
Gaussian Approximation and Multiplier Bootstrap for Stochastic Gradient Descent

Marina Sheshukova¹ Sergey Samsonov¹ Denis Belomestny^{2,1} Éric Moulines^{3,4}
Qi-Man Shao⁵ Zhuo-Song Zhang⁵ Alexey Naumov^{1,6}

¹HSE University, ²Duisburg-Essen University ³LRE EPITA ⁴MBZUAI CMS Division

⁵Southern University of Science and Technology ⁶Steklov Mathematical Institute of Russian Academy of Sciences

Abstract

In this paper, we establish the non-asymptotic validity of the multiplier bootstrap procedure for constructing the confidence sets using the Stochastic Gradient Descent (SGD) algorithm. Under appropriate regularity conditions, our approach avoids the need to approximate the limiting covariance of Polyak-Ruppert SGD iterates, which allows us to derive approximation rates in convex distance of order up to $1/\sqrt{n}$. Notably, this rate can be faster than the one that can be proven in the Polyak-Juditsky central limit theorem. To our knowledge, this provides the first fully non-asymptotic bound on the accuracy of bootstrap approximations in SGD algorithms. Our analysis builds on the Gaussian approximation results for nonlinear statistics of independent random variables.

1 INTRODUCTION

Stochastic Gradient Descent (SGD) is a widely used first-order optimization method well suited for large datasets and online learning. The algorithm has attracted significant attention; see, e.g. (Polyak and Juditsky, 1992; Nemirovski et al., 2009; Moulines and Bach, 2011). SGD aims to solve the optimization problem:

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

where ξ is a random variable defined on a measurable space $(\mathcal{Z}, \mathcal{Z})$. Instead of the exact gradient $\nabla f(\theta)$, the

algorithm accesses only unbiased stochastic estimates $F(\theta, \xi)$.

Throughout this work, we focus on strongly convex objective functions and denote by θ^* the unique minimizer of (1). The iterates θ_k , $k \in \mathbb{N}$, generated by SGD follow the recursive update:

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

where $\{\alpha_k\}_{k \in \mathbb{N}}$ is a sequence of step sizes (or learning rates), which may be either diminishing or constant, and $\{\xi_k\}_{k \in \mathbb{N}}$ is an i.i.d. sequence sampled from \mathbb{P}_ξ . Theoretical properties of SGD, particularly in the convex and strongly convex settings, have been extensively studied; see, e.g., (Nesterov, 2004; Moulines and Bach, 2011; Bubeck et al., 2015; Lan, 2020). Many optimization algorithms build on the recurrence (2) to accelerate the convergence of the sequence θ_k to θ^* . Notable examples include momentum acceleration (Qian, 1999), variance reduction techniques (Defazio et al., 2014; Schmidt et al., 2017), and averaging methods. In this work, we focus on Polyak-Ruppert averaging, originally proposed in (Ruppert, 1988) and (Polyak and Juditsky, 1992), which improves convergence by averaging the SGD iterates in (2). Specifically, the estimator is defined as

$$\bar{\theta}_n = \frac{1}{n} \sum_{i=0}^{n-1} \theta_i, \quad n \in \mathbb{N}. \quad (3)$$

It has been established (see (Polyak and Juditsky, 1992, Theorem 3)) that, under appropriate conditions on the objective function f , the noisy gradient estimates F , and the step sizes α_k , the sequence of averaged iterates $\{\bar{\theta}_n\}_{n \in \mathbb{N}}$ is asymptotically normal:

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

where \xrightarrow{d} denotes convergence in distribution, and $\mathcal{N}(0, \Sigma_\infty)$ is a zero-mean Gaussian distribution with covariance matrix Σ_∞ , defined later in Section 2.2. This result raises two key questions:

- (i) What is the rate of convergence in (4)?

- (ii) How can (4) be leveraged to construct confidence sets for θ^* , given that Σ_∞ is unknown in practice?

In our paper, we aim to answer both questions. To quantify convergence rates in (4), we use convex distance, defined for random vectors $X, Y \in \mathbb{R}^d$ as

$$d_C(X, Y) = \sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)|, \quad (5)$$

where $\mathcal{C}(\mathbb{R}^d)$ denotes the collection of convex subsets of \mathbb{R}^d . The supremum in (5) can be taken over different classes of sets, leading to various probability metrics. In particular, more restrictive classes such as rectangles give rise to metrics with logarithmic dependence on the dimension d (see, e.g., (Kojevnikov and Song, 2022; Chernozhukov et al., 2013)). The choice of the class of sets is often driven by the needs of a particular application and may affect the dependence of the resulting bounds on the dimension d . In particular, even for normal approximation of linear statistics, the dimensional dependence may vary depending on the chosen class. Results for convex distance can be found in (Bentkus, 2003), while results for rectangles are available in (Kojevnikov and Song, 2022; Chernozhukov et al., 2013).

Shao and Zhang (2022) derive Berry-Esseen-type bounds for $d_C(\sqrt{n}(\bar{\theta}_n - \theta^*), \mathcal{N}(0, \Sigma_n))$, where Σ_n is the covariance matrix of the linearized counterpart of (2); see the precise definition below in (14). We complement this result with the rates of convergence in (4). We also establish a lower bound on the convex distance $d_C(\sqrt{n}(\bar{\theta}_n - \theta^*), \mathcal{N}(0, \Sigma_\infty))$. This result shows that, for certain step-size sequences α_k in (2), the Gaussian approximation with covariance Σ_∞ is less accurate than an approximation based on another covariance, in particular Σ_n . A similar phenomenon has been reported in the bootstrap literature for i.i.d. data outside the context of gradient methods; see (Shao and Tu, 1995, Theorem 3.11).

One popular approach for solving (ii) is based on the *plug-in* methods (Chen et al., 2020, 2021), which aim to construct an estimator $\hat{\Sigma}_n$ of Σ_∞ directly. Theoretical guarantees for these methods typically focus on non-asymptotic bounds for how close $\hat{\Sigma}_n$ is to Σ_∞ , often in terms of $\mathbb{E}[\|\hat{\Sigma}_n - \Sigma_\infty\|]$. At the same time, the analysis of these methods bypasses item (i) and the issues related to the rate of convergence in (4). In our paper, we present, to our knowledge, the first fully non-asymptotic analysis of a procedure for constructing confidence intervals based on the bootstrap approach (Efron, 1992; Fang et al., 2018), which avoids direct approximation of Σ_∞ . Moreover, theoretical analysis of the underlying procedure, together with results on normal approximation with $\mathcal{N}(0, \Sigma_\infty)$ from (i), shows that the same approximation rate cannot be achieved

by plug-in methods, at least for a certain range of step sizes α_k in (2). Our key contributions are as follows:

- We establish the non-asymptotic validity of the multiplier bootstrap procedure introduced in (Fang et al., 2018). Under suitable regularity conditions, our bounds show that the distribution of $\sqrt{n}(\bar{\theta}_n - \theta^*)$ can be approximated, up to logarithmic factors, at rate $1/n^{\gamma-1/2}$ for step sizes of the form $\alpha_k = c_0/(k + k_0)^\gamma$ with $\gamma \in (1/2, 1)$. To our knowledge, this is the first bound on the accuracy of bootstrap approximations in SGD. Notably, this rate can be faster than the one obtained in (4). Our results improve upon recent works (Samsonov et al., 2024; Wu et al., 2024), which addressed the convergence rate in similar procedures for the LSA algorithm and TD learning, respectively.
- Our analysis of the multiplier bootstrap procedure reveals an important distinction: unlike plug-in estimators, the validity of the bootstrap method does not depend on approximating $\sqrt{n}(\bar{\theta}_n - \theta^*)$ by $\mathcal{N}(0, \Sigma_\infty)$. Instead, it requires approximation by $\mathcal{N}(0, \Sigma_n)$ with the matrix Σ_n being the covariance matrix of the linearized counterpart of (2). The structure of Σ_n is central to our analysis, both for the rate in (4) and for the non-asymptotic bootstrap validity. Precise definitions are provided in Section 2.2.
- We analyze the Polyak–Ruppert averaged SGD iterates (3) for strongly convex minimization problems and establish Gaussian approximation rates in (4) with respect to the convex distance. Specifically, we show that $d_C(\sqrt{n}(\bar{\theta}_n - \theta^*), \mathcal{N}(0, \Sigma_\infty))$ is of order $n^{-1/4}$ for step sizes $\alpha_k = c_0/(k + k_0)^{3/4}$ with appropriately chosen c_0 and k_0 . Our proof relies on the techniques of (Shao and Zhang, 2022) and (Wu et al., 2024). We further provide a lower bound showing that this rate of normal approximation with $\mathcal{N}(0, \Sigma_\infty)$ is tight in the regime $\alpha_k = c_0/(k + k_0)^\gamma$ with $\gamma \geq 3/4$.

Notations. Throughout this paper, we use the following notations. For a matrix $A \in \mathbb{R}^{d \times d}$ and a vector $x \in \mathbb{R}^d$, we denote by $\|A\|$ and $\|x\|$ their spectral norm and Euclidean norm, respectively. We also write $\|A\|_F$ for the Frobenius norm of matrix A . Given a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we write $\nabla f(\theta)$ and $\nabla^2 f(\theta)$ for its gradient and Hessian at a point θ . We use the standard abbreviations "i.i.d." for "independent and identically distributed" and "w.r.t." for "with respect to".

Literature review. The asymptotic behavior of SGD, including the asymptotic normality of the estimator $\bar{\theta}_n$ and its almost sure convergence, has been studied for smooth and strongly convex minimization problems (Polyak and Juditsky, 1992; Kushner and Yin, 2003; Benveniste et al., 2012). Optimal mean-squared error

(MSE) bounds for $\theta_n - \theta^*$ and $\bar{\theta}_n - \theta^*$ were first derived in (Nemirovski et al., 2009) for smooth and strongly convex objectives, and later refined in (Moulines and Bach, 2011). The constant step size regime for strongly convex problems has been analyzed in (Dieuleveut et al., 2020; Li et al., 2025). High-probability bounds for SGD iterates were established in (Rakhlin et al., 2012) and later extended in (Harvey et al., 2019), both addressing non-smooth and strongly convex minimization.

It is important to note that the results discussed above do not directly yield convergence rates for $\sqrt{n}(\theta_n - \theta^*)$ to $\mathcal{N}(0, \Sigma_\infty)$ in terms of convex or Wasserstein distance. Among the relevant contributions in this direction, we highlight recent works (Srikant, 2025; Samsonov et al., 2024; Wu et al., 2024), which provide quantitative bounds on the convergence rate in (4) for iterates of temporal difference (TD) learning and general linear stochastic approximation (LSA) schemes. These algorithms, however, do not necessarily reduce to SGD with a quadratic objective f , since the system matrix in LSA is not required to be symmetric. Convergence rates of order up to $1/\sqrt{n}$ were established in (Anastasiou et al., 2019) for a class of smooth test functions. The recent work (Agrawalla et al., 2023) establishes Berry–Esseen bounds of order up to $n^{-1/4}$ for the last iterate of SGD in the linear regression setting.

Bootstrap methods for i.i.d. observations were first introduced in (Efron, 1992). In the context of SGD, Fang et al. (2018) proposed the multiplier bootstrap procedure for constructing confidence intervals for θ^* and established its asymptotic validity. The same method was later analyzed in (Samsonov et al., 2024) for the LSA algorithm, where a rate of $n^{-1/4}$ was obtained for the accuracy of approximating the distribution of $\sqrt{n}(\bar{\theta}_n - \theta^*)$ by the distribution of its bootstrap-world counterpart in terms of convex distance.

A naive approach for constructing confidence intervals is to perform multiple independent runs of the algorithm to empirically acquire distributional information and then construct a confidence interval, as discussed in (Zhu et al., 2024). A popular class of methods for constructing confidence sets for θ^* relies on estimating the asymptotic covariance matrix Σ_∞ . Plug-in and batch-mean estimators of Σ_∞ have attracted a lot of attention (Chen et al., 2020, 2021; Chang et al., 2026), especially in the setting when the stochastic estimates of Hessian are available. Estimators of Σ_∞ based on the batch-mean method and its online variant were studied in (Chen et al., 2020) and (Zhu et al., 2023). Li et al. (2022b) studied the asymptotic validity of the plug-in estimator for Σ_∞ in the local SGD framework. The analysis in (Zhong et al., 2023) refined guarantees for both the multiplier bootstrap and batch-mean estimators of Σ_∞ in nonconvex problems. However, these

contributions typically establish recovery rates for Σ_∞ but only prove the asymptotic validity of the resulting confidence intervals. A notable exception is the recent work (Wu et al., 2024), which studied the temporal-difference (TD) learning algorithm. There, the authors provided a fully non-asymptotic analysis, obtaining the approximation rate $n^{-1/3}$ for the distribution of $\sqrt{n}(\bar{\theta}_n - \theta^*)$ in terms of convex distance.

