d\theta). \quad (54)$$

Moreover, using C1(6), we have

$$\left\| \int_{\mathbb{R}^d} G(\theta) \pi_\gamma(d\theta) \right\| \leq \gamma^{3/2} \frac{L D_{\text{last},6}^{3/2} \tau_6^3}{2\mu^{3/2}}. \quad (55)$$

Now we provide an explicit expression for the covariance matrix

$$\bar{\Sigma}_\gamma = \int_{\mathbb{R}^d} (\theta - \theta^*)^{\otimes 2} \pi_\gamma(d\theta). \quad (56)$$

Using the recurrence (11), we obtain that

$$\theta_1 - \theta^* = (I - \gamma H^*)(\theta_0 - \theta^*) - \gamma \varepsilon_1(\theta_0) - \gamma \eta(\theta_0),$$

where the function $\eta(\cdot)$ is defined in (28). Hence, taking second moment w.r.t. π_γ from both sides, we get that

$$\begin{aligned} \bar{\Sigma}_\gamma &= (I - \gamma H^*) \bar{\Sigma}_\gamma (I - \gamma H^*) + \gamma^2 \int_{\mathbb{R}^d} \mathcal{C}(\theta) \pi_\gamma(d\theta) + \gamma^2 \int_{\mathbb{R}^d} \{\eta(\theta)\}^{\otimes 2} \pi_\gamma(d\theta) \\ &\quad - \gamma \int_{\mathbb{R}^d} [(I - \gamma H^*)(\theta - \theta^*)\{\eta(\theta)\}^\top + \eta(\theta)(\theta - \theta^*)^\top (I - \gamma H^*)] \pi_\gamma(d\theta). \end{aligned} \quad (57)$$

In the above equation $\mathcal{C}(\theta)$ is defined in (49), and we additionally used that $\mathbb{E}[\varepsilon_1(\theta_0)|\mathcal{F}_0] = 0$. Using Taylor's expansion with integral remainder together with A2 and C1(6),

$$\begin{aligned} \gamma^2 \left\| \int_{\mathbb{R}^d} \{\eta(\theta)\}^{\otimes 2} \pi_\gamma(d\theta) \right\|_F &\leq \gamma^4 \frac{L^2 D_{\text{last},4}^2 \tau_4^4}{4\mu^2}, \\ \gamma \left\| \int_{\mathbb{R}^d} [(I - \gamma H^*)(\theta - \theta^*)\{\eta(\theta)\}^\top + \eta(\theta)(\theta - \theta^*)^\top (I - \gamma H^*)] \pi_\gamma(d\theta) \right\|_F &\leq \gamma^{5/2} \frac{L D_{\text{last},6}^{3/2} \tau_6^3}{2\mu^{3/2}} \end{aligned}$$

Moreover, (49) together with C1(6) imply that

$$\int_{\mathbb{R}^d} \mathcal{C}(\theta) \pi_\gamma(d\theta) = \mathcal{C}(\theta^*) + B\gamma^{1/2},$$

where $B \in \mathbb{R}^{d \times d}$ satisfies $\|B\| \leq C_2$. Using (57) together with C1(6), we obtain that $\bar{\Sigma}_\gamma$ is a solution to the matrix equation

$$H^* \bar{\Sigma}_\gamma + \bar{\Sigma}_\gamma H^* - \gamma H^* \bar{\Sigma}_\gamma H^* = \gamma \mathcal{C}(\theta^*) + B' \gamma^{3/2}, \quad (58)$$

where

$$\|B'\|_F \leq C_2 + \frac{L D_{\text{last},6}^{3/2} \tau_6^3}{2\mu^{3/2}} + \frac{L^{1/2} D_{\text{last},4}^2 \tau_4^4}{4\mu^2 C_{\text{step},6}^{3/2}}. \quad (59)$$

The matrix equation (58) can be written using vectorization operation as

$$\begin{aligned} \text{vec}(\bar{\Sigma}_\gamma) &= \gamma(\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*)^{-1} \text{vec}(\mathcal{C}(\theta^*)) \\ &\quad + \gamma^{3/2}(\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*)^{-1} \text{vec}(B'). \end{aligned}$$

Applying Lemma 13(c), we obtain that

$$(\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*)^{-1} = (\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^*)^{-1} + D,$$

where $D \in \mathbb{R}^{d^2 \times d^2}$ is a matrix which satisfies

$$\|D\| \leq \gamma L / \mu.$$

Thus,

$$\begin{aligned} \text{vec}(\bar{\Sigma}_\gamma) &= \gamma(\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^*)^{-1} \text{vec}(\mathcal{C}(\theta^*)) + \gamma^{3/2}(D/\sqrt{\gamma}) \text{vec}(\mathcal{C}(\theta^*)) \\ &\quad + \gamma^{3/2}(\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*)^{-1} \text{vec}(B'). \end{aligned}$$

We define the matrix B_2 such that

$$\text{vec}(B_2) = (D/\sqrt{\gamma}) \text{vec}(\mathcal{C}(\theta^*)) + (\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*)^{-1} \text{vec}(B') \quad (60)$$

Hence, using A3, (59), and Lemma 13, we get

$$\begin{aligned} \|B_2\| &\leq \|B_2\|_F = \|\text{vec}(B_2)\| \\ &\leq \|D/\sqrt{\gamma}\| \|\text{vec}(\mathcal{C}(\theta^*))\| + \|(\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*)^{-1}\| \|B'\|_F \\ &\leq \frac{\sqrt{\gamma} L \tau_2^2}{\mu} + \frac{1}{2} \left(C_2 + \frac{L D_{\text{last},6}^{3/2} \tau_6^3}{2\mu^{3/2}} + \frac{L^{1/2} D_{\text{last},4}^2 \tau_4^4}{4\mu^2 C_{\text{step},6}^{3/2}} \right) \\ &\leq \frac{\sqrt{L} \tau_2^2}{\sqrt{C_{\text{step},6}} \mu} + \frac{1}{2} \left(C_2 + \frac{L D_{\text{last},6}^{3/2} \tau_6^3}{2\mu^{3/2}} + \frac{L^{1/2} D_{\text{last},4}^2 \tau_4^4}{4\mu^2 C_{\text{step},6}^{3/2}} \right), \end{aligned}$$

where in the last inequality we use that $\gamma \leq 1/(L C_{\text{step},6})$. Combining the above bounds in (54), we arrive at the expansion formula (50). \square

Lemma 13. Assume A1 and A2. Then for any $\gamma \in (0, 1/(L C_{\text{step},6})]$ it holds

(a) All eigenvalues $\tilde{\lambda}_i$, $i \in \{1, \dots, d^2\}$ of the matrix $\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*$ satisfy

$$2\mu(1 - \gamma L/2) \leq \tilde{\lambda}_i \leq 2L(1 - \gamma\mu/2);$$

(b) $\|(\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*)^{-1}\| \leq 1/2$;

(c) In addition,

$$(\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*)^{-1} = (\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^*)^{-1} + D \text{ where } \|D\| \leq \gamma L / \mu.$$

Proof. Assumption A1 guarantees that the symmetric matrix \mathbf{H}^* is positive-definite. Let $u_1, \dots, u_d \in \mathbb{R}^d$ and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq \mu > 0$ be its eigenvectors and eigenvalues, respectively. Then we notice that

$$\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^* = \mathbf{H}^* \otimes (\mathbf{I} - (\gamma/2) \mathbf{H}^*) + (\mathbf{I} - (\gamma/2) \mathbf{H}^*) \otimes \mathbf{H}^*.$$

Hence, the latter operator is also diagonalizable in the orthogonal basis $u_i \otimes u_j \in \mathbb{R}^{d^2}$ with the respective eigenvalues being equal to $\lambda_i(1 - (\gamma/2)\lambda_j) + \lambda_j(1 - (\gamma/2)\lambda_i)$. Hence, we obtain the first part of lemma (a). To prove (b) it remains to note that for $\gamma \leq 1/L$ it holds $(2\mu(1 - \gamma L/2))^{-1} \leq 1/2$. Set now

$$\begin{aligned} S &= \mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* \in \mathbb{R}^{d^2 \times d^2} \\ R &= \mathbf{H}^* \otimes \mathbf{H}^* \in \mathbb{R}^{d^2 \times d^2}. \end{aligned} \quad (61)$$

Then it is easy to observe that

$$(S - \gamma R)^{-1} = S^{-1} + S^{-1} \sum_{k=1}^{\infty} \gamma^k (RS^{-1})^k,$$

provided that $\gamma \|RS^{-1}\| < 1$. Since R and S are diagonalizable in the same orthogonal basis $\{u_i \otimes u_j\}_{1 \leq i, j \leq d}$ with the eigenvalues $\lambda_i \lambda_j$ and $\lambda_i + \lambda_j$, respectively, the condition $\gamma \|RS^{-1}\| < 1$ holds provided that $\gamma < 2/L$. Hence, for $\gamma \leq 1/L$, it holds that

$$(\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^* - \gamma \mathbf{H}^* \otimes \mathbf{H}^*)^{-1} = (\mathbf{H}^* \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^*)^{-1} + D,$$

where $D \in \mathbb{R}^{d^2 \times d^2}$ satisfies

$$\|D\| \leq 2\gamma \|S^{-1}\| \|RS^{-1}\| \leq \frac{\gamma L}{\mu}.$$

□

We now state an auxiliary lemma about the function $\mathcal{C}(\theta)$ from (49).

Lemma 14. Assume A1, A2, A3(2), and C1(2). Then, for any $\gamma \in (0, 1/(L C_{\text{step},2})]$, it holds

$$\left\| \int_{\mathbb{R}^d} \mathcal{C}(\theta) \pi_{\gamma}(d\theta) - \mathcal{C}(\theta^*) \right\|_F \leq C_2 \gamma^{1/2},$$

where the constant C_2 is given by

$$C_2 = \left(\frac{L^2 D_{\text{last},2}}{\mu^{3/2}} + \frac{L \sqrt{D_{\text{last},2}}}{\sqrt{\mu}} \right) \tau_2^2. \quad (62)$$

Proof. Recall that

$$\varepsilon_1(\theta) = \nabla F(\theta, \xi_1) - \nabla f(\theta).$$

Hence, using the definition of $\mathcal{C}(\theta)$ in (49), we get, with $\theta \in \mathbb{R}^d$, that

$$\begin{aligned} \mathcal{C}(\theta) - \mathcal{C}(\theta^*) &= \mathbb{E}[(\varepsilon_1(\theta) - \varepsilon_1(\theta^*))(\varepsilon_1(\theta) - \varepsilon_1(\theta^*))^T] + \mathbb{E}[\varepsilon_1(\theta^*)(\varepsilon_1(\theta) - \varepsilon_1(\theta^*))^T] \\ &\quad + \mathbb{E}[(\varepsilon_1(\theta) - \varepsilon_1(\theta^*))\varepsilon_1(\theta^*)^T]. \end{aligned}$$

Using A3(2), we obtain

$$\mathbb{E}[\|\varepsilon_1(\theta) - \varepsilon_1(\theta^*)\|^2] \leq L \langle \nabla f(\theta) - \nabla f(\theta^*), \theta - \theta^* \rangle - \|\nabla f(\theta) - \nabla f(\theta^*)\|^2 \leq L^2 \|\theta - \theta^*\|^2.$$

Hence, combining the previous inequalities and using Hölder's inequality, we obtain for any $\theta \in \mathbb{R}^d$, that

$$\|\mathcal{C}(\theta) - \mathcal{C}(\theta^*)\|_F \leq L^2 \|\theta - \theta^*\|^2 + \tau_2 L \|\theta - \theta^*\|.$$

Applying now C1(2), we obtain

$$\left\| \int_{\mathbb{R}^d} \mathcal{C}(\theta) \pi_{\gamma}(d\theta) - \mathcal{C}(\theta^*) \right\|_F \leq \int_{\mathbb{R}^d} \|\mathcal{C}(\theta) - \mathcal{C}(\theta^*)\|_F \pi_{\gamma}(d\theta) \leq L^2 \frac{D_{\text{last},2} \gamma \tau_2^2}{\mu} + \tau_2 L \sqrt{\frac{D_{\text{last},2} \gamma \tau_2^2}{\mu}}.$$

We conclude the proof by using the fact that $\gamma \mu \leq 1$. □

Now we prove (24). We use synchronous coupling construction defined by the pair of recursions:

$$\begin{aligned} \theta_{k+1} &= \theta_k - \gamma \nabla F(\theta_k, \xi_{k+1}), \quad \theta_0 \sim \nu \\ \tilde{\theta}_{k+1} &= \tilde{\theta}_k - \gamma \nabla F(\tilde{\theta}_k, \xi_{k+1}), \quad \tilde{\theta}_0 \sim \pi_{\gamma}. \end{aligned}$$

Recall that the corresponding coupling kernel is denoted as $K_{\gamma}(\cdot, \cdot)$. Then we obtain

$$\begin{aligned} \mathbb{E}_{\nu}[\bar{\theta}_n] - \theta^* &= n^{-1} \sum_{k=n+1}^{2n} \mathbb{E}_{\nu, \pi_{\gamma}}^{K_{\gamma}}[\theta_k - \tilde{\theta}_k] + n^{-1} \sum_{k=n+1}^{2n} \mathbb{E}_{\pi_{\gamma}}[\tilde{\theta}_k - \theta^*] \\ &= n^{-1} \sum_{k=n+1}^{2n} \mathbb{E}_{\nu, \pi_{\gamma}}^{K_{\gamma}}[\theta_k - \tilde{\theta}_k] + (\bar{\theta}_{\gamma} - \theta^*). \end{aligned}$$

Using (47) and C1(2), we obtain

$$\begin{aligned}\|\mathbb{E}_{\nu, \pi_\gamma}^K[\theta_k - \tilde{\theta}_k]\| &\leq (1 - \gamma\mu)^{k/2} \{\mathbb{E}_{\nu, \pi_\gamma}^K \|\theta_0 - \tilde{\theta}_0\|^2\}^{1/2} \\ &\leq (1 - \gamma\mu)^{k/2} (\mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^2] + \frac{\sqrt{2\gamma\tau_2}}{\sqrt{\mu}}).\end{aligned}$$

Summing the above bounds for k from $n + 1$ to $2n$, we obtain (24).

A.3 PROOF OF PROPOSITION 5

Note that

$$\mathbb{E}_\nu[\bar{\theta}_n^{(RR)} - \theta^*] = 2\mathbb{E}_\nu[\bar{\theta}_n^\gamma - \theta^*] - \mathbb{E}_\nu[\bar{\theta}_n^{2\gamma} - \theta^*].$$

Applying (24), we obtain

$$\|\mathbb{E}_\nu[\bar{\theta}_n^{(RR)} - \theta^*]\| \leq \left(\frac{L}{\mu} C_1 + \frac{LD_{\text{last},6}^{3/2} \tau_6^3}{2\mu^{5/2}}\right) \gamma^{3/2} + \mathcal{R}_3(\theta_0 - \theta^*, \gamma, n), \quad (63)$$

where

$$\|\mathcal{R}_3(\theta_0 - \theta^*, \gamma, n)\| \lesssim \frac{(1 - \gamma\mu)^{(n+1)/2}}{n\gamma\mu} (\mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^2] + \frac{\sqrt{\gamma\tau_2}}{\sqrt{\mu}}), \quad (64)$$

and the statement follows.

B PROOF OF THEOREM 3

Theorem 15 (Version of Theorem 3 with explicit constants). *Assume A1, A2, A3(6), and C1(6). Then for any $\gamma \in (0, 1/(L C_{\text{step},6})]$, $n \in \mathbb{N}$, and initial distribution ν on \mathbb{R}^d , the sequence of Polyak-Ruppert estimates (14) satisfies*

$$\mathbb{E}_\nu^{1/2} [\|\mathbf{H}^*(\bar{\theta}_n^{(\gamma)} - \theta^*)\|^2] \leq \frac{\sqrt{\text{Tr} \Sigma_\varepsilon^*}}{\sqrt{n}} + \frac{C_2}{\gamma^{1/2} n} + C_3 \gamma + \frac{C_4 \gamma^{1/2}}{n^{1/2}} + \mathcal{R}_2(n, \gamma, \|\theta_0 - \theta^*\|),$$

where we have set

$$C_2 = c_0 D_{\text{last},2}^{1/2} \tau_2, \quad C_3 = c_0 \frac{L D_{\text{last},4} \tau_4^2}{2\mu}, \quad C_4 = c_0 L D_{\text{last},2}^{1/2} \tau_2. \quad (65)$$

and the remainder term $\mathcal{R}_2(n, \gamma, \|\theta_0 - \theta^*\|)$ is given by

$$\begin{aligned}\mathcal{R}_2(n, \gamma, \|\theta_0 - \theta^*\|) &= \frac{c_0 L (1 - \gamma\mu)^{(n+1)/2}}{\gamma\mu n} \mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^2] \\ &\quad + \frac{L c_0 (1 - \gamma\mu)^{n+1}}{2n\gamma\mu} \mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^4].\end{aligned} \quad (66)$$

Proof. Throughout the proof we omit upper index (γ) both for the elements of the sequence $\{\theta_k^{(\gamma)}\}_{k \in \mathbb{N}}$ and Polyak-Ruppert averaged estimates $\bar{\theta}_n^{(\gamma)}$. Instead, we write simply θ_k and $\bar{\theta}_n$, respectively. Summing the recurrence (30), we obtain that

$$\mathbf{H}^*(\bar{\theta}_n - \theta^*) = \frac{\theta_{n+1} - \theta^*}{\gamma n} - \frac{\theta_{2n} - \theta^*}{\gamma n} - \frac{1}{n} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta_k) - \frac{1}{n} \sum_{k=n+1}^{2n} \eta(\theta_k). \quad (67)$$

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

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

where $\nabla^3 f(\cdot) \in \mathbb{R}^{d \times d \times d}$. Using A2, we thus obtain that

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

Applying Minkowski's inequality to the decomposition (32) and to the last term therein, we get

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

Applying C1(2), we obtain that for $\gamma \in (0; 1/(L C_{\text{step},2})]$ it holds that

$$\mathbb{E}_\nu \|\theta_k - \theta^*\|^2 \lesssim (1 - \gamma\mu)^k \mathbb{E}_\nu [\|\theta_0 - \theta^*\|^2] + \frac{D_{\text{last},2} \gamma \tau_2^2}{\mu}. \quad (68)$$

Moreover, from $\gamma \in (0; 1/(L C_{\text{step},4})]$ it holds that

$$\mathbb{E}_\nu^{1/2} \|\theta_k - \theta^*\|^4 \lesssim (1 - \gamma\mu)^k \mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^4] + \frac{D_{\text{last},4} \gamma \tau_4^2}{\mu}. \quad (69)$$

Combining Lemma 16 with previous inequalities, we obtain

$$\begin{aligned} \mathbb{E}_\nu^{1/2}[\|\mathbf{H}^*(\bar{\theta}_n - \theta^*)\|^2] &\lesssim \frac{\sqrt{\text{Tr } \Sigma_\varepsilon^*}}{\sqrt{n}} + \frac{D_{\text{last},2}^{1/2} \tau_2}{\gamma^{1/2} n} + \frac{L D_{\text{last},4} \gamma \tau_4^2}{2\mu} + \frac{L D_{\text{last},2}^{1/2} \gamma^{1/2} \tau_2}{\mu^{1/2} n^{1/2}} \\ &\quad + \frac{(1 - \gamma\mu)^{(n+1)/2}}{\gamma n} \left(\frac{L}{\mu} + 1 \right) \mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^2] + \frac{L(1 - \gamma\mu)^{n+1}}{n\gamma\mu} \mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^4], \end{aligned}$$

and the result follows. \square

Below we provide an auxiliary lemma used in the proof of Theorem 3.