2 MAIN RESULTS

This section establishes the non-asymptotic validity of the multiplier bootstrap method proposed in (Fang et al., 2018). We restrict attention to smooth and strongly convex minimization problems, following the frameworks of (Moulines and Bach, 2011), (Anastasiou et al., 2019), and (Shao and Zhang, 2022). The procedure is based on perturbing the trajectory (2). Let $\mathcal{W}^{n-1} = \{w_\ell\}_{1 \leq \ell \leq n-1}$ be i.i.d. random variables with distribution \mathbb{P}_w , each satisfying $\mathbb{E}[w_1] = 1$ and $\text{Var}[w_1] = 1$. We assume that \mathcal{W}^{n-1} is independent of $\Xi^{n-1} = \{\xi_\ell\}_{1 \leq \ell \leq n-1}$. We then use \mathcal{W}^{n-1} to construct randomly perturbed trajectories of the SGD dynamics (2):

$$\begin{aligned} \theta_k^b &= \theta_{k-1}^b - \alpha_k w_k F(\theta_{k-1}^b, \xi_k), & \theta_0^b &= \theta_0 \in \mathbb{R}^d, \\ \bar{\theta}_n^b &= n^{-1} \sum_{k=0}^{n-1} \theta_k^b, & n &\geq 1. \end{aligned} \quad (6)$$

When generating different weights w_k , we obtain samples from the conditional distribution of $\bar{\theta}_n^b$ given the data Ξ^{n-1} . We further denote

$$\mathbb{P}^b = \mathbb{P}(\cdot \mid \Xi^{n-1}), \quad \mathbb{E}^b = \mathbb{E}(\cdot \mid \Xi^{n-1}).$$

The core principle behind the bootstrap procedure (6) is that the "bootstrap world" probabilities $\mathbb{P}^b(\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) \in B)$ are close to $\mathbb{P}(\sqrt{n}(\bar{\theta}_n - \theta^*) \in B)$ for $B \in \mathcal{C}(\mathbb{R}^d)$. Formally, we say that the procedure (6) is asymptotically valid if

$$\sup_{B \in \mathcal{C}(\mathbb{R}^d)} \left| \mathbb{P}^b(\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) \in B) - \mathbb{P}(\sqrt{n}(\bar{\theta}_n - \theta^*) \in B) \right|$$

converges to 0 in \mathbb{P} -probability as $n \rightarrow \infty$. This result was obtained in (Fang et al., 2018) under assumptions close to the original paper (Polyak and Juditsky, 1992). While an analytical expression for $\mathbb{P}^b(\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) \in B)$ is unavailable, it can be approximated via Monte Carlo simulations by generating M perturbed trajectories according to (6). Standard arguments (see, e.g., (Shao, 2003, Section 5.1)) suggest that the accuracy of this Monte Carlo approximation scales as $\mathcal{O}(M^{-1/2})$ when generating M parallel perturbed trajectories in (6).

Assumptions. We impose the following regularity conditions on the objective function f :

A1. The function f is twice continuously differentiable and L_1 -smooth on \mathbb{R}^d ; that is, there exists a constant $L_1 > 0$ such that for any $\theta, \theta' \in \mathbb{R}^d$,

$$\|\nabla f(\theta) - \nabla f(\theta')\| \leq L_1 \|\theta - \theta'\|.$$

Moreover, f is assumed to be μ -strongly convex on \mathbb{R}^d ; that is, there exists a constant $\mu > 0$ such that for any $\theta, \theta' \in \mathbb{R}^d$,

$$(\mu/2)\|\theta - \theta'\|^2 \leq f(\theta) - f(\theta') - \langle \nabla f(\theta'), \theta - \theta' \rangle.$$

A1 implies the following bound on the Hessian of f :

$$\mu \mathbf{I}_d \preceq \nabla^2 f(\theta) \preceq L_1 \mathbf{I}_d,$$

for all $\theta \in \mathbb{R}^d$. We next state the assumptions on $F(\theta, \xi)$. Specifically, we write

$$F(\theta_{k-1}, \xi_k) = \nabla f(\theta_{k-1}) + \zeta_k, \quad \zeta_k := \zeta(\theta_{k-1}, \xi_k),$$

so that $\{\zeta_k\}_{k \in \mathbb{N}}$ is a sequence of d -dimensional random vectors whose distribution may depend on θ_{k-1} . The SGD recursion (2) then takes the form

$$\theta_k = \theta_{k-1} - \alpha_k (\nabla f(\theta_{k-1}) + \zeta_k), \quad \theta_0 \in \mathbb{R}^d. \quad (7)$$

We impose the following assumption on the noise sequence ζ_k :

A 2. For each $k \geq 1$, ζ_k admits the decomposition $\zeta_k = \eta(\xi_k) + g(\theta_{k-1}, \xi_k)$, where

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

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

$$\mathbb{E}^{1/p}[\|g(\theta, \xi) - g(\theta', \xi)\|^p] \leq L_2 \|\theta - \theta'\|, \quad (8)$$

and $g(\theta^*, z) = 0$ for all $z \in \mathcal{Z}$.

(iii) There exists $C_\xi > 0$ such that for any $\theta \in \mathbb{R}^d$, the random vector $g(\theta, \xi) + \eta(\xi)$ is sub-Gaussian with variance proxy C_ξ^2 ; that is, for any $v \in \mathbb{R}^d$,

$$\mathbb{E}[\exp\{\langle g(\theta, \xi) + \eta(\xi), v \rangle\}] \leq \exp\{\|v\|^2 C_\xi^2 / 2\}.$$

Discussion. As an example of a sequence ζ_k satisfying conditions (i) and (ii) in A2, consider the case where the stochastic estimates $F(\theta, \xi)$ satisfy:

1. $\mathbb{E}[F(\theta, \xi)] = \nabla f(\theta)$ for all $\theta \in \mathbb{R}^d$;
2. $\mathbb{E}^{1/p}[\|F(\theta, \xi) - F(\theta', \xi)\|^p] \leq L \|\theta - \theta'\|$.

In this case, (i) and (ii) in A2 hold with $\eta(\xi) = F(\theta^*, \xi)$ and $g(\theta, \xi) = F(\theta, \xi) - F(\theta^*, \xi) - \nabla f(\theta)$. Condition (ii) in A2 is often imposed when studying averaged iterates;

see (Moulines and Bach, 2011, Assumption H2) and (Dieuleveut et al., 2020; Sheshukova et al., 2025). It is also possible to adapt our arguments to the weaker assumption

$$\mathbb{E}^{1/p}[\|g(\theta, \xi) - g(\theta', \xi)\|^p] \leq L_2 \|\theta - \theta'\|^\beta,$$

for some $1/2 < \beta < 1$, at the cost of a slower approximation rate. The boundary case $\beta = 1/2$, which arises for example in quantile regression, is not covered by our analysis and requires different techniques. We refer to (Chen et al., 2023, 2025; Cai et al., 2025) for non-asymptotic results in this setting. Furthermore, our proof strategy remains applicable if only a fixed number of moments $p \geq 4$ (rather than $p = \log n$) is available, at the cost of a slower rate with respect to n in the final estimates. Similarly, one can allow L_2 in (8) to depend polynomially on p . In this case, our proof approach remains valid, but yields additional polynomial dependence on $\log n$ in the final rate of Theorem 1.

The assumption A2-(iii) is crucial to establish high-order moment bounds:

$$\mathbb{E}^{1/p}[\|\theta_k - \theta^*\|^p] \quad \text{and} \quad \mathbb{E}^{1/p}[\|\theta_k^b - \theta^*\|^p], \quad (9)$$

see Lemma 14 in the Appendix. Our proof generalizes the argument of (Harvey et al., 2019, Theorem 4.1), which requires the noise variables ζ_k in (7) to be almost surely bounded. This argument was later generalized in (Madden et al., 2024) to the setting where ζ_k is conditionally sub-Gaussian given $\mathcal{F}_{k-1} = \sigma(\theta_i, i \leq k-1)$, with variance proxy uniformly (in θ) bounded by a constant factor. This is exactly the setting considered in A2-(iii). Such an assumption is widely used in the literature; see (Nemirovski et al., 2009; Hazan and Kale, 2014) and the remarks in (Harvey et al., 2019). We believe that this assumption can be generalized for sub-Weibull noise setting (see (Madden et al., 2024, Definition 5)), at the expense of additional technicalities (see (Madden et al., 2024, Theorem 9)). Some of the authors considering bounds of type (9), such as Rakhlin et al. (2012), imposed the stronger assumption that $\sup_{\theta \in \mathbb{R}^d} \|F(\theta, \xi)\|$ is almost surely bounded. A less restrictive approach might be to consider updates involving gradient clipping; see, e.g., (Sadiev et al., 2023), or a version of SGD algorithm with projections as in (Rakhlin et al., 2012). However, both approaches change the limiting covariance matrix Σ_∞ and the leading term in Gaussian approximation for $\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n)$. Both approaches would introduce an additional level of technical difficulties, so we leave a detailed study of these schemes for the future work.

We further impose the following condition on the Hessian matrix $\nabla^2 f(\theta)$ at θ^* :

A 3. *There exist constants $L_3, \beta > 0$ such that for all θ with $\|\theta - \theta^*\| \leq \beta$, we have*

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

A 3 ensures that the Hessian of f is Lipschitz continuous in a neighborhood of θ^* . Similar assumptions have been considered in (Shao and Zhang, 2022; Anastasiou et al., 2019), as well as in other works on first-order optimization methods; see, e.g., (Li et al., 2022a). Several studies on the non-asymptotic analysis of SGD impose stronger smoothness conditions, such as bounded derivatives of f up to order 4, as in (Dieuleveut et al., 2020). We also impose the following assumption on the bootstrap weights w_i used in the algorithm:

A 4. *There exist constants $0 < W_{\min} < W_{\max} < +\infty$ such that $W_{\min} \leq w_1 \leq W_{\max}$ a.s.*

The original work (Fang et al., 2018) also required positive bootstrap weights w_i . We impose boundedness of w_i in order to derive the high-probability bound in Lemma 14. An explicit example of a distribution satisfying **A 4** and the constraints $\mathbb{E}[w_1] = 1$ and $\text{Var}[w_1] = 1$ is given in Appendix A.3. Finally, we impose the following condition on the step sizes α_k and sample size n :

A 5. *Let $\alpha_k = c_0(k_0 + k)^{-\gamma}$, where $\gamma \in (1/2, 1)$ and c_0, k_0 satisfy $2c_0W_{\max}^2L_1^2 \leq 1$,*

$$k_0 \geq \max \left\{ \left(\frac{2\gamma}{\mu c_0 W_{\min}} \right)^{1/(1-\gamma)}, \left(\frac{1}{\mu W_{\min}} \right)^{1/\gamma} \right\}.$$

A 6. *The number of observations n is assumed to be sufficiently large. Precise expressions for n are provided in Appendix A (see **A' 6**).*

The particular bound on k_0 in **A 5** arises from the high-order moment bounds (see Lemma 14 in the Appendix). We further discuss the lower bound on the number of observations imposed in **A 6** in the proof of Theorem 1.

Discussion of assumptions Most of the theoretical assumptions used in our analysis (namely, **A 1**, **A 2(i,ii)**, **A 3**) are standard in the stochastic approximation literature and are of the same nature as the conditions appearing in the classical Polyak–Juditsky CLT (Polyak and Juditsky, 1992). Similar assumptions are also routinely used in several recent works on non-asymptotic analysis of stochastic approximation (Li et al., 2022a; Shao and Zhang, 2022; Sheshukova et al., 2025). The additional assumptions **A 2(iii)**, **A 4**, and **A 5** are required to obtain high-probability guarantees for the last SGD iterate, while **A 6** ensures that $\lambda_{\min}(\Sigma_n^b)$ is separated from zero with high probability.

2.1 Non-asymptotic multiplier bootstrap validity

Theorem 1. *Assume **A 1–A 6**. Then, with \mathbb{P} -probability at least $1 - 2/n$, we have*

$$\begin{aligned} \sup_{B \in \mathcal{C}(\mathbb{R}^d)} & \left| \mathbb{P}^b(\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) \in B) - \mathbb{P}(\sqrt{n}(\bar{\theta}_n - \theta^*) \in B) \right| \\ & \leq \frac{C_1 \sqrt{\log n}}{n^{1/2}} + \frac{C_2 \log n}{n^{\gamma-1/2}} + \frac{C_3 (\log n)^{3/2}}{n^{\gamma/2}}, \end{aligned}$$

where C_1, C_2, C_3 are given in Appendix A.1, equation (25).