Lemma 16. Assume A1, A2, A3(2), and C1(2). Then for any $\gamma \in (0; 1/(L C_{\text{step},2})]$ and any $n \in \mathbb{N}$, it holds

$$\mathbb{E}_\nu^{1/2}[\|\sum_{k=n+1}^{2n} \{\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\}\|^2] \lesssim \frac{L D_{\text{last},2}^{1/2} \sqrt{\gamma n} \tau_2}{\mu^{1/2}} + \frac{L(1 - \gamma\mu)^{(n+1)/2}}{\gamma\mu} \mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^2]. \quad (70)$$

Moreover, let $p \geq 2$, and assume A1, A2, A3(p), and C1(p). Then for any $\gamma \in (0; 1/(L C_{\text{step},p})]$ and $n \in \mathbb{N}$ it holds that

$$\begin{aligned} \mathbb{E}_\nu^{1/p}[\|\sum_{k=n+1}^{2n} \{\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\}\|^p] &\lesssim \frac{L D_{\text{last},p}^{1/2} \sqrt{\gamma n p} \tau_p}{\mu^{1/2}} \\ &\quad + \frac{L p (1 - \gamma\mu)^{(n+1)/2}}{\mu^{1/2} \gamma^{1/2}} \mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^p]. \end{aligned} \quad (71)$$

Proof. Since $\{\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\}$ is a martingale-difference sequence with respect to \mathcal{F}_k , we have

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

where $\varepsilon_{k+1}(\theta^*) = \nabla F(\theta^*, \xi_{k+1})$ uses the same noise variable ξ_{k+1} as $F(\theta_k, \xi_{k+1})$. Note that

$$\begin{aligned} \mathbb{E}_\nu [\|\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\|^2] &= \mathbb{E}_\nu [\|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\theta^*, \xi_{k+1})\|^2] \\ &\quad - 2\mathbb{E}_\nu [\langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\theta^*, \xi_{k+1}), \nabla f(\theta_k) - \nabla f(\theta^*) \rangle] + \|\nabla f(\theta_k) - \nabla f(\theta^*)\|^2. \end{aligned}$$

Using A2, A3(2), and taking conditional expectation with respect to \mathcal{F}_k , we obtain

$$\begin{aligned} \mathbb{E}_\nu [\|\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\|^2] &\leq \mathbb{E}_\nu [L \langle \nabla f(\theta_k) - \nabla f(\theta^*), \theta_k - \theta^* \rangle - \|\nabla f(\theta_k) - \nabla f(\theta^*)\|^2] \\ &\leq L^2 \mathbb{E}_\nu [\|\theta_k - \theta^*\|^2]. \end{aligned}$$

Thus, we obtain that

$$\mathbb{E}_\nu[\|\sum_{k=n+1}^{2n} \{\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\}\|^2] \leq L^2 \sum_{k=n+1}^{2n} \mathbb{E}_\nu[\|\theta_k - \theta^*\|^2],$$

and the statement (70) follows from the assumption C1(2). In order to prove (71), we apply Burkholder's inequality Osekowski (2012, Theorem 8.6) and obtain

$$\begin{aligned} \mathbb{E}_\nu^{1/p}[\|\sum_{k=n+1}^{2n} \{\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\}\|^p] &\leq p \mathbb{E}_\nu^{1/p}[(\sum_{k=n+1}^{2n} \|\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\|^2)^{p/2}] \\ &\leq p \left(\sum_{k=n+1}^{2n} \mathbb{E}_\nu^{2/p}[\|\varepsilon_{k+1}(\theta_k) - \varepsilon_{k+1}(\theta^*)\|^p] \right)^{1/2} \\ &\lesssim p L \left(\sum_{k=n+1}^{2n} \mathbb{E}_\nu^{2/p}[\|\theta_k - \theta^*\|^p] \right)^{1/2} \\ &\stackrel{(a)}{\lesssim} \frac{L D_{\text{last},p}^{1/2} \sqrt{\gamma n p \tau_p}}{\mu^{1/2}} + \frac{L p (1 - \gamma \mu)^{(n+1)/2}}{\mu^{1/2} \gamma^{1/2}} \mathbb{E}_\nu^{1/p}[\|\theta_0 - \theta^*\|^p], \end{aligned}$$

where in (a) we have additionally used C1(p). \square

C PROOF OF THEOREM 6

Within this section we often use the definition of the function $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ from (37):

$$\psi(\theta) = (1/2) \nabla^3 f(\theta^*)(\theta - \theta^*)^{\otimes 2} \quad (72)$$

Theorem 17 (Version of Theorem 6 with explicit constants). *Assume A1, A2, A3(6), and C1(6). Then for any $\gamma \in (0, 1/(L C_{\text{step},6})]$, initial distribution ν , and $n \in \mathbb{N}$, the Richardson-Romberg estimator $\bar{\theta}_n^{(RR)}$ defined in (36) satisfies*

$$\begin{aligned} \mathbb{E}_\nu^{1/2}[\|\mathbf{H}^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^2] &\leq \frac{\sqrt{\text{Tr } \Sigma_\varepsilon^*}}{n^{1/2}} + \frac{C_{\text{RR},1} \gamma^{1/2}}{n^{1/2}} + \frac{C_{\text{RR},2}}{\gamma^{1/2} n} + C_{\text{RR},3} \gamma^{3/2} + \frac{C_{\text{RR},4} \gamma}{n^{1/2}} \\ &\quad + \mathcal{R}_4(n, \gamma, \|\theta_0 - \theta^*\|), \end{aligned}$$

where we have set

$$\begin{aligned} C_{\text{RR},1} &= \frac{c_0 D_{\text{last},4} L \tau_4^2}{\mu^{3/2}} + \frac{c_0 L D_{\text{last},2}^{1/2} \tau_2}{\mu^{1/2}}, \quad C_{\text{RR},2} = \frac{c_0 D_{\text{last},2}^{1/2} \tau_2}{\mu^{1/2}} \\ C_{\text{RR},3} &= c_0 \left(\frac{L D_{\text{last},6}^{3/2} \tau_6^3}{\mu^{3/2}} + C_1 \right), \quad C_{\text{RR},4} = \frac{c_0 D_{\text{last},4} L \tau_4^2}{\mu}, \end{aligned} \quad (73)$$

C_1 is defined in (51), and the remainder term $\mathcal{R}_4(n, \gamma, \|\theta_0 - \theta^*\|)$ is given by

$$\begin{aligned} \mathcal{R}_4(n, \gamma, \|\theta_0 - \theta^*\|) &= \frac{c_0 L (1 - \gamma \mu)^{(n+1)/2}}{n \gamma \mu} \\ &\quad \times \left(\mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^6] + \mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^4] + \mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^2] + \frac{D_{\text{last},4} \gamma \tau_4^2}{\mu} \right). \end{aligned} \quad (74)$$

Proof. Using the recursion (30), we obtain that

$$\begin{aligned} \mathbf{H}^*(\bar{\theta}_n^{(RR)} - \theta^*) &= \frac{2(\theta_{n+1}^{(\gamma)} - \theta^*)}{\gamma n} - \frac{2(\theta_{2n}^{(\gamma)} - \theta^*)}{\gamma n} - \frac{\theta_{n+1}^{(2\gamma)} - \theta^*}{2\gamma n} + \frac{\theta_{2n}^{(2\gamma)} - \theta^*}{2\gamma n} \\ &\quad - \frac{1}{n} \sum_{k=n+1}^{2n} [2\varepsilon_{k+1}(\theta_k^{(\gamma)}) - \varepsilon_{k+1}(\theta_k^{(2\gamma)})] - \frac{1}{n} \sum_{k=n+1}^{2n} [2\eta(\theta_k^{(\gamma)}) - \eta(\theta_k^{(2\gamma)})]. \end{aligned} \quad (75)$$

Therefore, applying Minkowski's inequality to the decomposition (75), we obtain for any initial distribution ν that

$$\begin{aligned}
\mathbb{E}_\nu^{1/2}[\|\mathbf{H}^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^2] &\leq \underbrace{\frac{1}{n} \mathbb{E}_\nu^{1/2}[\|\sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*)\|^2]}_{T_1} + \underbrace{\frac{2}{\gamma n} \mathbb{E}_\nu^{1/2}[\|\theta_{n+1}^{(\gamma)} - \theta^*\|^2] + \frac{2}{\gamma n} \mathbb{E}_\nu^{1/2}[\|\theta_{2n}^{(\gamma)} - \theta^*\|^2]}_{T_2} \\
&\quad + \underbrace{\frac{1}{2\gamma n} \mathbb{E}_\nu^{1/2}[\|\theta_{n+1}^{(2\gamma)} - \theta^*\|^2] + \frac{1}{2\gamma n} \mathbb{E}_\nu^{1/2}[\|\theta_{2n}^{(2\gamma)} - \theta^*\|^2]}_{T_3} \\
&\quad + \underbrace{\frac{2}{n} \mathbb{E}_\nu^{1/2}[\|\sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta_k^{(\gamma)}) - \varepsilon_{k+1}(\theta^*)\|^2]}_{T_4} \\
&\quad + \underbrace{\frac{1}{n} \mathbb{E}_\nu^{1/2}[\|\sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta_k^{(2\gamma)}) - \varepsilon_{k+1}(\theta^*)\|^2]}_{T_5} + \underbrace{\|2\pi_\gamma(\psi) - \pi_{2\gamma}(\psi)\|}_{T_6} \\
&\quad + \underbrace{\frac{2}{n} \mathbb{E}_\nu^{1/2}[\|\sum_{k=n+1}^{2n} \eta(\theta_k^{(\gamma)}) - \pi_\gamma(\psi)\|^2] + \frac{1}{n} \mathbb{E}_\nu^{1/2}[\|\sum_{k=n+1}^{2n} \eta(\theta_k^{(2\gamma)}) - \pi_{2\gamma}(\psi)\|^2]}_{T_7}.
\end{aligned}$$

Now we upper bound the terms in the right-hand side of the above bound separately. First, we note that

$$T_1 = \frac{\sqrt{\text{Tr } \Sigma_\varepsilon^*}}{\sqrt{n}}.$$

Using C1(2), we get

$$T_2 + T_3 \lesssim \frac{(1 - \gamma\mu)^{n+1/2}}{\gamma n} \mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^2] + \frac{D_{\text{last},2}^{1/2} \tau_2}{\mu^{1/2} \gamma^{1/2} n}.$$

Applying Lemma 16, we get

$$T_4 + T_5 \lesssim \frac{L D_{\text{last},2}^{1/2} \gamma^{1/2} \tau_2}{\mu^{1/2} n^{1/2}} + \frac{L(1 - \gamma\mu)^{(n+1)/2}}{\mu \gamma n} \mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^2].$$

Now we proceed with the term T_6 . Applying the recurrence (11), we obtain that

$$\theta_1^{(\gamma)} - \theta^* = (\mathbf{I} - \gamma \mathbf{H}^*)(\theta_0^{(\gamma)} - \theta^*) - \gamma \varepsilon_1(\theta_0^{(\gamma)}) - \gamma \eta(\theta_0^{(\gamma)}). \quad (76)$$

Thus, taking expectation w.r.t. π_γ in both sides above, we get

$$\mathbf{H}^*(\bar{\theta}_\gamma - \theta^*) = \mathbb{E}_{\pi_\gamma}[\eta(\theta_0^{(\gamma)})] = \pi_\gamma(\psi) + \pi_\gamma(G),$$

where $G(\theta)$ is defined in (37) and writes as

$$G(\theta) = \frac{1}{2} \left(\int_0^1 t^2 \nabla^4 f(t\theta^* + (1-t)\theta) dt \right) (\theta - \theta^*)^{\otimes 3}.$$

Hence, applying A2 together with Proposition 2, we obtain that

$$T_6 = \|2\pi_\gamma(\psi) - \pi_{2\gamma}(\psi)\| \lesssim C_1 \gamma^{3/2}. \quad (77)$$

Finally, using Lemma 21, Lemma 20, and Lemma 18, we obtain that

$$\begin{aligned}
T_7 &\lesssim \frac{D_{\text{last},4} L \gamma \tau_4^2}{\mu n^{1/2}} + \frac{D_{\text{last},4} L \gamma^{1/2} \tau_4^2}{\mu^{3/2} n^{1/2}} + \frac{L D_{\text{last},6}^{3/2} \gamma^{3/2} \tau_6^3}{\mu^{3/2}} \\
&\quad + \frac{L(1 - \gamma\mu)^{(n+1)/2}}{n \gamma \mu} \left(\mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^6] + \mathbb{E}_\nu^{1/2}[\|\theta_0 - \theta^*\|^4] + \frac{D_{\text{last},4} \gamma \tau_4^2}{\mu} \right).
\end{aligned}$$

Combining the bounds above completes the proof. \square

Below we provide some auxiliary technical lemmas.

Lemma 18. Assume [A1](#), [A2](#), [A3](#)(4), and [C1](#)(4). Then for any $\gamma \in (0; 1/(\mathbf{L} \mathbf{C}_{\text{step},4})]$ and any $n \in \mathbb{N}$ it holds

$$n^{-1} \mathbb{E}_{\pi_\gamma}^{1/2} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \right\|^2 \right] \lesssim \frac{\mathbf{D}_{\text{last},4} \mathbf{L}_3 \gamma \tau_4^2}{\mu n^{1/2}} + \frac{\mathbf{D}_{\text{last},4} \mathbf{L}_3 \gamma^{1/2} \tau_4^2}{\mu^{3/2} n^{1/2}}. \quad (78)$$

Proof. Using the fact that π_γ is a stationary distribution, we obtain that

$$\begin{aligned} \mathbb{E}_{\pi_\gamma} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \right\|^2 \right] &= n \mathbb{E}_{\pi_\gamma} [\|\psi(\theta_0) - \pi_\gamma(\psi)\|^2] \\ &\quad + \sum_{k=1}^{n-1} (n-k) \mathbb{E}_{\pi_\gamma} [(\psi(\theta_0) - \pi_\gamma(\psi))^T (\psi(\theta_k) - \pi_\gamma(\psi))] \end{aligned}$$

Using the Markov property, Cauchy–Schwartz inequality, [Proposition 1](#), and [Lemma 22](#), we obtain

$$\mathbb{E}_{\pi_\gamma} [(\psi(\theta_0) - \pi_\gamma(\psi))^T (\psi(\theta_k) - \pi_\gamma(\psi))] \quad (79)$$

$$= \mathbb{E}_{\pi_\gamma} [(\psi(\theta_0) - \pi_\gamma(\psi))^T (\mathbf{Q}_\gamma^k \psi(\theta_0) - \pi_\gamma(\psi))] \quad (80)$$

$$\stackrel{(a)}{\lesssim} (1/2)^{k/m(\gamma)} \mathbf{L}_3 \mathbb{E}_{\pi_\gamma} [\|\psi(\theta_0) - \pi_\gamma(\psi)\| \int c(\theta_0, \vartheta) d\pi_\gamma(\vartheta)], \quad (81)$$

where in (a) we additionally used the fact that

$$\mathbf{W}_c(\delta_{\theta_0}, \pi_\gamma) = \int c(\theta_0, \vartheta) d\pi_\gamma(\vartheta).$$

Using [C1](#)(4), we get

$$\mathbb{E}_{\pi_\gamma} [\|\psi(\theta_0) - \pi_\gamma\|^2] \leq \mathbb{E}_{\pi_\gamma} [\|\psi(\theta_0)\|^2] \leq \mathbf{L}_3^2 \mathbb{E}_{\pi_\gamma} [\|\theta_0 - \theta^*\|^4] \leq \frac{\mathbf{L}_3^2 \mathbf{D}_{\text{last},4} \gamma^2 \tau_4^4}{\mu^2}, \quad (82)$$

and, using [C1](#)(2) and [C1](#)(4), we get

$$\int \int c^2(\theta_0, \vartheta) d\pi_\gamma(\vartheta) d\pi_\gamma(\theta_0) \quad (83)$$

$$\leq \int \int \|\theta_0 - \vartheta\|^2 \left(\|\theta_0 - \theta^*\| + \|\vartheta - \theta^*\| + \frac{2^{3/2} \gamma^{1/2} \tau_2}{\mu^{1/2}} \right)^2 d\pi_\gamma(\vartheta) d\pi_\gamma(\theta_0) \quad (84)$$

$$\lesssim \int \int (\|\theta_0 - \theta^*\|^4 + \|\vartheta - \theta^*\|^4) + \frac{\gamma \tau_2^2}{\mu} (\|\theta_0 - \theta^*\|^2 + \|\vartheta - \theta^*\|^2) d\pi_\gamma(\vartheta) d\pi_\gamma(\theta_0) \quad (85)$$

$$\lesssim \frac{\mathbf{D}_{\text{last},4} \gamma^2 \tau_4^2}{\mu^2} + \frac{\mathbf{D}_{\text{last},2} \gamma^2 \tau_2^4}{\mu^2} \lesssim \frac{\mathbf{D}_{\text{last},4} \gamma^2 \tau_4^4}{\mu^2}. \quad (86)$$

Using (82), (83), and Cauchy–Schwartz inequality for (79), we obtain

$$\mathbb{E}_{\pi_\gamma} [(\psi(\theta_0) - \pi_\gamma(\psi))^T (\psi(\theta_k) - \pi_\gamma(\psi))] \lesssim (1/2)^{k/m(\gamma)} \frac{\mathbf{L}_3 \mathbf{D}_{\text{last},4} \gamma^2 \tau_4^4}{\mu^2}.$$

Combining the inequalities above and using that $m(\gamma) = \lceil 2 \frac{\log 4}{\gamma \mu} \rceil \leq \frac{2 \log 4 + 1}{\gamma \mu}$, we get

$$\begin{aligned} n^{-1} \mathbb{E}_{\pi_\gamma}^{1/2} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \right\|^2 \right] &\leq \left(\frac{\mathbf{D}_{\text{last},4} \mathbf{L}_3^2 \gamma^2 \tau_4^4}{\mu^2 n} + \frac{\mathbf{D}_{\text{last},4} m(\gamma) \mathbf{L}_3^2 \gamma^2 \tau_4^4}{\mu^2 n} \right)^{1/2} \\ &\lesssim \frac{\mathbf{D}_{\text{last},4} \mathbf{L}_3 \gamma \tau_4^2}{\mu n^{1/2}} + \frac{\mathbf{D}_{\text{last},4} \mathbf{L}_3 \gamma^{1/2} \tau_4^2}{\mu^{3/2} n^{1/2}}. \end{aligned}$$

□

Lemma 19. Assume A1, A2, A3(4). Then for any $\gamma \in (0, \frac{2}{11L}]$, and any $k \in \mathbb{N}$ it holds that

$$\mathbb{E}[\|\theta_{k+1} - \tilde{\theta}_{k+1}\|^4 | \mathcal{F}_k] \leq (1 - \gamma\mu)^2 \|\theta_k - \tilde{\theta}_k\|^4. \quad (87)$$