Remark 1. *The constants C_1, C_2, C_3 in Theorem 1 depend on the problem dimension d , as well as on $\|\theta_0 - \theta^*\|$. To make the dependence on d explicit, we assume the natural scaling $\|\theta_0 - \theta^*\| \lesssim \sqrt{d}$. These dimension dependence aligns with the one assumed in (Shao and Zhang, 2022). We also assume that C_ξ from **A 2-(iii)** is dimension-free. Then Theorem 2 yields*

$$\begin{aligned} \sup_{B \in \mathcal{C}(\mathbb{R}^d)} & \left| \mathbb{P}^b(\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) \in B) - \mathbb{P}(\sqrt{n}(\bar{\theta}_n - \theta^*) \in B) \right| \\ & \lesssim \frac{(d^2 + d^{3/2} \log(2d))\sqrt{\log n}}{n^{1/2}} + \frac{d^{3/2} \log n}{n^{\gamma-1/2}} \\ & \quad + \frac{d^{3/2} (\log n)^{3/2}}{n^{\gamma/2}}. \end{aligned}$$

Here, the notation \lesssim indicates that the inequality holds up to constants independent of n and d .

Remark 2. *The result of Theorem 1 can also be established for the step size $\alpha_k = c_0/(k + k_0)$. The required Gaussian approximation with covariance matrix Σ_n was proved in (Shao and Zhang, 2022). The only differences compared to Theorem 1 are the additional $\log n$ factors in the bound and modified conditions on c_0 and k_0 in **A 5**.*

Remark 3. *In Appendix H we provide numerical results on logistic and linear regressions, illustrating the coverage probabilities achieved by the multiplier bootstrap approach and the batch-mean approach (Roy and Balasubramanian, 2023).*

Proof sketch of Theorem 1. The proof of non-asymptotic bootstrap validity is based on the Gaussian approximation performed both in the "real" world and bootstrap world together with an appropriate Gaussian comparison inequality: Here, Σ and Σ^b are covariance matrices to be specified later. To understand the origin of the Gaussian approximation, we linearize the statistics $\sqrt{n}(\bar{\theta}_n - \theta^*)$ and $\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n)$. We detail the derivation for $\sqrt{n}(\bar{\theta}_n - \theta^*)$, while analogous arguments for $\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n)$ are provided in Section 2.3.

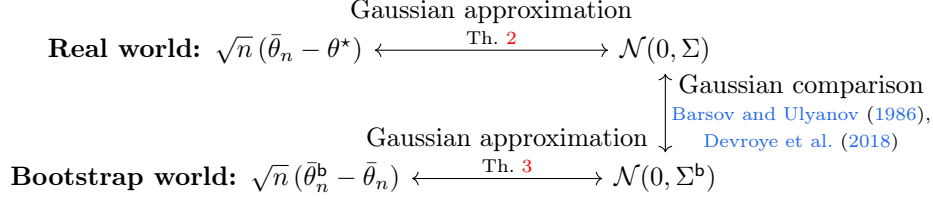


Figure 1: Gaussian approximations in the real and bootstrap worlds, and their comparison.

Let $G = \nabla^2 f(\theta^*)$. We expand $\sqrt{n}(\bar{\theta}_n - \theta^*)$ into a weighted sum of independent random vectors plus lower-order terms. By the Newton–Leibniz formula,

$$\nabla f(\theta) = G(\theta - \theta^*) + H(\theta), \quad (10)$$

where

$$H(\theta) = \int_0^1 (\nabla^2 f(\theta^* + t(\theta - \theta^*)) - G)(\theta - \theta^*) dt.$$

Note that $H(\theta)$ is of order $\|\theta - \theta^*\|^2$ (see Lemma 18). The recursion for the SGD error (7) writes as

$$\begin{aligned} \theta_k - \theta^* &= (\mathbf{I}_d - \alpha_k G)(\theta_{k-1} - \theta^*) \\ &\quad - \alpha_k (\eta(\xi_k) + g(\theta_{k-1}, \xi_k) + H(\theta_{k-1})). \end{aligned} \quad (11)$$

For $i \in \{0, \dots, n-1\}$, define

$$Q_i = \alpha_i \sum_{j=i}^{n-1} \prod_{k=i+1}^j (\mathbf{I}_d - \alpha_k G), \quad (12)$$

with the convention that empty products equal \mathbf{I}_d . Averaging (11) and rearranging yields $\sqrt{n}(\bar{\theta}_n - \theta^*) = W + D$ with

$$W = -\frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i \eta(\xi_i), \quad D = \sqrt{n}(\bar{\theta}_n - \theta^*) - W. \quad (13)$$

Here, W is a weighted sum of i.i.d. mean-zero random vectors with covariance

$$\Sigma_n = n^{-1} \sum_{k=1}^{n-1} Q_k \Sigma_\xi Q_k^\top, \quad (14)$$

and D is the remainder term, defined in Appendix B, equation (28). Moreover, in Appendix D.1 we show that Q_i can be approximated by $G^{-\top}$ and Σ_n approximates

$$\Sigma_\infty = G^{-1} \Sigma_\xi G^{-\top}.$$

We expect D not to affect the asymptotic distribution of the linear statistic W , which should be Gaussian by the central limit theorem. A key issue is the choice of the approximating Gaussian distribution $\mathcal{N}(0, \Sigma)$, with $\Sigma = \Sigma_n$ or Σ_∞ , and its bootstrap analogue Σ^b . This choice does not alter the bootstrap recursion (6), but only influences the rates of bootstrap approximation.

Fang et al. (2018) selected $\mathcal{N}(0, \Sigma_\infty)$ for their asymptotic analysis, and a similar approach was used in (Samsonov et al., 2024, Theorem 3) for the LSA setting. However, as shown in Theorem 4, this choice implies that the rate of normal approximation in bootstrap world is not faster than $n^{-1/4}$. In contrast, Theorem 2 and Theorem 3 demonstrate that rates up to $n^{-1/2}$ are attainable by taking $\Sigma = \Sigma_n$ in diagram 1 and using the corresponding bootstrap analogue. To conclude the proof, we apply the Gaussian comparison inequality; see Appendix A.1. A full proof of Theorem 1 is given in Appendix A. \square

Discussion. In (Samsonov et al., 2024), a counterpart of Theorem 1 was established for the LSA algorithm with an approximation rate of order $n^{-1/4}$, up to logarithmic factors. This rate is suboptimal, since the authors employed $\mathcal{N}(0, \Sigma_\infty)$ as the Gaussian approximation when proving bootstrap validity. A more recent work (Wu et al., 2024) improved the rate to $n^{-1/3}$ for the temporal-difference (TD) learning procedure with linear function approximation. Their approach relies on a direct estimate of Σ_∞ and yields a rate of order $n^{-1/3}$ when approximating the distribution of $\sqrt{n}(\bar{\theta}_n - \theta^*)$ by $\mathcal{N}(0, \hat{\Sigma}_n)$ with a suitably constructed estimator $\hat{\Sigma}_n$; see (Wu et al., 2024, Theorems 3.4 and 3.5). Chen et al. (2020) proposed a plug-in estimator $\hat{\Sigma}_n$ of Σ_∞ and showed that

$$\mathbb{E}[\|\hat{\Sigma}_n - \Sigma_\infty\|] \lesssim Cn^{-\gamma/2}, \quad \gamma \in (1/2, 1),$$

under weaker assumptions than those adopted in the present section. However, this bound is insufficient to establish an analogue of the Gaussian comparison result Lemma 1 for $\mathcal{N}(0, \hat{\Sigma}_n)$ and $\mathcal{N}(0, \Sigma_\infty)$ on a set of large \mathbb{P} -probability. Achieving such a result would require high-probability bounds on $\|\hat{\Sigma}_n - \Sigma_\infty\|$, which are more delicate and may necessitate additional assumptions beyond those in (Chen et al., 2020).

Furthermore, controlling the error between the distribution of $\sqrt{n}(\bar{\theta}_n - \theta^*)$ and $\mathcal{N}(0, \hat{\Sigma}_n)$ requires a Berry–Esseen type bound on the distance between $\sqrt{n}(\bar{\theta}_n - \theta^*)$ and $\mathcal{N}(0, \Sigma_\infty)$. As shown in Theorem 4, the approximation rate in this problem deteriorates as $\gamma \rightarrow 1$, creating an additional trade-off in the analysis of plug-in procedures based on estimating Σ_∞ .

This phenomenon highlights the fundamental difference between the multiplier bootstrap approach and the plug-in method of (Chen et al., 2020).

2.2 Gaussian approximation in the real world

In this section, we relax the assumptions A2 and A5. We introduce a family of assumptions, denoted A7(p) with $p \geq 2$, on the noise sequence ζ_k , and A8 on the step sizes α_k :

A7 (p). Conditions (i) and (ii) in A2 hold. In addition, there exists $\sigma_p > 0$ such that $\mathbb{E}^{1/p}[\|\eta(\xi_1)\|^p] \leq \sigma_p$.

A8. $\alpha_k = c_0/(k_0 + k)^\gamma$, where $\gamma \in (1/2, 1)$, $k_0 \geq 1$, and c_0 satisfies $2c_0L_1 \leq 1$.

Clearly, A5 implies A8, and A2 implies A7(p) for all $p \geq 2$ with $\sigma_p \lesssim C_\xi(\sqrt{d} + \sqrt{p})$. We also note that, for the results in this section, it suffices to assume that A2(ii) holds only for any $p \leq 4$.

The main result of this section establishes a Gaussian approximation for $\sqrt{n}(\bar{\theta}_n - \theta^*)$ with $\mathcal{N}(0, \Sigma_n)$. This refines the bounds of (Shao and Zhang, 2022, Theorem 3.4) and serves as a key step toward analyzing normal approximation with $\mathcal{N}(0, \Sigma_\infty)$ in Section 2.4.

Theorem 2. Assume A1, A3, A7(4), and A8. Then, for $Y \sim \mathcal{N}(0, \text{Id})$,

$$\text{d}_C\left(\sqrt{n}\Sigma_n^{-1/2}(\bar{\theta}_n - \theta^*), Y\right) \leq \frac{C_4}{\sqrt{n}} + \frac{C_5}{n^{\gamma-1/2}} + \frac{C_6}{n^{\gamma/2}}, \quad (15)$$

where C_4, C_5, C_6 are defined in Appendix B, equation (29). Moreover, since Σ_n is non-degenerate and the image of a convex set under a non-degenerate linear transformation is convex,

$$\begin{aligned} \text{d}_C\left(\sqrt{n}\Sigma_n^{-1/2}(\bar{\theta}_n - \theta^*), Y\right) \\ = \text{d}_C\left(\sqrt{n}(\bar{\theta}_n - \theta^*), \Sigma_n^{1/2}Y\right). \end{aligned}$$

Remark 4. The constants C_4, C_5, C_6 in Theorem 2 depend on the problem dimension d and on the parameters specified in A1–A7(4)–A3–A8. Moreover, C_5 depends on $\|\theta_0 - \theta^*\|$. To make the dependence on d explicit, we assume the natural scaling $\sigma_2, \sigma_4, \|\theta_0 - \theta^*\| \lesssim \sqrt{d}$. Under this assumption, Theorem 2 implies

$$\text{d}_C\left(\sqrt{n}\Sigma_n^{-1/2}(\bar{\theta}_n - \theta^*), Y\right) \lesssim \frac{d^2}{\sqrt{n}} + \frac{d^{3/2}}{n^{\gamma-1/2}} + \frac{d^{3/2}}{n^{\gamma/2}}.$$

Here, the notation \lesssim indicates that the inequality holds up to constants independent of n and d .

Remark 5. When $\gamma \rightarrow 1$, the correction terms in (15) scale as $\mathcal{O}(1/\sqrt{n})$, yielding an overall approximation rate approaching $1/\sqrt{n}$. For $\gamma \in (0, 1)$, we have

$1/n^{\gamma/2} < 1/n^{\gamma-1/2}$, so the term $C_5/n^{\gamma-1/2}$ dominates. We retain both terms in (15), as they correspond to the moments of the statistics

$$\frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i H(\theta_{i-1}) \quad \text{and} \quad \frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i g(\theta_{i-1}, \xi_i),$$

respectively. The first of these has nonzero mean, since $H(\theta_{i-1})$ is quadratic in $\|\theta_{i-1} - \theta^*\|^2$. In the case of constant step-size SGD, this term can be corrected using the Richardson–Romberg technique (Dieuleveut et al., 2020; Sheshukova et al., 2025). However, it is unclear if this technique can be extended to the diminishing step-size regime.

Proof sketch of Theorem 2. The starting point is the decomposition (13), which expresses $\sqrt{n}(\bar{\theta}_n - \theta^*) = W + D$ as the sum of a linear term W and a nonlinear remainder D . Establishing a Gaussian approximation for this statistic therefore reduces to analyzing the interplay between the leading linear part and the error term. Our analysis follows the framework of (Shao and Zhang, 2022). Consider independent random variables X_1, \dots, X_n taking values in a measurable space \mathcal{X} , and a d -dimensional statistic $T = T(X_1, \dots, X_n)$ that admits the decomposition

$$W = \sum_{\ell=1}^n Z_\ell, \quad D = T - W,$$

with $Z_\ell = r_\ell(X_\ell)$ for measurable maps $r_\ell : \mathcal{X} \rightarrow \mathbb{R}^d$. The component W captures the linear structure, while D accounts for the nonlinear correction which is treated as an error term, assumed to be “small” relative to W in an appropriate sense.