Moreover, let $p \geq 2$ and assume A1, A2, and A3(2p). Then for any $\gamma \in (0, \frac{p}{4 \cdot 3^p L}]$ and any $k \in \mathbb{N}$ it holds that

$$\mathbb{E}[\|\theta_{k+1} - \tilde{\theta}_{k+1}\|^{2p} | \mathcal{F}_k] \leq (1 - \gamma\mu)^p \|\theta_k - \tilde{\theta}_k\|^{2p}. \quad (88)$$

Proof. Recall that the sequences $\{\theta_k\}_{k \in \mathbb{N}}$ and $\{\tilde{\theta}_k\}_{k \in \mathbb{N}}$ are defined by the recurrences

$$\theta_{k+1} = \theta_k - \gamma \nabla F(\theta_k, \xi_{k+1}), \quad \theta_0 = \theta \in \mathbb{R}^d, \quad (89)$$

$$\tilde{\theta}_{k+1} = \tilde{\theta}_k - \gamma \nabla F(\tilde{\theta}_k, \xi_{k+1}), \quad \tilde{\theta}_0 = \tilde{\theta} \in \mathbb{R}^d. \quad (90)$$

Expanding the brackets, we obtain that

$$\begin{aligned} \|\theta_{k+1} - \tilde{\theta}_{k+1}\|^4 &= \|\theta_k - \tilde{\theta}_k\|^4 + \gamma^4 \|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^4 \\ &\quad + 4\gamma^2 \langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle^2 \\ &\quad + 2\gamma^2 \|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^2 \|\theta_k - \tilde{\theta}_k\|^2 \\ &\quad - 4\gamma \langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle \|\theta_k - \tilde{\theta}_k\|^2 \\ &\quad - 4\gamma^3 \langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle \|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^2 \end{aligned}$$

Using A3(4) and Cauchy-Schwartz inequality, we get

$$\begin{aligned} \mathbb{E}[\|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^4 | \mathcal{F}_k] &\leq L^3 \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \|\theta_k - \tilde{\theta}_k\|^2, \\ \mathbb{E}[\langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle^2 | \mathcal{F}_k] &\leq L \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \|\theta_k - \tilde{\theta}_k\|^2, \\ \mathbb{E}[\|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^2 \|\theta_k - \tilde{\theta}_k\|^2 | \mathcal{F}_k] &\leq L \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \|\theta_k - \tilde{\theta}_k\|^2, \\ \mathbb{E}[\langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle \|\theta_k - \tilde{\theta}_k\|^2 | \mathcal{F}_k] &= \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \|\theta_k - \tilde{\theta}_k\|^2. \end{aligned}$$

Similarly,

$$\begin{aligned} \mathbb{E}[\langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle \|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^2 | \mathcal{F}_k] \\ \leq L^2 \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \|\theta_k - \tilde{\theta}_k\|^2 \end{aligned}$$

Combining all inequalities above, we obtain

$$\begin{aligned} \mathbb{E}[\|\theta_{k+1} - \tilde{\theta}_{k+1}\|^4 | \mathcal{F}_k] &\leq \|\theta_k - \tilde{\theta}_k\|^4 \\ &\quad - (4\gamma - \gamma^4 L^3 - 4\gamma^2 L - 2\gamma^2 L - 4\gamma^3 L^2) \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \|\theta_k - \tilde{\theta}_k\|^2 \end{aligned}$$

Using A1 and since $1 - \gamma^3 L^3 / 4 - 3\gamma L / 2 - \gamma^2 L^2 \geq 1 - 11\gamma L / 4$, we get

$$\begin{aligned} \mathbb{E}[\|\theta_{k+1} - \tilde{\theta}_{k+1}\|^4 | \mathcal{F}_k] &\leq (1 - 4\gamma\mu(1 - 11\gamma L / 4)) \|\theta_k - \tilde{\theta}_k\|^4 \\ &\leq (1 - 2\gamma\mu(1 - 11\gamma L / 4))^2 \|\theta_k - \tilde{\theta}_k\|^4. \end{aligned}$$

Since $1 - 11\gamma L / 4 \geq 1/2$ for $\gamma \leq 2/(11L)$, we complete the proof.

For simplicity of proof for second part of lemma we define $\delta_{k+1} = \|\theta_{k+1} - \theta'_{k+1}\|$. Then we have

$$\begin{aligned} \mathbb{E}[\delta_{k+1}^{2p} | \mathcal{F}_k] &= \mathbb{E}[(\delta_k^2 - 2\gamma \langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle + \gamma^2 \|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^2)^p | \mathcal{F}_k] \\ &= \mathbb{E}[\sum_{\substack{i+j+l=p \\ i,j,l \in \{0, \dots, p\}}} \frac{p!}{i!j!l!} \delta_k^{2i} (-2\gamma \langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle)^j \gamma^{2l} \|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^{2l} | \mathcal{F}_k]. \end{aligned}$$

Now we bound each term in the sum above.

1. First, for $i = p, j = 0, l = 0$ the corresponding term in the sum is equal to δ_k^{2p} .

2. Second, for $i = p - 1, j = 1, l = 0$, we have

$$\mathbb{E}[(-2\gamma\langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle) | \mathcal{F}_k] = -2\gamma\langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle.$$

3. Third, for $l \geq 1$ or $j \geq 2$ we use Cauchy-Schwartz inequality and get

$$\begin{aligned} & (2\gamma\langle \nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1}), \theta_k - \tilde{\theta}_k \rangle)^j \gamma^{2l} \|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^{2l} \\ & \leq 2^j \gamma^{j+2l} \delta_k^j \|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^{2l+j}. \end{aligned}$$

Moreover using A3(2p), we get

$$\begin{aligned} & \mathbb{E}[2^j \gamma^{j+2l} \delta_k^{2i+j} \|\nabla F(\theta_k, \xi_{k+1}) - \nabla F(\tilde{\theta}_k, \xi_{k+1})\|^{2l+j} | \mathcal{F}_k] \\ & \leq 2^j \gamma^{j+2l} \delta_k^{2p-2} L^{2l+j-1} \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle. \end{aligned}$$

Combining all inequalities above, we obtain

$$\begin{aligned} \mathbb{E}[\delta_{k+1}^{2p} | \mathcal{F}_k] & \leq \delta_k^{2p} - 2p\gamma \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \delta_k^{2p-2} \\ & \quad + \left(\sum_{\substack{i+j+l=p \\ i,j,l \in \{0, \dots, p\} \\ j+2l \geq 2}} \frac{p!}{i!j!l!} 2^j \gamma^{j+2l} L^{2l+j-1} \right) \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \delta_k^{2p-2}. \end{aligned}$$

Since $\gamma \leq \frac{p}{3^p 4L}$, we have

$$\mathbb{E}[\delta_{k+1}^{2p} | \mathcal{F}_k] \leq \delta_k^{2p} - \gamma p \langle \nabla f(\theta_k) - \nabla f(\tilde{\theta}_k), \theta_k - \tilde{\theta}_k \rangle \delta_k^{2p-2}.$$

It remains to apply A1 together with an elementary bound $(1 - p\mu\gamma) \leq (1 - \gamma\mu)^p$. \square

Further, without loss of generality, we assume that $C_{\text{step}, 2p} = \min(C_{\text{step}, 2p}, \frac{p}{3^p 4L})$ and $C_{\text{step}, 4} = \min(C_{\text{step}, 4}, \frac{2}{11L})$.

Lemma 20. Assume A1, A2, A3(4), and C1(4). Then for any $\gamma \in (0; 1/(L C_{\text{step}, 4})]$, any $n \in \mathbb{N}$ and initial distribution ν it holds

$$\begin{aligned} n^{-1} \mathbb{E}_\nu^{1/2} [\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \|^2] & \lesssim n^{-1} \mathbb{E}_{\pi_\gamma}^{1/2} [\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \|^2] \\ & \quad + \frac{L_3(1 - \gamma\mu)^{(n+1)/2}}{n\gamma\mu} \left(\mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^4] + \frac{D_{\text{last}, 4} \gamma \tau_4^2}{\mu} \right). \end{aligned}$$

Proof. Using the synchronous coupling construction defined in (46) and the corresponding coupling kernel K_γ , we obtain that

$$\begin{aligned} \mathbb{E}_\nu^{1/2} [\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \|^2] & = (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} [\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \|^2])^{1/2} \\ & \leq \mathbb{E}_{\pi_\gamma}^{1/2} [\| \sum_{k=n+1}^{2n} \{\psi(\tilde{\theta}_k) - \pi_\gamma(\psi)\} \|^2] + (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} [\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \psi(\tilde{\theta}_k)\} \|^2])^{1/2} \end{aligned} \tag{91}$$

Applying Minkowski's inequality to the last term and using Lemma 22, we get

$$\begin{aligned} (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} [\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \psi(\tilde{\theta}_k)\} \|^2])^{1/2} & \leq \sum_{k=n+1}^{2n} (\mathbb{E}_{\nu, \pi_\gamma}^K [\| \{\psi(\theta_k) - \psi(\tilde{\theta}_k)\} \|^2])^{1/2} \\ & \leq \frac{L_3}{2} \sum_{k=n+1}^{2n} (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} [c^2(\theta_k, \tilde{\theta}_k)])^{1/2}. \end{aligned}$$