Suppose $\mathbb{E}[Z_\ell] = 0$ and $\sum_{\ell=1}^n \mathbb{E}[Z_\ell Z_\ell^\top] = \text{Id}$. Let $\Upsilon_n = \sum_{\ell=1}^n \mathbb{E}[\|Z_\ell\|^3]$. Then, for $Y \sim \mathcal{N}(0, \text{Id})$, we have

$$\begin{aligned} \text{d}_C(T, Y) &\leq 259d^{1/2}\Upsilon_n + 2\mathbb{E}[\|W\|\|D\|] \\ &\quad + 2 \sum_{\ell=1}^n \mathbb{E}[\|Z_\ell\| \|D - D^{(\ell)}\|], \quad (16) \end{aligned}$$

where $D^{(\ell)} = D(X_1, \dots, X_{\ell-1}, X'_\ell, X_{\ell+1}, \dots, X_n)$ and X'_ℓ is an independent copy of X_ℓ . This result follows from (Shao and Zhang, 2022, Theorem 2.1). The bound (16) extends to the case $\sum_{\ell=1}^n \mathbb{E}[Z_\ell Z_\ell^\top] = \Sigma \succ 0$, as shown in (Shao and Zhang, 2022, Corollary 2.3). To apply this result, we take $X_i = \xi_i$, $Z_\ell = h(X_\ell)$, and let ξ'_i be an i.i.d. copy of ξ_i . The key step is to bound $\mathbb{E}^{1/2}[\|D(\xi_1, \dots, \xi_{n-1})\|^2]$ and $\mathbb{E}^{1/2}[\|D - D'_i\|^2]$. Detailed proof is given in Appendix B. \square

2.3 Gaussian approximation in the bootstrap world

In the main result of this section, we study the Gaussian approximation for $\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n)$ w.r.t. \mathbb{P}^b , where the

target distribution is an appropriately chosen normal law. Although this result is similar in its nature to Theorem 2, it requires to address additional challenges that arise in the "bootstrap world". Our first steps are the same as in (11) and (6):

$$\begin{aligned} \theta_k^b - \theta_k &= (I - \alpha_k G)(\theta_{k-1}^b - \theta_{k-1}) \\ &\quad - \alpha_k (H(\theta_{k-1}^b) + g(\theta_{k-1}^b, \xi_k) \\ &\quad \quad - H(\theta_{k-1}) - g(\theta_{k-1}, \xi_k)) \\ &\quad - \alpha_k (w_k - 1)(G(\theta_{k-1}^b - \theta^*) + \eta(\xi_k) \\ &\quad \quad + g(\theta_{k-1}^b, \xi_k) + H(\theta_{k-1}^b)). \end{aligned}$$

Taking average of the above identity and rearranging the terms, we obtain a counterpart of (13):

$$\begin{aligned} \sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) &= W^b + D^b, \\ W^b &= -\frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} (w_i - 1) Q_i \eta(\xi_i), \\ D^b &= \sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) - W^b. \end{aligned} \quad (17)$$

Here W^b is a weighted sum of i.i.d. random variables Ξ^{n-1} , such that $\mathbb{E}^b[W^b] = 0$, $\mathbb{E}^b[W^b \{W^b\}^\top] = \Sigma_n^b$, where

$$\Sigma_n^b = n^{-1} \sum_{i=1}^{n-1} Q_i \eta(\xi_i) \eta(\xi_i)^\top Q_i^\top, \quad (18)$$

and D^b is a non-linear statistic of Ξ^{n-1} . The main difficulty arises when analyzing the conditional distribution of $\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n)$ given the data Ξ^{n-1} . The approach of (Shao and Zhang, 2022) requires controlling the second moments of D^b and $D^b - \{D^b\}^{(i)}$ with respect to the bootstrap measure \mathbb{P}^b , on a high-probability event under the original measure \mathbb{P} . At the same time, the martingale structure of the summands in D^b is lost unless we condition on an extended filtration.

$$\tilde{\mathcal{F}}_i = \sigma(w_1, \dots, w_i, \xi_1, \dots, \xi_i), 1 \leq i \leq n-1. \quad (19)$$

Therefore, it is not clear whether the approach of (Shao and Zhang, 2022), discussed in Section 2.2, can be applied directly. Instead, we rely on a linearization method that exploits high-order moment bounds for the remainder term D^b ; see Proposition 1 in Appendix A. This necessity motivates the strong bounded-noise assumption in A2. We now state the main result of this section:

Theorem 3. *Assume A1 - A6. Then with \mathbb{P} -probability at least $1 - 2/n$, it holds*

$$\begin{aligned} \sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\sqrt{n}\{\Sigma_n^b\}^{-\frac{1}{2}}(\bar{\theta}_n^b - \bar{\theta}_n) \in B) - \mathbb{P}^b(Y^b \in B)| \\ \leq \frac{M_{3,1}^b}{n^{1/2}} + \frac{M_{3,2}^b \log n}{n^{\gamma-1/2}} + \frac{M_{3,3}^b \log^{3/2} n}{n^{\gamma/2}}, \end{aligned}$$

where $Y^b \sim \mathcal{N}(0, I)$ under \mathbb{P}^b and $\{M_{3,i}^b\}_{i=1}^3$ are defined in Appendix C, equation (37).

Proof sketch of Theorem 3. We apply the bound

$$\begin{aligned} \sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\{\Sigma_n^b\}^{-\frac{1}{2}}(W^b + D^b) \in B) - \mathbb{P}^b(Y^b \in B)| \\ \leq \sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\{\Sigma_n^b\}^{-\frac{1}{2}}W^b \in B) - \mathbb{P}^b(Y^b \in B)| \\ \quad + 2c_d (\mathbb{E}^b[\|\{\Sigma_n^b\}^{-\frac{1}{2}}D^b\|^p])^{\frac{1}{1+p}}, \end{aligned} \quad (20)$$

where $c_d \leq 4d^{1/4}$ is the isoperimetric constant of the class of convex sets, see e.g. (Bentkus, 2003). The proof of (20) is provided in Proposition 1 in Appendix C. We first control $\mathbb{E}[\|D^b\|^p]$ using Burkholder's inequality, where \mathbb{E} denotes expectation with respect to the product measure $\mathbb{P}_\xi^{\otimes n} \otimes \mathbb{P}_w^{\otimes n}$. Applying Markov's inequality then yields \mathbb{P} -high-probability bounds for the behavior of $\mathbb{E}^b[\|D^b\|^p]$. This step requires bounds on $\mathbb{E}^{1/p}[\|\theta_k - \theta^*\|^p]$ and $\mathbb{E}^{1/p}[\|\theta_k^b - \theta^*\|^p]$, $k \in \{1, \dots, n-1\}$, where $p = c \log n$ for some absolute constant c , with polynomial dependence on p .

To control the second term on the right-hand side of (20), we note that Σ_n^b concentrates around Σ_n by the matrix Bernstein inequality (see Lemma 2). Hence, there exists a set Ω_1 with $\mathbb{P}(\Omega_1) \geq 1 - 1/n$ such that $\lambda_{\min}(\Sigma_n^b) > 0$ on Ω_1 . On this event, one can apply Berry-Esseen type bounds for non-i.i.d. sums of random vectors. The full proof is presented in Appendix C. \square

Remark 6. *The constants $M_{3,1}^b, M_{3,2}^b, M_{3,3}^b$ in Theorem 3 depend on the problem dimension d , as well as $\|\theta_0 - \theta^*\|$. To make the dependence on d explicit, we assume the natural scaling $\|\theta_0 - \theta^*\| \lesssim \sqrt{d}$. We also assume that C_ξ from A2-(iii) is dimension-free. Under this assumption, Theorem 3 implies*

$$\begin{aligned} \sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\sqrt{n}\{\Sigma_n^b\}^{-\frac{1}{2}}(\bar{\theta}_n^b - \bar{\theta}_n) \in B) - \mathbb{P}^b(Y^b \in B)| \\ \lesssim \frac{d^2}{\sqrt{n}} + \frac{d^{5/4} \log n}{n^{\gamma-1/2}} + \frac{d^{3/4} (\log n)^{3/2}}{n^{\gamma/2}}. \end{aligned}$$

Here, the notation \lesssim indicates that the inequality holds up to constants independent of n and d .

2.4 Rate of convergence in the Polyak-Juditsky central limit theorem

We finally address the change from Σ_n to Σ_∞ and derive convergence rates in the Polyak-Juditsky result (4). Our argument builds on Theorem 2 together with the following lemma.

Lemma 1. *Assume A1 and A8. Let $Y \sim \mathcal{N}(0, I_d)$. Then the convex distance between the distributions of $\Sigma_n^{1/2}Y$ and $\Sigma_\infty^{1/2}Y$ satisfies*

$$d_{\mathcal{C}}(\Sigma_n^{1/2}Y, \Sigma_\infty^{1/2}Y) \leq C_\infty n^{\gamma-1},$$

where the constant C_∞ is defined in (49).

Combining Theorem 2 with Lemma 1 and applying the triangle inequality yields the following result on the closeness of $\sqrt{n}(\bar{\theta}_n - \theta^*)$ to $\mathcal{N}(0, \Sigma_\infty)$.

Theorem 4. *Assume A1, A3, A7(4), A8. Then, with $Y \sim \mathcal{N}(0, \mathbf{I}_d)$ it holds that*

$$d_{\mathcal{C}}(\sqrt{n}(\bar{\theta}_n - \theta^*), \Sigma_\infty^{1/2} Y) \leq \frac{C_4}{\sqrt{n}} + \frac{C_5 + C_6}{n^{\gamma-1/2}} + \frac{C_\infty}{n^{1-\gamma}}, \quad (21)$$

where C_4, C_5 and C_6 are given in Theorem 2.

Discussion. Theorem 2 shows that the normal approximation with $\mathcal{N}(0, \Sigma_n)$ improves as the step sizes α_k become less aggressive, i.e., as $\gamma \rightarrow 1$. At the same time, Theorem 4 highlights a trade-off: the rate at which Σ_n converges to Σ_∞ also influences the overall approximation quality. Optimizing the bound in (21) with respect to γ yields the optimal choice $\gamma = 3/4$, which leads to the approximation rate

$$\begin{aligned} d_{\mathcal{C}}(\sqrt{n}(\bar{\theta}_n - \theta^*), \Sigma_\infty^{1/2} Y) \\ \leq \frac{C'_1}{n^{1/4}} + \frac{C'_2}{\sqrt{n}} (\|\theta_0 - \theta^*\| + \|\theta_0 - \theta^*\|^2), \end{aligned}$$

where C'_1 and C'_2 are instance-dependent constants (independent of $\|\theta_0 - \theta^*\|$) that can be derived from Theorem 4. This result enables a non-asymptotic assessment of methods for constructing confidence intervals based on direct estimation of Σ_∞ , such as those studied in (Chen et al., 2020; Zhu et al., 2023).

Lower bounds We provide a lower bound showing that the bound in Theorem 4 is tight in certain regimes of step size decay power $\gamma \in (1/2, 1)$. To this end, we consider the minimization problem (1) with $f(\theta) = \theta^2/2$ and $\theta_0 = 0$. In this case, $\theta^* = 0$. We use an additive noise model, where the stochastic gradient oracles $F(\theta, \xi)$ are given by $F(\theta, \xi) = \theta + \xi$, $\xi \sim \mathcal{N}(0, 1)$. Unrolling (2), we get

$$\sqrt{n}\bar{\theta}_n = -n^{-1/2} \sum_{j=1}^{n-1} Q_j \xi_j \quad (22)$$

with $Q_j = \alpha_j \sum_{k=j}^{n-1} \prod_{\ell=j+1}^k (1 - \alpha_\ell)$, showing that $\sqrt{n}(\bar{\theta}_n - \theta^*) \sim \mathcal{N}(0, \sigma_{n,\gamma}^2)$ with $\sigma_{n,\gamma}^2 = n^{-1} \sum_{j=1}^{n-1} Q_j^2$. Lemma 1 (see also (48) in the Appendix), we have $G = 1, \Sigma_\infty = 1$, and $\sigma_{n,\gamma}^2 \rightarrow 1$ as $n \rightarrow \infty$. Moreover, the following lower bound holds:

Theorem 5. *Consider the sequence $\{\theta_k\}_{k \in \mathbb{N}}$ defined by the recurrence (22) with $\alpha_j = c_0/(1+j)^\gamma$. Then it holds, for the number of observations n sufficiently large, that*

$$|\sigma_{n,\gamma}^2 - 1| > \frac{C_1(\gamma, c_0)}{n^{1-\gamma}}, \quad (23)$$

where the constant $C_1(\gamma, c_0)$ depends only upon c_0 and γ . Moreover, for n large enough

$$d_{\mathcal{C}}(\sqrt{n}(\bar{\theta}_n - \theta^*), \mathcal{N}(0, 1)) > \frac{C_2(\gamma, c_0)}{n^{1-\gamma}}. \quad (24)$$

Discussion. Proof of Theorem 5 is provided in Appendix E, along with simple numerical simulations indicating the tightness of the lower bound (23). Note that the bound (24) shows that the distribution of $\sqrt{n}(\bar{\theta}_n - \theta^*)$ cannot be approximated by $\mathcal{N}(0, \Sigma_\infty)$ at a rate faster than $1/n^{1-\gamma}$. Moreover, it shows that the rate of normal approximation in Theorem 4 cannot be improved when $\gamma \in [3/4, 1)$. This fact is extremely important when considering the bootstrap validity result in Theorem 1 and the normal approximation in Theorem 2. Indeed, both results suggest that normal approximation rates of order up to $1/\sqrt{n}$ can be achieved as $\gamma \rightarrow 1$, but they require using a different covariance matrix Σ_n , corresponding to the linearized recurrence in (14). At the same time, in the regime $\gamma \rightarrow 1$, the approximation by $\mathcal{N}(0, \Sigma_\infty)$ can be too slow. It is an interesting and, to the best of our knowledge, open question to provide lower bounds analogous to Theorem 5 which show the tightness of other summands in Theorem 4 in the regime $1/2 < \gamma < 3/4$.