Using Hölder's and Minkowski's inequality and applying Lemma 19, (68) and (69), we obtain

$$\begin{aligned}
& (\mathbb{E}_{\nu, \pi_\gamma}^{\mathbf{K}_\gamma} [c^2(\theta_k, \tilde{\theta}_k)])^{1/2} \\
& \leq (\mathbb{E}_{\nu, \pi_\gamma}^{\mathbf{K}_\gamma} [\|\theta_k - \tilde{\theta}_k\|^4])^{1/4} (\mathbb{E}_{\pi_\gamma}^{1/4} [\|\tilde{\theta}_k - \theta^*\|^4] + \mathbb{E}_\nu^{1/4} [\|\theta_k - \theta^*\|^4 + \frac{\gamma^{1/2} \tau_2}{\mu^{1/2}}]) \\
& \leq (1 - \gamma\mu)^{k/2} (\mathbb{E}_{\nu, \pi_\gamma}^{\mathbf{K}_\gamma} [\|\theta_0 - \tilde{\theta}_0\|^4])^{1/4} (\mathbb{E}_\nu^{1/4} [\|\theta_0 - \theta^*\|^4] + \frac{D_{\text{last},4}^{1/2} \gamma^{1/2} \tau_4}{\mu^{1/2}} + \frac{\gamma^{1/2} \tau_2}{\mu^{1/2}}) \\
& \lesssim (1 - \gamma\mu)^{k/2} \left(\frac{D_{\text{last},4} \gamma \tau_4^2}{\mu} + \mathbb{E}_\nu^{1/2} \|\theta_0 - \theta^*\|^4 \right)
\end{aligned}$$

Combining all inequalities above, we get

$$(\mathbb{E}_{\nu, \pi_\gamma}^{\mathbf{K}_\gamma} [\|\sum_{k=n+1}^{2n} \{\psi(\theta_k) - \psi(\theta'_k)\}\|^2])^{1/2} \lesssim \frac{L_3(1 - \gamma\mu)^{(n+1)/2}}{\gamma\mu} \left(\mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^4] + \frac{D_{\text{last},4} \gamma \tau_4^2}{\mu} \right).$$

Substituting the last inequality into (91) we complete the proof. \square

Lemma 21. Assume A1, A2, A3(6), and C1(6). Then for any $\gamma \in (0; 1/(L C_{\text{step},6})]$, $n \in \mathbb{N}$, and initial distribution ν , it holds that

$$\begin{aligned}
n^{-1} \mathbb{E}_\nu^{1/2} \left[\sum_{k=n+1}^{2n} \|\eta(\theta_k) - \pi_\gamma(\psi)\|^2 \right] & \leq n^{-1} \mathbb{E}_\nu^{1/2} \left[\sum_{k=n+1}^{2n} \|\psi(\theta_k) - \pi_\gamma(\psi)\|^2 \right] \\
& + \frac{L_4(1 - \gamma\mu)^{(n+1)/2}}{n\gamma\mu} \mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^6] + \frac{L_4 D_{\text{last},6}^{3/2} \gamma^{3/2} \tau_6^3}{3\mu^{3/2}}.
\end{aligned} \tag{92}$$

Proof. Applying the 4-rd order Taylor expansion with integral remainder, we get that

$$\eta(\theta) = \psi(\theta) + \frac{1}{2} \left(\int_0^1 t^2 \nabla^4 f(t\theta^* + (1-t)\theta) dt \right) (\theta - \theta^*)^{\otimes 3}, \tag{93}$$

and using A2, we obtain

$$(1/2) \left\| \left(\int_0^1 t^2 \nabla^4 f(t\theta^* + (1-t)\theta) dt \right) (\theta - \theta^*)^{\otimes 3} \right\| \leq L_4 \|\theta - \theta^*\|^3. \tag{94}$$

Therefore, combining (93), A2, and applying Minkowski's inequality, we get

$$\begin{aligned}
\mathbb{E}_\nu^{1/2} \left[\sum_{k=n+1}^{2n} \|\eta(\theta_k) - \pi_\gamma(\psi)\|^2 \right] & \leq \mathbb{E}_\nu^{1/2} \left[\sum_{k=n+1}^{2n} \|\psi(\theta_k) - \pi_\gamma(\psi)\|^2 \right] \\
& + \frac{L_4}{6} \sum_{k=n+1}^{2n} \mathbb{E}_\nu^{1/2} [\|\theta_k - \theta^*\|^6]
\end{aligned} \tag{95}$$

Applying C1(6) for the last term of (95), we get

$$\begin{aligned}
\mathbb{E}_\nu^{1/2} \left[\sum_{k=n+1}^{2n} \|\eta(\theta_k) - \pi_\gamma(\psi)\|^2 \right] & \lesssim \mathbb{E}_\nu^{1/2} \left[\sum_{k=n+1}^{2n} \|\psi(\theta_k) - \pi_\gamma(\psi)\|^2 \right] + \frac{L_4 n D_{\text{last},6}^{3/2} \gamma^{3/2} \tau_6^3}{\mu^{3/2}} \\
& + \frac{L_4(1 - \gamma\mu)^{3(n+1)/2}}{1 - (1 - \gamma\mu)^{3/2}} \mathbb{E}_\nu^{1/2} [\|\theta_0 - \theta^*\|^6].
\end{aligned} \tag{96}$$

It remains to notice that $(1 - \gamma\mu)^{3/2} \leq (1 - \gamma\mu)$, and the statement follows. \square

We conclude this section with a technical statement on the properties of the function ψ from (72).

Lemma 22. Let $\psi(\cdot)$ be a function defined in (72). Then for any $\theta, \theta' \in \mathbb{R}^d$, it holds that

$$\|\psi(\theta) - \psi(\theta')\| \leq \frac{1}{2} L_3 c(\theta, \theta').$$

Proof. For simplicity, let us denote $T = \nabla^3 f(\theta^*)$. Hence,

$$\|\psi(\theta) - \psi(\theta')\| \leq \frac{1}{2} \|T(\theta - \theta^*)^{\otimes 2} - T(\theta' - \theta^*)^{\otimes 2}\|. \quad (97)$$

Note that

$$\|T\| = \sup_{x \neq 0, y \neq 0, z \neq 0} \frac{\sum_{i,j,k} T_{ijk} x_i y_j z_k}{\|x\| \|y\| \|z\|} \geq \sup_{x \neq 0, y \neq 0, z \neq 0} \sup_k \frac{z_k \sum_{i,j} T_{ijk} x_i y_j}{\|z\| \|y\| \|x\|} = \sup_{x \neq 0, y \neq 0} \frac{\|t(x, y)\|}{\|y\| \|x\|}, \quad (98)$$

where $t(x, y)_k = \sum_{i,j} T_{ijk} x_i y_j$. Therefore, for any $x, y \in \mathbb{R}^d$, it holds that

$$\|t(x, y)\| \leq \|x\| \|y\| \|T\| \quad (99)$$

We denote $v = T x^{\otimes 2} - T y^{\otimes 2}$. Then

$$\begin{aligned} v_k &= \sum_{i,j} T_{ijk} (x_i x_j - y_i y_j) = \sum_{i,j} T_{ijk} ((x_i - y_i) x_j + (x_i - y_i) y_j) = \\ &= \sum_{i,j} T_{ijk} (x_i - y_i) x_j + \sum_{i,j} T_{ijk} (x_i - y_i) y_j, \end{aligned} \quad (100)$$

where the first inequality is true since $T_{ijk} = T_{jik}$ by definition of T . Combining (99) and (C) and using triangle inequality, we obtain

$$\|v\| \leq \|T\| \|x - y\| (\|x\| + \|y\|) \leq \|T\| \|x - y\| (\|x\| + \|y\| + \frac{2\sqrt{2}\tau_2\sqrt{\gamma}}{\sqrt{\mu}}).$$

We complete the proof setting $x = \theta - \theta^*, y = \theta' - \theta^*$ □

D PROOF OF THEOREM 9

Theorem 23 (Version of Theorem 9 with explicit constants). Let $p \geq 2$ and assume A1, A2, A3(3p), and C1(3p). Then for any $\gamma \in (0; 1/(L C_{\text{step}, 3p})]$, initial distribution ν , and $n \in \mathbb{N}$, the estimator $\bar{\theta}_n^{(RR)}$ defined in (36) satisfies

$$\begin{aligned} \mathbb{E}_\nu^{1/p} [\|H^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^p] &\leq \frac{c_1 \sqrt{\text{Tr} \Sigma_\varepsilon^*} p^{1/2}}{n^{1/2}} + \frac{c_2 p \tau_p}{n^{1-1/p}} + \frac{C_{RR,5}}{n \gamma^{1/2}} + \frac{C_{RR,6} \gamma^{1/2}}{n^{1/2}} + C_{RR,7} \gamma^{3/2} \\ &\quad + \frac{C_{RR,8}}{n} + \mathcal{R}_5(n, \gamma, \|\theta_0 - \theta^*\|), \end{aligned}$$

where we have set

$$\begin{aligned} C_{RR,5} &= \frac{c_0 D_{\text{last}, p}^{1/2} \tau_p}{\mu^{1/2}}, \quad C_{RR,6} = \frac{c_0 L D_{\text{last}, p}^{1/2} p \tau_p}{\mu^{1/2}} + \frac{c_0 L D_{\text{last}, 2p} p \tau_{2p}^2}{\mu^{3/2}}, \\ C_{RR,7} &= c_0 \left(C_1 + \frac{L D_{\text{last}, 3p}^3 \tau_{3p}^3}{\mu^{3/2}} \right), \quad C_{RR,8} = \frac{c_0 L D_{\text{last}, 2p} \tau_{2p}}{\mu^2}, \end{aligned} \quad (101)$$

C_1 is defined in (51), and the remainder term $\mathcal{R}_5(n, \gamma, \|\theta_0 - \theta^*\|)$ is given by

$$\begin{aligned} \mathcal{R}_5(n, \gamma, \|\theta_0 - \theta^*\|) &= \frac{c_0 (1 - \gamma \mu)^{(n+1)/2}}{\gamma n} \mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^p] + \frac{c_0 L p (1 - \gamma \mu)^{(n+1)/2}}{\mu^{1/2} \gamma^{1/2} n} \mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^p] \\ &\quad + \frac{c_0 L (1 - \gamma \mu)^{(n+1)/2}}{\gamma \mu n} \left(\mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^{2p}] + \frac{D_{\text{last}, 2p} \gamma \tau_{2p}^2}{\mu} \right) + \frac{c_0 L (1 - \gamma \mu)^{(3/2)n}}{\gamma \mu} \mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^{3p}] \end{aligned} \quad (102)$$

Proof. Using the decomposition (75), we obtain that for any $p \geq 2$, it holds that

$$\begin{aligned}
\mathbb{E}_\nu^{1/p}[\|\mathbf{H}^*(\bar{\theta}_n^{(RR)} - \theta^*)\|^p] &\lesssim \underbrace{\frac{1}{n} \mathbb{E}_\nu^{1/p}[\|\sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*)\|^p]}_{T_1} + \underbrace{\frac{1}{\gamma n} \mathbb{E}_\nu^{1/p}[\|\theta_{n+1}^{(\gamma)} - \theta^*\|^p] + \frac{1}{\gamma n} \mathbb{E}_\nu^{1/p}[\|\theta_{2n}^{(\gamma)} - \theta^*\|^p]}_{T_2} \\
&+ \underbrace{\frac{1}{\gamma n} \mathbb{E}_\nu^{1/p}[\|\theta_{n+1}^{(2\gamma)} - \theta^*\|^p] + \frac{1}{\gamma n} \mathbb{E}_\nu^{1/p}[\|\theta_{2n}^{(2\gamma)} - \theta^*\|^p]}_{T_3} \\
&+ \underbrace{\frac{1}{n} \mathbb{E}_\nu^{1/p}[\|\sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta_k^{(\gamma)}) - \varepsilon_{k+1}(\theta^*)\|^p]}_{T_4} \\
&+ \underbrace{\frac{1}{n} \mathbb{E}_\nu^{1/p}[\|\sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta_k^{(2\gamma)}) - \varepsilon_{k+1}(\theta^*)\|^p]}_{T_5} + \underbrace{\|2\pi_\gamma(\psi) - \pi_{2\gamma}(\psi)\|}_{T_6} \\
&+ \underbrace{\frac{1}{n} \mathbb{E}_\nu^{1/p}[\|\sum_{k=n+1}^{2n} \psi(\theta_k^{(\gamma)}) - \pi_\gamma(\psi)\|^p] + \frac{1}{n} \mathbb{E}_\nu^{1/p}[\|\sum_{k=n+1}^{2n} \psi(\theta_k^{(2\gamma)}) - \pi_{2\gamma}(\psi)\|^p]}_{T_7} \\
&+ \underbrace{\frac{1}{n} \sum_{k=n+1}^{2n} \mathbb{E}_\nu^{1/p}[\|G(\theta_k^{(\gamma)})\|^p] + \frac{1}{n} \sum_{k=n+1}^{2n} \mathbb{E}_\nu^{1/p}[\|G(\theta_k^{(2\gamma)})\|^p]}_{T_8}.
\end{aligned}$$

Now we upper bounds the terms above separately. Applying first the Pinelis version of Rosenthal inequality (Pinelis, 1994) together with A3(p), we obtain that

$$T_1 \leq \frac{c_1 \sqrt{\text{Tr} \Sigma_\varepsilon^*} p^{1/2}}{n^{1/2}} + \frac{c_2 p \tau_p}{n^{1-1/p}}.$$

Applying C1(p) (which is implied by C1(3p)), we obtain that

$$T_2 + T_3 \lesssim \frac{D_{\text{last},p}^{1/2} \tau_p}{\mu^{1/2} n \gamma^{1/2}} + \frac{(1 - \gamma \mu)^{(n+1)/2}}{\gamma n} \mathbb{E}_\nu^{1/p}[\|\theta_0 - \theta^*\|^p].$$

Applying Lemma 16 (see the bound (71)), we get that

$$T_4 + T_5 \lesssim \frac{L D_{\text{last},p}^{1/2} \gamma^{1/2} p \tau_p}{\mu^{1/2} n^{1/2}} + \frac{L p (1 - \gamma \mu)^{(n+1)/2}}{\mu^{1/2} \gamma^{1/2} n} \mathbb{E}_\nu^{1/p}[\|\theta_0 - \theta^*\|^p].$$

Using the bounds (76) and (77), we obtain

$$T_6 \lesssim C_1 \gamma^{3/2}.$$

Applying Proposition 8, we get

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

Using this bound and Lemma 24, we obtain that

$$T_7 \lesssim \frac{L D_{\text{last},2p} p \tau_{2p}^2 \gamma^{1/2}}{\mu^{3/2} n^{1/2}} + \frac{L D_{\text{last},2p} \tau_{2p}}{\mu^2 n} + \frac{L (1 - \gamma \mu)^{(n+1)/2}}{\gamma \mu n} \left(\mathbb{E}_\nu^{1/p}[\|\theta_0 - \theta^*\|^{2p}] + \frac{D_{\text{last},2p} \gamma \tau_{2p}^2}{\mu} \right)$$

Finally, applying the definition of $G(\theta)$ in (37) together with C1(3p), we obtain that

$$\begin{aligned}
T_8 &\lesssim \frac{L D_{\text{last},3p}^{3/2} \gamma^{3/2} \tau_{3p}^3}{\mu^{3/2}} + \frac{L}{n} \sum_{k=n+1}^{2n} (1 - \gamma \mu)^{(3/2)k} \mathbb{E}_\nu^{1/p}[\|\theta_0 - \theta^*\|^{3p}] \\
&\lesssim \frac{L D_{\text{last},3p}^{3/2} \gamma^{3/2} \tau_{3p}^3}{\mu^{3/2}} + \frac{L (1 - \gamma \mu)^{(3/2)n}}{\gamma \mu} \mathbb{E}_\nu^{1/p}[\|\theta_0 - \theta^*\|^{3p}].
\end{aligned}$$

To complete the proof it remains to combine the bounds for T_1 to T_8 . \square

D.1 PROOF OF PROPOSITION 8

In the proof below we use the notation

$$\bar{\psi}(\theta) = \psi(\theta) - \pi_\gamma(\psi).$$

We proceed with the blocking technique. Indeed, let us set the parameter

$$m = m(\gamma) = \left\lceil \frac{2 \log 4}{\gamma \mu} \right\rceil. \quad (103)$$

Our choice of parameter $m(\gamma)$ is due to Proposition 1. For notation conciseness we write it simply as m , dropping its dependence upon γ . Using Minkowski's inequality, we obtain that

$$\mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=0}^{n-1} \bar{\psi}(\theta_k) \right\|^p \right] \leq \mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=0}^{\lfloor n/m \rfloor m - 1} \bar{\psi}(\theta_k) \right\|^p \right] + m \mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \bar{\psi}(\theta_0) \right\|^p \right]. \quad (104)$$

Now we consider the Poisson equation, associated with \mathbb{Q}_γ^m and function $\bar{\psi}$, that is,

$$g_m(\theta) - \mathbb{Q}_\gamma^m g_m(\theta) = \bar{\psi}(\theta). \quad (105)$$

The function

$$g_m(\theta) = \sum_{k=0}^{\infty} \mathbb{Q}_\gamma^{km} \bar{\psi}(\theta) \quad (106)$$

is well-defined under the assumptions **A1**, **A2**, **A3**($2p$), and **C1**($2p$). Moreover, g_m is a solution of the Poisson equation (105). Define $q := \lfloor n/m \rfloor$, then we have

$$\sum_{k=0}^{qm-1} \bar{\psi}(\theta_k) = \sum_{r=0}^{m-1} B_{m,r}, \quad \text{with} \quad B_{m,r} = \sum_{k=0}^{q-1} \{g_m(\theta_{km+r}) - \mathbb{Q}_\gamma^m g_m(\theta_{km+r})\}. \quad (107)$$

Using Minkowski's inequality, we get from (104), that

$$\mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=0}^{n-1} \bar{\psi}(\theta_k) \right\|^p \right] \leq m \mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=1}^q \{g_m(\theta_{km}) - \mathbb{Q}_\gamma^m g_m(\theta_{(k-1)m})\} \right\|^p \right] + 2m \mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \bar{\psi}(\theta_0) \right\|^p \right] \quad (108)$$

Now we upper bound both terms of (108) separately. Under assumption **A2**, and applying **C1**($2p$), we get

$$\mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \bar{\psi}(\theta_0) \right\|^p \right] \leq \frac{L}{2} \mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \theta_0 - \theta^* \right\|^{2p} \right] \leq \frac{L D_{\text{last}, 2p} \gamma \tau_{2p}^2}{2\mu}. \quad (109)$$

To proceed with the first term, we apply Burkholder's inequality (Osekowski, 2012, Theorem 8.6), and obtain that

$$\begin{aligned} \mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=1}^q \{g_m(\theta_{km}) - \mathbb{Q}_\gamma^m g_m(\theta_{(k-1)m})\} \right\|^p \right] \\ \leq p \mathbb{E}_{\pi_\gamma}^{1/p} \left[\left(\sum_{k=1}^q \left\| \{g_m(\theta_{km}) - \mathbb{Q}_\gamma^m g_m(\theta_{(k-1)m})\} \right\|^2 \right)^{p/2} \right]. \end{aligned} \quad (110)$$

Applying now Minkowski's inequality again, we get

$$\begin{aligned} \mathbb{E}_{\pi_\gamma}^{2/p} \left[\left(\sum_{k=1}^q \left\| \{g_m(\theta_{km}) - \mathbb{Q}_\gamma^m g_m(\theta_{(k-1)m})\} \right\|^2 \right)^{p/2} \right] &\leq q \mathbb{E}_{\pi_\gamma}^{2/p} \left[\left\| \{g_m(\theta_{km}) - \mathbb{Q}_\gamma^m g_m(\theta_{(k-1)m})\} \right\|^p \right] \\ &\lesssim q \left(\mathbb{E}_{\pi_\gamma}^{2/p} \left[\left\| g_m(\theta_0) \right\|^p \right] + \mathbb{E}_{\pi_\gamma}^{2/p} \left[\left\| \mathbb{Q}_\gamma^m g_m(\theta_0) \right\|^p \right] \right) \\ &\lesssim q \mathbb{E}_{\pi_\gamma}^{2/p} \left[\left\| g_m(\theta_0) \right\|^p \right]. \end{aligned}$$