3 CONCLUSION

In our paper, we performed the fully non-asymptotic analysis of the multiplier bootstrap procedure for SGD applied to strongly convex minimization problems. We showed that the algorithm can achieve approximation rates in convex distances of order up to $1/\sqrt{n}$. We highlight the fact that the validity of the multiplier bootstrap procedure does not require one to consider Berry-Esseen bounds with the asymptotic covariance matrix Σ_∞ , which is in sharp contrast to the methods that require direct estimation of Σ_∞ .

Acknowledgment

This work is an output of a research project HSE-BR-2025-019 implemented as part of the Basic Research Program at HSE University. This research was supported in part through computational resources of HPC facilities at HSE University (Kostenetskiy et al., 2021).

References

- Bhavya Agrawalla, Krishnakumar Balasubramanian, and Promit Ghosal. High-dimensional Central Limit Theorems for Linear Functionals of Online Least-Squares SGD. *arXiv preprint arXiv:2302.09727*, 2023.
- Andreas Anastasiou, Krishnakumar Balasubramanian, and Murat A. Erdogdu. Normal Approximation for Stochastic Gradient Descent via Non-Asymptotic Rates of Martingale CLT. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 115–137. PMLR, 25–28 Jun 2019. URL <https://proceedings.mlr.press/v99/anastasiou19a.html>.
- G. D. Anderson and S.-L. Qiu. A monotoneity property of the gamma function. *Proceedings of the American Mathematical Society*, 125(11):3355–3362, 1997. ISSN 00029939, 10886826. URL <http://www.jstor.org/stable/2162408>.
- Keith Ball. The reverse isoperimetric problem for Gaussian measure. *Discrete & Computational Geometry*, 10(4):411–420, October 1993. ISSN 1432-0444. doi: 10.1007/BF02573986. URL <https://doi.org/10.1007/BF02573986>.
- S. Barsov and V. Ulyanov. Estimates for the closeness of Gaussian measures. *Dokl. Akad. Nauk SSSR*, 291(2):273–277, 1986.
- V. Bentkus. On the dependence of the Berry–Esseen bound on dimension. *Journal of Statistical Planning and Inference*, 113(2):385–402, 2003. ISSN 0378-3758. doi: [https://doi.org/10.1016/S0378-3758\(02\)00094-0](https://doi.org/10.1016/S0378-3758(02)00094-0). URL <https://www.sciencedirect.com/science/article/pii/S0378375802000940>.
- A. Benveniste, M. Métivier, and P. Priouret. *Adaptive algorithms and stochastic approximations*, volume 22. Springer Science & Business Media, 2012.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Leheng Cai, Qirui Hu, Juntao Sun, and Shuyuan Wu. Time-uniform and Asymptotic Confidence Sequence of Quantile under Local Differential Privacy. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- Xiangyu Chang, Xi Chen, Zehua Lai, He Li, Zhihong Liu, and Yichen Zhang. Online Statistical Inference for Contextual Bandits via Stochastic Gradient Descent. *Journal of the American Statistical Association*, 0(ja):1–24, 2026. doi: 10.1080/01621459.2026.2621503. URL <https://doi.org/10.1080/01621459.2026.2621503>.
- Haoyu Chen, Wenbin Lu, and Rui Song. Statistical inference for online decision making via stochastic gradient descent. *Journal of the American Statistical Association*, 116(534):708–719, 2021.
- Likai Chen, Georg Keilbar, and Wei Biao Wu. Recursive quantile estimation: Non-asymptotic confidence bounds. *Journal of Machine Learning Research*, 24(91):1–25, 2023.
- Likai Chen, Georg Keilbar, and Wei Biao Wu. Smoothed SGD for quantiles: Bahadur representation and Gaussian approximation. *arXiv preprint arXiv:2505.13299*, 2025.
- Xi Chen, Jason D. Lee, Xin T. Tong, and Yichen Zhang. Statistical inference for model parameters in stochastic gradient descent. *The Annals of Statistics*, 48(1):251 – 273, 2020. doi: 10.1214/18-AOS1801. URL <https://doi.org/10.1214/18-AOS1801>.
- Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Ann. Statist.*, 41(6):2786–2819, 2013. ISSN 0090-5364. doi: 10.1214/13-AOS1161. URL <https://doi.org/10.1214/13-AOS1161>.
- 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.
- Luc Devroye, Abbas Mehrabian, and Tommy Reddad. The total variation distance between high-dimensional gaussians with the same mean. *arXiv preprint arXiv:1810.08693*, 2018.
- 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>.
- Alain Durmus, Eric Moulines, Alexey Naumov, Sergey Samsonov, and Hoi-To Wai. On the Stability of Random Matrix Product with Markovian Noise: Application to Linear Stochastic Approximation and TD learning. In Mikhail Belkin and Samory Kpotufe, editors, *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pages 1711–1752. PMLR, 15–19 Aug 2021. URL <https://proceedings.mlr.press/v134/durmus21a.html>.
- Bradley Efron. Bootstrap methods: another look at the jackknife. In *Breakthroughs in statistics: Methodology and distribution*, pages 569–593. Springer, 1992.
- Yixin Fang, Jinfeng Xu, and Lei Yang. Online bootstrap confidence intervals for the stochastic gradient

- descent estimator. *Journal of Machine Learning Research*, 19(78):1–21, 2018. URL <http://jmlr.org/papers/v19/17-370.html>.
- M. V. Fedoryuk. *Metod perevala*. Izdat. “Nauka”, Moscow, 1977.
- James M. Flegal and Galin L. Jones. Batch means and spectral variance estimators in Markov chain Monte Carlo. *Ann. Statist.*, 38(2):1034–1070, 04 2010. doi: 10.1214/09-AOS735. URL <https://doi.org/10.1214/09-AOS735>.
- Wolfgang Gabcke. Neue herleitung und explizite restabschätzung der riemann-siegel-formel. 2015.
- Friedrich Götze, Alexey Naumov, Vladimir Spokoiny, and Vladimir Ulyanov. Large ball probabilities, Gaussian comparison and anti-concentration. *Bernoulli*, 25(4A):2538–2563, 2019. ISSN 1350-7265. doi: 10.3150/18-BEJ1062. URL <https://doi.org/10.3150/18-BEJ1062>.
- Nicholas JA Harvey, Christopher Liaw, Yaniv Plan, and Sikander Randhawa. Tight analyses for non-smooth stochastic gradient descent. In *Conference on Learning Theory*, pages 1579–1613. PMLR, 2019.
- Elad Hazan and Satyen Kale. Beyond the regret minimization barrier: optimal algorithms for stochastic strongly-convex optimization. *The Journal of Machine Learning Research*, 15(1):2489–2512, 2014.
- Chi Jin, Praneeth Netrapalli, Rong Ge, Sham M Kakade, and Michael I Jordan. A short note on concentration inequalities for random vectors with subgaussian norm. *arXiv preprint arXiv:1902.03736*, 2019.
- Denis Kojevnikov and Kyungchul Song. A Berry–Esseen bound for vector-valued martingales. *Statistics & Probability Letters*, 186:109448, 2022.
- PS Kostenetskiy, RA Chulkevich, and VI Kozyrev. Hpc resources of the higher school of economics. In *Journal of Physics: Conference Series*, volume 1740, page 012050. IOP Publishing, 2021.
- Harold Kushner and G George Yin. *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media, 2003.
- Guanghui Lan. *First-order and Stochastic Optimization Methods for Machine Learning*. 01 2020. ISBN 978-3-030-39567-4. doi: 10.1007/978-3-030-39568-1.
- Chris Junchi Li, Wenlong Mou, Martin Wainwright, and Michael Jordan. ROOT-SGD: Sharp Nonasymptotics and Asymptotic Efficiency in a Single Algorithm. In Po-Ling Loh and Maxim Raginsky, editors, *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 909–981. PMLR, 02–05 Jul 2022a. URL <https://proceedings.mlr.press/v178/li22a.html>.
- Jiaqi Li, Zhipeng Lou, Johannes Schmidt-Hieber, and Wei Biao Wu. Statistical Guarantees for High-Dimensional Stochastic Gradient Descent. *arXiv preprint arXiv:2510.12013*, 2025.
- Xiang Li, Jiadong Liang, Xiangyu Chang, and Zhihua Zhang. Statistical estimation and online inference via local sgd. In Po-Ling Loh and Maxim Raginsky, editors, *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 1613–1661. PMLR, 02–05 Jul 2022b. URL <https://proceedings.mlr.press/v178/li22b.html>.
- Liam Madden, Emiliano Dall’Anese, and Stephen Becker. High probability convergence bounds for non-convex stochastic gradient descent with sub-weibull noise. *Journal of Machine Learning Research*, 25(241):1–36, 2024.
- Marc S. Meketon and Bruce W. Schmeiser. Overlapping Batch Means: Something for Nothing? In *Proceedings of the 1984 Winter Simulation Conference*, pages 227–230, Washington, DC, USA, 1984. IEEE Press. URL https://informs-sim.org/wsc84papers/1984_0041.pdf.
- Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. *Advances in neural information processing systems*, 24:451–459, 2011.
- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
- Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Applied Optimization. Springer, 2004. ISBN 9781402075537. URL <http://books.google.fr/books?id=VyYLeM-13CgC>.
- Frank W. J. Olver. *Asymptotics and special functions*. AKP Classics. A K Peters, Ltd., Wellesley, MA, 1997. ISBN 1-56881-069-5. Reprint of the 1974 original [Academic Press, New York; MR0435697 (55 #8655)].
- Adam Osekowski. *Sharp martingale and semimartingale inequalities*, volume 72. Springer Science & Business Media, 2012.
- 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.

- Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. In *Proceedings of the 29th International Conference on Machine Learning*, pages 1571–1578, 2012.
- Abhishek Roy and Krishnakumar Balasubramanian. Online covariance estimation for stochastic gradient descent under Markovian sampling. *arXiv preprint arXiv:2308.01481*, 2023.
- David Ruppert. Efficient estimations from a slowly convergent Robbins-Monro process. Technical report, Cornell University Operations Research and Industrial Engineering, 1988.
- Abdurakhmon Sadiev, Marina Danilova, Eduard Gorbunov, Samuel Horváth, Gauthier Gidel, Pavel Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability bounds for stochastic optimization and variational inequalities: the case of unbounded variance. In *International Conference on Machine Learning*, pages 29563–29648. PMLR, 2023.
- Sergey Samsonov, Eric Moulines, Qi-Man Shao, Zhuo-Song Zhang, and Alexey Naumov. Gaussian Approximation and Multiplier Bootstrap for Polyak-Ruppert Averaged Linear Stochastic Approximation with Applications to TD Learning. In *Advances in Neural Information Processing Systems*, volume 37, pages 12408–12460. Curran Associates, Inc., 2024.
- Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162:83–112, 2017.
- Jun Shao. *Mathematical statistics*. Springer Science & Business Media, 2003.
- Jun Shao and Dongsheng Tu. *The Jackknife and Bootstrap*. Springer Series in Statistics. Springer New York, NY, 1 edition, 1995. ISBN 978-0-387-94515-6. doi: 10.1007/978-1-4612-0795-5. URL <https://doi.org/10.1007/978-1-4612-0795-5>. Springer Book Archive, © Springer Science+Business Media New York 1995.
- 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.
- Marina Sheshukova, Denis Belomestny, Alain Durmus, Eric Moulines, Alexey Naumov, and Sergey Samsonov. Nonasymptotic analysis of stochastic gradient descent with the richardson-romberg extrapolation. *International Conference on Learning Representations (ICLR)*, 2025.
- R Srikant. Rates of Convergence in the Central Limit Theorem for Markov Chains, with an Application to TD Learning. *Mathematics of Operations Research*, 2025. doi: 10.1287/moor.2024.0444.
- Joel A Tropp et al. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.
- Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018.
- Weichen Wu, Gen Li, Yuting Wei, and Alessandro Rinaldo. Statistical Inference for Temporal Difference Learning with Linear Function Approximation. *arXiv preprint arXiv:2410.16106*, 2024.
- Yanjie Zhong, Todd Kuffner, and Soumendra Lahiri. Online Bootstrap Inference with Nonconvex Stochastic Gradient Descent Estimator. *arXiv preprint arXiv:2306.02205*, 2023.
- Wanrong Zhu, Xi Chen, and Wei Biao Wu. Online Covariance Matrix Estimation in Stochastic Gradient Descent. *Journal of the American Statistical Association*, 118(541):393–404, 2023. doi: 10.1080/01621459.2021.1933498. URL <https://doi.org/10.1080/01621459.2021.1933498>.
- Wanrong Zhu, Zhipeng Lou, Ziyang Wei, and Wei Biao Wu. High confidence level inference is almost free using parallel stochastic optimization. *arXiv preprint arXiv:2401.09346*, 2024.