It remains to upper bound the moment $\mathbb{E}_{\pi_\gamma}^{2/p} \left[\left\| g_m(\theta_0) \right\|^p \right]$. In order to do this, we first note that due to the duality theorem (Douc et al., 2018, Theorem 20.1.2.), we get that for any $k \in \mathbb{N}$,

$$\left\| \mathbb{Q}^{km} \psi(\theta) - \pi_\gamma(\psi) \right\| \leq \frac{1}{2} L_3 \mathbf{W}_c(\delta_\theta \mathbb{Q}_\gamma^{km}, \pi_\gamma) \leq 2 L_3 (1/2)^k \mathbf{W}_c(\delta_\theta, \pi_\gamma),$$

where the last inequality is due to Proposition 1. Hence, applying the definition of $g_m(\theta)$ in (106), we obtain that

$$\mathbb{E}_{\pi_\gamma}^{1/p} [\|g_m(\theta_0)\|^p] \leq \sum_{k=0}^{\infty} \mathbb{E}_{\pi_\gamma}^{1/p} [\|Q_\gamma^{km} \bar{\psi}(\theta)\|^p] \leq 2L_3 \sum_{k=0}^{\infty} (1/2)^k \mathbb{E}_{\pi_\gamma}^{1/p} [\{\mathbf{W}_c(\delta_\theta, \pi_\gamma)\}^p].$$

To control the latter term, we simply apply the definition of $\mathbf{W}_c(\delta_\theta, \pi_\gamma)$ and a cost function $c(\theta, \theta')$ together with C1(2p), we get

$$\begin{aligned} \mathbb{E}_{\pi_\gamma}^{1/p} [\{\mathbf{W}_c(\delta_\theta, \pi_\gamma)\}^p] &\lesssim \left(\int_{\mathbb{R}^d \times \mathbb{R}^d} \|\theta - \theta'\|^p \left(\|\theta - \theta^*\| + \|\theta' - \theta^*\| + \frac{\tau_2 \sqrt{\gamma}}{\sqrt{\mu}} \right)^p \pi_\gamma(d\theta) \pi_\gamma(d\theta') \right)^{1/p} \\ &\leq \left(\int \|\theta - \theta'\|^{2p} \pi_\gamma(d\theta) \pi_\gamma(d\theta') \right)^{1/2p} \left(\int \left(\|\theta - \theta^*\| + \|\theta' - \theta^*\| + \frac{\tau_2 \sqrt{\gamma}}{\sqrt{\mu}} \right)^{2p} \pi_\gamma(d\theta) \pi_\gamma(d\theta') \right)^{1/2p} \\ &\lesssim \frac{D_{\text{last}, 2p} \tau_{2p}^2 \gamma}{\mu}. \end{aligned}$$

Combining now the bounds above in (110), we get that

$$\mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=1}^q \{g_m(\theta_{km}) - Q_\gamma^m g_m(\theta_{(k-1)m})\} \right\|^p \right] \lesssim \frac{D_{\text{last}, 2p} L_3 \tau_{2p}^2 \gamma \sqrt{q}}{\mu}, \quad (111)$$

and, hence, substituting into (104), we get

$$\mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=0}^{n-1} \bar{\psi}(\theta_k) \right\|^p \right] \lesssim \frac{D_{\text{last}, 2p} L_3 \tau_{2p}^2 \gamma \sqrt{q} m}{\mu} + \frac{L D_{\text{last}, 2p} \tau_{2p}^2 \gamma m}{2\mu}. \quad (112)$$

Now the statement follows from the definition of $m = m(\gamma)$ in (103) and $q = \lfloor n/m \rfloor \leq n/m$.

D.2 VERSION OF PROPOSITION 8 FOR ARBITRARY INITIAL DISTRIBUTION ν .

In order to prove Theorem 23, we need a generalization of Proposition 8 for arbitrary initial distribution ν . We provide this result below.

Lemma 24. *Under assumptions of Proposition 8 for any $\gamma \in (0, 1/(LC_{\text{step}, 6})]$, $n \in \mathbb{N}$ and initial distribution ν , it holds that*

$$\begin{aligned} \mathbb{E}_\nu^{1/p} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \right\|^p \right] &\leq \mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\tilde{\theta}_k) - \pi_\gamma(\psi)\} \right\|^p \right] \\ &\quad + \frac{L_3(1 - \gamma\mu)^{(n+1)/2}}{\gamma\mu} \left(\mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^{2p}] + \frac{D_{\text{last}, 2p} \gamma \tau_{2p}^2}{\mu} \right). \end{aligned}$$

Proof. We consider the synchronous coupling contraction defined in (46) and denote by K_γ the corresponding coupling kernel. Hence, we have

$$\begin{aligned} \mathbb{E}_\nu^{1/p} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \right\|^p \right] &= (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \pi_\gamma(\psi)\} \right\|^p \right])^{1/p} \\ &\leq \mathbb{E}_{\pi_\gamma}^{1/p} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\tilde{\theta}_k) - \pi_\gamma(\psi)\} \right\|^p \right] + (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \psi(\tilde{\theta}_k)\} \right\|^p \right])^{1/p}. \end{aligned}$$

It remains to bound the last term in the inequality above. Applying Minkowski's inequality together with Lemma 22, we get

$$(\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \psi(\tilde{\theta}_k)\} \right\|^p \right])^{1/p} \leq \frac{L_3}{2} \sum_{k=n+1}^{2n} (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} [c^p(\theta_k, \tilde{\theta}_k)])^{1/p}.$$

Using Hölder's and Minkowski's inequalities together with **C1**($2p$) and Lemma 19, we obtain

$$\begin{aligned}
& (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} [c^p(\theta_k, \tilde{\theta}_k)])^{1/p} \\
& \leq (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} [\|\theta_k - \tilde{\theta}_k\|^{2p}])^{1/(2p)} (\mathbb{E}_{\pi_\gamma}^{1/(2p)} [\|\tilde{\theta}_k - \theta^*\|^{2p}] + \mathbb{E}_\nu^{1/(2p)} [\|\theta_k - \theta^*\|^{2p} + \frac{2^{3/2}\gamma^{1/2}\tau_2}{\mu^{1/2}}]) \\
& \leq (1 - \gamma\mu)^{k/2} (\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} [\|\theta_0 - \tilde{\theta}_0\|^{2p}])^{1/(2p)} (\mathbb{E}_\nu^{1/(2p)} [\|\theta_0 - \theta^*\|^{2p}] + \frac{2D_{\text{last}, 2p}^{1/2}\gamma^{1/2}\tau_{2p}}{\mu^{1/2}} + \frac{2^{3/2}\gamma^{1/2}\tau_2}{\mu^{1/2}}) \\
& \lesssim (1 - \gamma\mu)^{k/2} \left(\frac{D_{\text{last}, 2p}\gamma\tau_{2p}^2}{\mu} + \mathbb{E}_\nu^{1/p} \|\theta_0 - \theta^*\|^{2p} \right)
\end{aligned}$$

Combining all inequalities above, we get

$$\left(\mathbb{E}_{\nu, \pi_\gamma}^{K_\gamma} \left[\left\| \sum_{k=n+1}^{2n} \{\psi(\theta_k) - \psi(\theta'_k)\} \right\|^p \right] \right)^{1/p} \lesssim \frac{L_3(1 - \gamma\mu)^{(n+1)/2}}{\gamma\mu} \left(\mathbb{E}_\nu^{1/p} [\|\theta_0 - \theta^*\|^{2p}] + \frac{D_{\text{last}, 2p}\gamma\tau_{2p}^2}{\mu} \right),$$

and the statement follows. \square

E EXPERIMENTAL DETAILS

We recall the error representation (38), and obtain with simple algebra:

$$\begin{aligned}
\mathbf{H}^*(\bar{\theta}_n^{(\gamma)} - \theta^*) + n^{-1} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*) &= \frac{\theta_{n+1}^{(\gamma)} - \theta^*}{\gamma n} - \frac{\theta_{2n}^{(\gamma)} - \theta^*}{\gamma n} \\
&- \frac{1}{n} \sum_{k=n+1}^{2n} \{\varepsilon_{k+1}(\theta_k^{(\gamma)}) - \varepsilon_{k+1}(\theta^*)\} - \frac{1}{n} \sum_{k=n+1}^{2n} \psi(\theta_k^{(\gamma)}) - \frac{1}{n} \sum_{k=n+1}^{2n} G(\theta_k^{(\gamma)}). \quad (113)
\end{aligned}$$

Under **A3**(6), the statistics $\frac{1}{n} \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*)$ is a sum of independent random variables, and

$$n^{-2} \mathbb{E} \left[\left\| \sum_{k=n+1}^{2n} \varepsilon_{k+1}(\theta^*) \right\|^2 \right] = \frac{\text{Tr } \Sigma_\varepsilon^*}{n}.$$

Hence, in order to trace the rate of the second-order terms in (40), it is enough to find the decay rate of the right-hand side in (113). We select different sample sizes $n = 250 \times 2^k$, where $k = 0, \dots, 14$, and run the SGD procedure (2) based on the constant step sizes γ and 2γ , selecting $\gamma = 1/\sqrt{n}$. Then we construct the associated estimates $\bar{\theta}_n^{(\gamma)}$ and $\bar{\theta}_n^{(2\gamma)}$. We conduct $M = 320$ independent parallel runs to approximate the expectations. Code to reproduce experiments is provided at https://github.com/svsamsonov/richardson_romberg_example.