Checklist

1. For all models and algorithms presented, check if you include:
 - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. Answer: Yes. All main results of this submission are supported with rigorous assumptions and proofs. Assumptions are presented in Section 2 and Section 2.2, as well as the statements of the main theorems. Proofs are provided in Appendix.
 - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. Answer: Yes. The paper is devoted to the analysis of SGD algorithm and rates of normal approximation for the SGD algorithm. All theoretical results are mathematically justified. At the same time, the paper does not provide a new algorithm.
 - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. Answer: Yes. Code to reproduce experiments will be presented at the anonymous github.

2. For any theoretical claim, check if you include:
 - (a) Statements of the full set of assumptions of all theoretical results. Answer: Yes. All theoretical results are stated with explicit pointers to the underlying assumptions (see Section 2 and Section 2.2)
 - (b) Complete proofs of all theoretical results. Answer: Yes. Detailed proof of each theorem stated in the main text is provided in the appendix with clear references to the relevant sections.
 - (c) Clear explanations of any assumptions. Answer: Yes. A detailed discussion of each of the assumptions is provided in Section 2.
 - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. Answer: Not applicable.
3. For all figures and tables that present empirical results, check if you include:
 - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). Answer: Yes. All code is open source, link to an anonymous GitHub repository is included.
 - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). Answer: Yes. Numerical results are stated with a complete description of the environments that are used, as well as the precise sets of hyperparameters that we used. The code (in Python) is provided as supplementary with the paper, making it easy for one to reproduce our numerical experiments.
 - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). Answer: Yes.
 - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). Answer: Yes. All necessary information to reproduce experiments is provided in Appendix.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
 - (a) Citations of the creator If your work uses existing assets. Answer: Not applicable.
 - (b) The license information of the assets, if applicable. Answer: Not applicable.
 - (c) New assets either in the supplemental material or as a URL, if applicable. Answer: Not applicable.
 - (d) Information about consent from data providers/curators. Answer: Not Applicable.
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
 - (a) The full text of instructions given to participants and screenshots. Answer: Not applicable.
 - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. Answer: Not applicable.
 - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. Answer: Not applicable.

Contents

| | | |
|----------|--|-----------|
| 1 | INTRODUCTION | 1 |
| 2 | MAIN RESULTS | 3 |
| 2.1 | Non-asymptotic multiplier bootstrap validity | 5 |
| 2.2 | Gaussian approximation in the real world | 7 |
| 2.3 | Gaussian approximation in the bootstrap world | 7 |
| 2.4 | Rate of convergence in the Polyak–Juditsky central limit theorem | 8 |
| 3 | CONCLUSION | 9 |
| A | Proof of Theorem 1 | 15 |
| A.1 | Proof of Theorem 1 | 15 |
| A.2 | Matrix Bernstein inequality for Σ_n^b and Gaussian comparison | 16 |
| A.3 | Example of distribution satisfying A4 | 17 |
| B | Proof of Theorem 2 | 17 |
| B.1 | Bounds for D_n | 19 |
| B.2 | Bounds for $D_n - D_n^{(i)}$ | 20 |
| B.3 | Bounds for $\theta_k^{(i)} - \theta_k$ | 22 |
| C | Proof of Theorem 3 | 23 |
| C.1 | From non-linear to linear statistics | 24 |
| C.2 | Bounds for D^b | 25 |
| D | Proof of Theorem 4 | 27 |
| D.1 | Proof of Lemma 1 | 27 |
| E | Lower bounds | 31 |
| E.1 | Numerical demonstration | 32 |
| F | Results for the last iterate | 32 |
| F.1 | Last iterate bound | 32 |
| F.2 | High probability bounds on the last iterate | 41 |
| G | Technical bounds | 43 |
| G.1 | Proof of Lemma 3 | 45 |
| G.2 | Properties of sub-Gaussian random vectors | 46 |
| G.3 | Gaussian comparison lemma | 46 |
| H | Numerical Experiments | 47 |

where $M_{1,1}^b, M_{2,1}^b$ and $C_{Q,\xi}$ are defined in (43), (44), (26) respectively. The first set ensures that Σ_n^b concentrates around Σ_n , while the second set guarantees that the remainder term D^b is small.

Therefore applying Theorem 3 (see the proof in C) and Corollary 1 we get that on the set $\Omega_0 \cap \Omega_1$ with $\mathbb{P}(\Omega_0 \cap \Omega_1) \geq 1 - 2/n$, it holds

$$\sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) \in B) - \mathbb{P}^b(\{\Sigma_n^b\}^{1/2} Y^b \in B)| \leq \frac{M_{3,1}^b}{\sqrt{n}} + \frac{M_{3,2}^b \log n}{n^{\gamma-1/2}} + \frac{M_{3,3}^b \log^{3/2} n}{n^{\gamma/2}},$$

where $\{M_{3,i}^b\}_{i=1}^3$ are defined in equation (37).

To complete the proof, we need to bound the convex distance between the two Gaussians. By Lemma 3, we have $\|\Sigma_n^{-1/2}\| \leq C_\Sigma$. Hence, due to Lemma 2, on the set Ω_1 with $\mathbb{P}(\Omega_1) \geq 1 - 1/n$ we have

$$\text{Tr}\{(\Sigma_n^{-1/2} \Sigma_n^b \Sigma_n^{-1/2} - I_p)^2\} \leq d \|(\Sigma_n^{-1/2} \Sigma_n^b \Sigma_n^{-1/2} - I_p)^2\|^2 \leq d C_\Sigma^2 \|\Sigma_n^b - \Sigma_n\|^2 \leq \delta^2.$$

where we have set

$$\delta = \frac{10C_{Q,\xi} C_\Sigma^2 \sqrt{d \log(2dn)}}{3\sqrt{n}}$$

Applying Lemma 20 with probability at least $1 - 1/n$ we get

$$T_3 \leq \frac{5C_{Q,\xi} C_\Sigma^2 \sqrt{d \log(2dn)}}{\sqrt{n}}.$$

Collecting previous bounds and using triangle inequality we get that on the set $\Omega_0 \cap \Omega_1$ with $\mathbb{P}(\Omega_0 \cap \Omega_1) \geq 1 - 2/n$, it holds:

$$\sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\sqrt{n}(\bar{\theta}_n^b - \bar{\theta}_n) \in B) - \mathbb{P}(\sqrt{n}(\bar{\theta}_n - \theta^*) \in B)| \leq \frac{C_1 \sqrt{\log n}}{n^{1/2}} + \frac{C_2 \log n}{n^{\gamma-1/2}} + \frac{C_3 \log^{3/2} n}{n^{\gamma/2}},$$

where

$$C_1 = C_4 + M_{3,1}^b + 5C_{Q,\xi} C_\Sigma^2 \sqrt{d \log(2d)}, \quad C_2 = C_5 + M_{3,2}^b, \quad C_3 = C_6 + M_{3,3}^b. \quad (25)$$

A.2 Matrix Bernstein inequality for Σ_n^b and Gaussian comparison

Lemma 2. *Under assumptions A1, A2, A5, A6, there is a set $\Omega_1 \in \mathcal{F}_{n-1}$, such that $\mathbb{P}(\Omega_1) \geq 1 - 1/n$ and on Ω_1 it holds that*

$$\|\Sigma_n^b - \Sigma_n\| \leq \frac{10C_{Q,\xi} \sqrt{\log(2dn)}}{3\sqrt{n}}$$

where the constant $C_{Q,\xi}$ is given by

$$C_{Q,\xi} := C_Q^2 (C_{1,\xi}^2 + \lambda_{\max}(\Sigma_\xi)), \quad (26)$$

and $C_{1,\xi}, C_Q$ are defined in A2 and Lemma 3, respectively.

Proof. Note that

$$\Sigma_n^b - \Sigma_n = \frac{1}{n} \sum_{i=1}^{n-1} Q_i(\eta(\xi_i)\eta(\xi_i)^\top - \Sigma_\xi) Q_i^\top.$$

For simplicity we denote $A_i = Q_i(\eta(\xi_i)\eta(\xi_i)^\top - \Sigma_\xi) Q_i^\top$. Note that for any $i \in \{1, \dots, n-1\}$ it holds that

$$\mathbb{E}[A_i] = 0, \quad \|A_i\| \leq C_{Q,\xi}, \quad \left\| \sum_{i=1}^{n-1} \mathbb{E}[A_i A_i^\top] \right\| \leq n C_{Q,\xi}^2.$$

Then, using matrix Bernstein inequality (Tropp et al., 2015, Chapter 6), we obtain

$$\mathbb{P}\left(\frac{1}{n} \left\| \sum_{i=1}^{n-1} A_i \right\| \geq t\right) \leq 2d \exp\left\{\frac{-t^2 n^2 / 2}{nC_{Q,\xi}^2 + nC_{Q,\xi} t / 3}\right\}.$$

Taking $t_\delta = \frac{4C_{Q,\xi} \log(2d/\delta)}{3n} + \frac{2C_{Q,\xi} \sqrt{\log(2d/\delta)}}{\sqrt{n}}$, we obtain that with probability at least $1 - \delta$, it holds

$$\frac{1}{n} \left\| \sum_{i=1}^{n-1} A_i \right\| \leq t_\delta .$$

Setting $\delta = 1/n$ and applying A6 completes the proof. \square

Corollary 1. *Under assumptions A1, A2, A5, A6, on Ω_1 it holds that*

$$\lambda_{\min}(\Sigma_n^b) \geq \frac{1}{2C_\Sigma^2} .$$

Proof. Using eigenvalue stability (Lidski's) inequality, we obtain

$$\lambda_{\min}(\Sigma_n^b) \geq \lambda_{\min}(\Sigma_n) - \|\Sigma_n - \Sigma_n^b\| .$$

Note that on Ω_1 , we have

$$\|\Sigma_n - \Sigma_n^b\| \leq \frac{10C_{Q,\xi} \sqrt{\log(2dn)}}{3\sqrt{n}} \leq \frac{1}{2C_\Sigma^2} ,$$

where in the last inequality we use A6. \square

A.3 Example of distribution satisfying A4

To construct examples of distributions satisfying the above assumption, one can use the beta distribution, which is defined on $[0, 1]$, and then shift and scale it. Set $W = a + bX$ where $X \sim \text{Beta}(\alpha, \beta)$ and $a, b > 0$. We have $\mathbb{E}[X] = \frac{\alpha}{\alpha + \beta}$, $\text{Var}(X) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$ and $a \leq W \leq a + b$ a.s. By solving (for a and b) the equations $\mathbb{E}[W] = a + b\mathbb{E}[X] = 1$ and $\text{Var}(W) = b^2 \text{Var}(X) = 1$, we derive $b = 1/\sqrt{\text{Var}(X)}$ and $a = 1 - \mathbb{E}[X]/\sqrt{\text{Var}(X)}$. Note that $a > 0$ provided $\alpha + \beta + 1 < \beta/\alpha$.

B Proof of Theorem 2

We first provide details of the expansion (13). Recall that the error of SGD approximation may be rewritten as follows

$$\theta_k - \theta^* = (\mathbf{I} - \alpha_k G)(\theta_{k-1} - \theta^*) - \alpha_k (H(\theta_{k-1}) + \eta(\xi_k) + g(\theta_{k-1}, \xi_k)) . \quad (27)$$

Iteratively spinning this expression out we get

$$\theta_k - \theta^* = \prod_{j=1}^k (\mathbf{I} - \alpha_j G)(\theta_0 - \theta^*) - \sum_{j=1}^k \alpha_j \prod_{i=j+1}^k (\mathbf{I} - \alpha_i G)(H(\theta_{j-1}) + \eta(\xi_j) + g(\theta_{j-1}, \xi_j)) .$$

Taking average of (27) and changing the order of summation, we obtain

$$\sqrt{n}(\bar{\theta}_n - \theta^*) = \frac{1}{\sqrt{n}\alpha_0} Q_0(\theta_0 - \theta^*) - \frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i (H(\theta_{i-1}) + \eta(\xi_i) + g(\theta_{i-1}, \xi_i)) ,$$

where Q_i is defined in (12). Finally, we obtain

$$\begin{aligned} \sqrt{n}(\bar{\theta}_n - \theta^*) &= W + D, \\ D &= \frac{1}{\sqrt{n}\alpha_0} Q_0(\theta_0 - \theta^*) - \frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i g(\theta_{i-1}, \xi_i) - \frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i H(\theta_{i-1}) , \\ W &= -\frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i \eta(\xi_i) . \end{aligned} \quad (28)$$

Under suitable assumptions on the step size and the Hessian, we can show that the matrices Q_i , defined in (12), are uniformly bounded for all i . This fact will be used in the proof of Theorem 2. Moreover, we also establish that $\lambda_{\min}(Q_i)$ is bounded away from zero, which in turn implies a lower bound on $\lambda_{\min}(\Sigma_n)$.

Lemma 3. *Assume A1 and A8. Then for any $i \in \{0, \dots, n-1\}$ it holds that*

$$\lambda_{\max}(Q_i) \leq C_Q ,$$

where the constant C_Q is defined in (67) Moreover,

$$\lambda_{\min}(Q_i) \geq C_Q^{\min} , \text{ and } \|\Sigma_n^{-1/2}\| \leq C_\Sigma ,$$

where the matrix Σ_n is defined in (14), and C_Q^{\min}, C_Σ are defined in (69) and (70) respectively.

The version of this lemma with proof and with explicit constants is given in G.1.

Proof of Theorem 2. We normalize the both parts of (13) by $\Sigma_n^{1/2}$ and obtain

$$\sqrt{n}\Sigma_n^{-\frac{1}{2}}(\bar{\theta}_n - \theta^*) = \sum_{i=1}^{n-1} \underbrace{\frac{\Sigma_n^{-\frac{1}{2}}}{\sqrt{n}} Q_i \eta(\xi_i)}_{w_i} + D_{n,1} + D_{n,2} + D_{n,3} ,$$

where we have set

$$\begin{aligned} D_{n,1} &= \frac{\Sigma_n^{-\frac{1}{2}}}{\sqrt{n}\alpha_0} Q_0(\theta_0 - \theta^*) , \\ D_{n,2} &= -\frac{\Sigma_n^{-\frac{1}{2}}}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i H(\theta_{i-1}) , \\ D_{n,3} &= -\frac{\Sigma_n^{-\frac{1}{2}}}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i g(\theta_{i-1}, \xi_i) . \end{aligned}$$

Also, for any $1 \leq i \leq n-1$ we construct

$$\begin{aligned} D_{n,1}^{(i)} &= \frac{\Sigma_n^{-1/2}}{\sqrt{n}\alpha_0} Q_0(\theta_0^{(i)} - \theta^*) , \\ D_{n,2}^{(i)} &= -\frac{\Sigma_n^{-1/2}}{\sqrt{n}} \sum_{j=1}^{n-1} Q_j H(\theta_{j-1}^{(i)}) , \\ D_{n,3}^{(i)} &= -\frac{\Sigma_n^{-1/2}}{\sqrt{n}} \sum_{j=1}^{n-1} Q_j g(\theta_{j-1}^{(i)}, \tilde{\xi}_j^{(i)}) , \end{aligned}$$

where we set

$$\tilde{\xi}_j^{(i)} = \begin{cases} \xi_j , & \text{if } j \neq i \\ \xi'_j , & \text{if } j = i . \end{cases}$$

Define $D_n = D_{n,1} + D_{n,2} + D_{n,3}$, $D_n^{(i)} = D_{n,1}^{(i)} + D_{n,2}^{(i)} + D_{n,3}^{(i)}$, $W_n = \sum_{i=1}^{n-1} w_i$ and $\Upsilon_n = \sum_{i=1}^n \mathbb{E}[\|\omega_i\|^3]$ (we keep the same notations as in the unnormalized setting for simplicity). Let $Y \sim \mathcal{N}(0, I_d)$. Then, using (Shao and Zhang, 2022, Theorem 2.1), we have

$$d_C(\sqrt{n}\Sigma_n^{-1/2}(\bar{\theta}_n - \theta^*), Y) \leq 259d^{1/2}\Upsilon_n + 2 \underbrace{\mathbb{E}\{\|W_n\|\|D_n\|\}}_{R_1} + 2 \underbrace{\sum_{i=1}^{n-1} \mathbb{E}[\|\omega_i\|\|D_n - D_n^{(i)}\|]}_{R_2} .$$

Applying Hölder's inequality, we get

$$\begin{aligned} R_1 &\leq \mathbb{E}^{1/2}[\|W_n\|^2] \mathbb{E}^{1/2}[\|D_n\|^2], \\ R_2 &\leq \sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|\omega_i\|^2] \mathbb{E}^{1/2}[\|D_n - D_n^{(i)}\|^2]. \end{aligned}$$

Note that $\mathbb{E}^{1/2}[\|W_n\|^2] = \sqrt{d}$. Applying Lemma 3, we get $\mathbb{E}^{1/2}\|w_i\|^2 \leq \frac{1}{\sqrt{n}} C_\Sigma C_Q \sigma_2$ and

$$\Upsilon_n \leq \frac{1}{\sqrt{n}} (C_\Sigma C_Q \sigma_4)^3.$$

To complete the proof, it remains to bound $\mathbb{E}^{1/2}[\|D_n\|^2]$ and $\sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|D_n - D_n^{(i)}\|^2]$. The first term can be bounded using Lemma 4, while the second is controlled via Lemma 5. Combining these results, we obtain the following bound:

$$d_C(\sqrt{n} \Sigma_n^{-1/2}(\bar{\theta}_n - \theta^*), Y) \leq \frac{\sqrt{d} M_{3,1}}{\sqrt{n}} + \frac{M_{3,2}}{\sqrt{n}} (\|\theta_0 - \theta^*\| + \|\theta_0 - \theta^*\|^2 + \sigma_2 + \sigma_4^2) + M_{3,3} n^{1/2-\gamma} + M_{3,4} n^{-\gamma/2},$$

where

$$\begin{aligned} M_{3,1} &= 259(C_\Sigma C_Q \sigma_4)^3, \\ M_{3,2} &= 2\sqrt{d} M_{1,1} + C_\Sigma C_Q \sigma_2 M_{2,1}, \\ M_{3,3} &= 2\sqrt{d} M_{1,2} \sigma_4^2, \\ M_{3,4} &= (2\sqrt{d} M_{1,3} + M_{2,3} C_\Sigma C_Q \sigma_2) \sigma_2 + C_\Sigma C_Q M_{2,2} \sigma_4^2 \sigma_2. \end{aligned}$$

Constants $M_{1,1}, M_{1,2}, M_{1,3}$ are defined in (31) and $M_{2,1}, M_{2,2}, M_{3,3}$ are defined in (32). We simplify the last inequality and get the statement of the theorem with

$$\begin{aligned} C_4 &= \sqrt{d} M_{3,1} + M_{3,2} \|\theta_0 - \theta^*\| + \|\theta_0 - \theta^*\|^2 + \sigma_2 + \sigma_4^2, \\ C_5 &= M_{3,3}, \\ C_6 &= M_{3,4}. \end{aligned} \tag{29}$$

□

B.1 Bounds for D_n

For simplicity of notations we define

$$\begin{aligned} T_1(A) &= 1 + \frac{1}{A^{1/(1-\gamma)}(1-\gamma)} \Gamma\left(\frac{1}{1-\gamma}\right), \\ T_2(A) &= 1 + \max\left(\exp\left\{\frac{1}{1-\gamma}\right\} \frac{1}{A^{1/(1-\gamma)}(1-\gamma)} \Gamma\left(\frac{1}{1-\gamma}\right), \frac{1}{A(1-\gamma)^2}\right). \end{aligned} \tag{30}$$

Lemma 4. *Assume A1, A3, A7(4) and A8. Then it holds that*

$$\mathbb{E}^{1/2}[\|D_n\|^2] \leq \frac{M_{1,1}}{\sqrt{n}} (\|\theta_0 - \theta^*\| + \|\theta_0 - \theta^*\|^2 + \sigma_2 + \sigma_4^2) + M_{1,2} \sigma_4^2 n^{1/2-\gamma} + M_{1,3} \sigma_2 n^{-\gamma/2},$$

where

$$\begin{aligned} M_{1,1} &= C_\Sigma C_Q \left(T_1\left(\frac{\mu c_0}{4}\right) (L_2 + L_H) \max(\sqrt{C_{4,1}}, \sqrt{C_1}) + k_0^2/c_0 \right) \\ M_{1,2} &= C_\Sigma C_Q L_H \sqrt{C_{4,2}} c_0 \frac{1}{1-\gamma} \\ M_{1,3} &= C_\Sigma C_Q L_2 \sqrt{C_2} \sqrt{c_0} \sqrt{\frac{1}{1-\gamma}}, \end{aligned} \tag{31}$$

where $C_{4,1}$ and $C_{4,2}$ are defined in Corollary 4, C_1 and C_2 are defined in Lemma 12 and $T_1(\cdot)$ is defined in eq. (30).

Proof. Using Minkowski's inequality and the definition of D_n , we obtain

$$\mathbb{E}^{1/2}[\|D_n\|^2] \leq \mathbb{E}^{1/2}[\|D_{n,1}\|^2] + \mathbb{E}^{1/2}[\|D_{n,2}\|^2] + \mathbb{E}^{1/2}[\|D_{n,3}\|^2],$$

and consider each of the terms $D_{n,1}, D_{n,2}, D_{n,3}$ separately. Applying Lemma 3, we get

$$\mathbb{E}^{1/2}[\|D_{n,1}\|^2] \leq \frac{C_\Sigma C_Q k_0^\gamma}{\sqrt{nc_0}} \|\theta_0 - \theta^*\|.$$

Now we consider the term $D_{n,2}$. Applying Minkowski's inequality, Lemma 3 and Lemma 18, we have

$$\mathbb{E}^{1/2}[\|D_{n,2}\|^2] \leq \frac{C_\Sigma C_Q}{\sqrt{n}} \sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|H(\theta_{i-1})\|^2] \leq \frac{C_\Sigma C_Q L_H}{\sqrt{n}} \sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|\theta_{i-1} - \theta^*\|^4].$$

For $D_{n,3}$ we note that $\{g(\theta_{i-1}, \xi_i)\}_{i=1}^{n-1}$ is a martingale difference with respect to \mathcal{F}_i . Hence, using Lemma 3 and A7, we get

$$\mathbb{E}^{1/2}[\|D_{n,3}\|^2] \leq \frac{C_\Sigma C_Q}{\sqrt{n}} \left(\sum_{i=1}^{n-1} \mathbb{E}[\|g(\theta_{i-1}, \xi_i)\|^2] \right)^{1/2} \leq \frac{C_\Sigma C_Q L_2}{\sqrt{n}} \left(\mathbb{E} \left[\sum_{i=1}^{n-1} \|\theta_{i-1} - \theta^*\|^2 \right] \right)^{1/2}.$$

Hence, it is enough to upper bound $\mathbb{E}[\|\theta_i - \theta^*\|^{2p}]$ for $p = 1$ and $p = 2$ and $i \in \{0, \dots, n-2\}$. Using Lemma 12 and Lemma 17, we obtain

$$\begin{aligned} \left(\sum_{i=0}^{n-2} \mathbb{E}[\|\theta_i - \theta^*\|^2] \right)^{1/2} &\leq \left(\sum_{i=0}^{n-2} C_1 \exp\left\{ -\frac{\mu c_0}{4} (i + k_0)^{1-\gamma} \right\} [\|\theta_0 - \theta^*\|^2 + \sigma_2^2] + C_2 \sigma_2^2 \alpha_i \right)^{1/2} \\ &\leq \sqrt{C_1} \sqrt{T_1\left(\frac{\mu c_0}{4}\right)} [\|\theta_0 - \theta^*\| + \sigma_2] + \sqrt{C_2} \sigma_2 \sqrt{c_0} \left(\frac{(n-2+k_0)^{1-\gamma} - (k_0-1)^{1-\gamma}}{1-\gamma} \right)^{1/2}, \end{aligned}$$

where $T_1(\cdot)$ is defined in (30). Using Corollary 4 and Lemma 17, we get

$$\begin{aligned} \sum_{i=0}^{n-2} \mathbb{E}^{1/2}[\|\theta_i - \theta^*\|^4] &\leq \sum_{i=0}^{n-2} \sqrt{C_{4,1}} \exp\left\{ -\frac{\mu c_0}{4} i^{1-\gamma} \right\} [\|\theta_0 - \theta^*\|^2 + \sigma_4^2] + \sqrt{C_{4,2}} \sigma_4^2 \alpha_i \\ &\leq \sqrt{C_{4,1}} T_1\left(\frac{\mu c_0}{4}\right) [\|\theta_0 - \theta^*\|^2 + \sigma_4^2] + \sqrt{C_{4,2}} \sigma_4^2 c_0 \left(\frac{(n-2+k_0)^{1-\gamma} - (k_0-1)^{1-\gamma}}{1-\gamma} \right). \end{aligned}$$

We finish the proof, using simple inequality $(n-2+k_0)^{1-\gamma} - (k_0-1)^{1-\gamma} \leq n^{1-\gamma}$ □

B.2 Bounds for $D_n - D_n^{(i)}$

Lemma 5. *Assume A1, A3, A7(4) and A8. Then it holds that*

$$\sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|D_n - D_n^{(i)}\|^2] \leq \frac{M_{2,1}}{\sqrt{n}} (\|\theta_0 - \theta^*\| + \|\theta_0 - \theta^*\|^2 + \sigma_2 + \sigma_4^2) + M_{2,2} \sigma_4^2 n^{1/2-\gamma} + M_{2,3} \sigma_2 n^{1/2-\gamma/2},$$

where

$$\begin{aligned} M_{2,1} &= C_\Sigma C_Q T_1\left(\frac{\mu c_0}{8}\right) T_2\left(\frac{\mu c_0}{1-\gamma}\right) (L_2 + L_H) \max(\sqrt{2(C_1 + c_0^2 k_0^{-\gamma} R_1 R_2)}, c^2 k_0^{-\gamma} \sqrt{R_{4,1} R_{4,2}}) \\ M_{2,2} &= C_\Sigma C_Q L_H c_0 \sqrt{R_{4,1} R_{4,3}} T_2\left(\frac{\mu c_0}{1-\gamma}\right) \frac{1}{1-\gamma} \\ M_{2,3} &= \sqrt{2} C_\Sigma C_Q L_2 \sqrt{C_2 + R_1 R_3 c_0} T_2\left(\frac{\mu c_0}{1-\gamma}\right) \frac{1}{1-\gamma/2}. \end{aligned} \tag{32}$$

Constants R_1, R_2, R_3 are defined in (34) and constants $R_{4,1}, R_{4,2}, R_{4,3}$ are defined (35).

Proof. Using Minkowski's inequality and the definition of D_n and $D_n^{(i)}$, we obtain

$$\sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|D_n - D_n^{(i)}\|^2] \leq \sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|D_{n,2} - D_{n,2}^{(i)}\|^2] + \sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|D_{n,3} - D_{n,3}^{(i)}\|^2]$$

Define $\mathcal{F}_j^{(i)} = \mathcal{F}_j$ if $j \leq i$ and $\mathcal{F}_j^{(i)} = \sigma(\mathcal{F}_j \vee \sigma(\xi'_i))$ otherwise. Then $\{g(\theta_{j-1}, \xi_j) - g(\theta_{j-1}^{(i)}, \tilde{\xi}_j)\}_{j=1}^{n-1}$ is a martingale difference with respect to $\mathcal{F}_j^{(i)}$. Hence, we have, using Lemma 3 and the fact that $\theta_{j-1} = \theta_{j-1}^{(i)}$ for $j \leq i$, we obtain that

$$\begin{aligned} \mathbb{E}[\|D_{n,3} - D_{n,3}^{(i)}\|^2] &= \mathbb{E}\left\|\frac{\Sigma_n^{-1/2}}{\sqrt{n}} \sum_{j=1}^{n-1} Q_j(g(\theta_{j-1}, \xi_j) - g(\theta_{j-1}^{(i)}, \tilde{\xi}_j))\right\|^2 \\ &\leq \frac{C_\Sigma^2 C_Q^2}{n} \mathbb{E}[\|g(\theta_{i-1}, \xi_i) - g(\theta_{i-1}, \xi'_i)\|^2] + \frac{C_\Sigma^2 C_Q^2}{n} \sum_{j=i+1}^{n-1} \mathbb{E}[\|g(\theta_{j-1}, \xi_j) - g(\theta_{j-1}^{(i)}, \xi_j)\|^2]. \end{aligned}$$

Using A7 and Lemma 3, we get

$$\mathbb{E}[\|D_{n,3} - D_{n,3}^{(i)}\|^2] \leq \frac{2C_\Sigma^2 C_Q^2 L_2^2}{n} \mathbb{E}[\|\theta_{i-1} - \theta^*\|^2] + \frac{C_\Sigma^2 C_Q^2 L_2^2}{n} \sum_{j=i+1}^{n-1} \mathbb{E}[\|\theta_{j-1} - \theta_{j-1}^{(i)}\|^2].$$

Using Lemma 6 and Lemma 17, we obtain

$$\begin{aligned} \sum_{j=i+1}^{n-1} \mathbb{E}[\|\theta_{j-1} - \theta_{j-1}^{(i)}\|^2] &\leq R_1 R_2 \exp\left\{-\frac{\mu c_0}{4}(i+k_0-1)^{1-\gamma}\right\} \alpha_i^2 (\|\theta_0 - \theta^*\|^2 + \sigma_2^2) T_2\left(\frac{\mu c_0}{1-\gamma}\right) (i+k_0)^\gamma \\ &\quad + R_1 R_3 \sigma_2^2 \alpha_i^2 T_2\left(\frac{\mu c_0}{1-\gamma}\right) (i+k_0)^\gamma \\ &\leq R_1 R_3 \sigma_2^2 c_0 T_2\left(\frac{\mu c_0}{1-\gamma}\right) \alpha_i + R_1 R_2 c_0^2 k_0^{-\gamma} T_2\left(\frac{\mu c_0}{1-\gamma}\right) \exp\left\{-\frac{\mu c_0}{4}(i+k_0-1)^{1-\gamma}\right\} (\|\theta_0 - \theta^*\|^2 + \sigma_2^2). \end{aligned}$$

Combining inequalities above, we get

$$\begin{aligned} \sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|D_{n,3} - D_{n,3}^{(i)}\|^2] &\leq \frac{\sqrt{2} C_\Sigma C_Q L_2}{\sqrt{n}} \sqrt{C_1 + c_0^2 k_0^{-\gamma} R_1 R_2 T_2\left(\frac{\mu c_0}{1-\gamma}\right) T_1\left(\frac{\mu c_0}{8}\right) (\|\theta_0 - \theta^*\| + \sigma_2)} \\ &\quad + \frac{\sqrt{2} C_\Sigma C_Q L_2}{\sqrt{n}} \sqrt{C_2 + R_1 R_3 c_0 T_2\left(\frac{\mu c_0}{1-\gamma}\right) \sigma_2 \left(\frac{(n+k_0-2)^{1-\gamma/2} - (k_0-1)^{1-\gamma/2}}{1-\gamma/2}\right)}. \end{aligned}$$

We now proceed with $\sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|D_{n,2} - D_{n,2}^{(i)}\|^2]$. Using Minkowski's inequality together with Lemma 3 and Lemma 18, we get

$$\mathbb{E}^{1/2}[\|D_{n,2} - D_{n,2}^{(i)}\|^2] \leq \frac{C_\Sigma C_Q L_H}{\sqrt{n}} \sum_{j=i+1}^{n-1} \mathbb{E}^{1/2}[\|\theta_{j-1} - \theta_{j-1}^{(i)}\|^4].$$

Applying Lemma 7 and Lemma 17, we get using that $\alpha_i^2 (i+k_0)^\gamma \leq \alpha_0^2 k_0^{-\gamma}$ that

$$\begin{aligned} \sum_{j=i+1}^{n-1} \mathbb{E}^{1/2}[\|\theta_{j-1} - \theta_{j-1}^{(i)}\|^4] &\leq c_0^2 k_0^{-\gamma} \sqrt{R_{4,1} R_{4,2}} T_2\left(\frac{\mu c_0}{1-\gamma}\right) \exp\left\{-\frac{\mu c_0}{4}(i+k_0-1)^{1-\gamma}\right\} (\|\theta_0 - \theta^*\|^2 + \sigma_4^2) \\ &\quad + \alpha_i c_0 \sqrt{R_{4,1} R_{4,3}} T_2\left(\frac{\mu c_0}{1-\gamma}\right) \sigma_4^2. \end{aligned}$$

Finally, applying Lemma 17, we get

$$\begin{aligned} \sum_{i=1}^{n-1} \mathbb{E}^{1/2}[\|D_{n,2} - D_{n,2}^{(i)}\|^2] &\leq \frac{C_\Sigma C_Q L_H}{\sqrt{n}} c_0^2 k_0^{-\gamma} \sqrt{R_{4,1} R_{4,2}} T_2\left(\frac{\mu c_0}{1-\gamma}\right) T_1\left(\frac{\mu c_0}{4}\right) (\|\theta_0 - \theta^*\|^2 + \sigma_4^2) \\ &\quad + \frac{C_\Sigma C_Q L_H}{\sqrt{n}} c_0 \sqrt{R_{4,1} R_{4,3}} T_2\left(\frac{\mu c_0}{1-\gamma}\right) \sigma_4^2 \left(\frac{(n+k_0-2)^{1-\gamma} - (k_0-1)^{1-\gamma}}{1-\gamma}\right). \end{aligned}$$

We finish the proof, using that $(n-2+k_0)^\beta - (k_0-1)^\beta \leq n^\beta$ for $\beta \in (0, 1)$. \square

B.3 Bounds for $\theta_k^{(i)} - \theta_k$

Let $(\xi'_1, \dots, \xi'_{n-1})$ be an independent copy of $(\xi_1, \dots, \xi_{n-1})$. For each $1 \leq i \leq n-1$, we construct the sequence $\theta_k^{(i)}$, $1 \leq k \leq n-1$, as follows:

$$\theta_k^{(i)} = \begin{cases} \theta_k, & \text{if } k < i \\ \theta_{k-1}^{(i)} - \alpha_k(\nabla f(\theta_{k-1}^{(i)}) + g(\theta_{k-1}^{(i)}, \xi'_k) + \eta(\xi'_k)), & \text{if } k = i \\ \theta_{k-1}^{(i)} - \alpha_k(\nabla f(\theta_{k-1}^{(i)}) + g(\theta_{k-1}^{(i)}, \xi_k) + \eta(\xi_k)), & \text{if } k > i. \end{cases} \quad (33)$$

Lemma 6. *Assume A1, A3, A7(2) and A8. Then for any $k \in \mathbb{N}$ and $1 \leq i \leq n-1$ it holds*

$$\mathbb{E}[\|\theta_k^{(i)} - \theta_k\|^2] \leq \alpha_i^2 R_1 \exp\left\{-2\mu \sum_{j=i+1}^k \alpha_j\right\} \left(R_2 \exp\left\{-\frac{\mu C_0}{4}(i+k_0-1)^{1-\gamma}\right\} (\|\theta_0 - \theta^*\|^2 + \sigma_2^2) + R_3 \sigma_2^2\right),$$

where we have set

$$R_1 = 4 \exp\left\{\frac{2c_0^2(L_1 + L_2)^2}{2\gamma - 1}\right\}, \quad R_2 = L_2^2 C_1, \quad R_3 = (1 + C_2 L_2). \quad (34)$$

And constant C_1 and C_2 are defined in Lemma 12.

Proof. By construction (33), we have

$$\theta_k^{(i)} - \theta_k = \begin{cases} 0, & \text{if } k < i \\ -\alpha_k(g(\theta_{k-1}, \xi'_k) + \eta(\xi'_k) - g(\theta_{k-1}, \xi_k) - \eta(\xi_k)), & \text{if } k = i \\ \theta_{k-1}^{(i)} - \theta_{k-1} - \alpha_k(\nabla f(\theta_{k-1}^{(i)}) - \nabla f(\theta_{k-1}) + g(\theta_{k-1}^{(i)}, \xi_k) - g(\theta_{k-1}, \xi_k)), & \text{if } k > i \end{cases}$$

Since ξ'_i is independent copy of ξ_i , we obtain

$$\begin{aligned} \mathbb{E}[\|\theta_i^{(i)} - \theta_i\|^2] &\stackrel{(a)}{\leq} 4\alpha_i^2 (L_2^2 \mathbb{E}[\|\theta_{i-1} - \theta^*\|^2] + \sigma_2^2) \\ &\stackrel{(b)}{\leq} 4\alpha_i^2 \left(L_2^2 C_1 \exp\left\{-\frac{\mu C_0}{4}(i+k_0-1)^{1-\gamma}\right\} (\|\theta_0 - \theta^*\|^2 + \sigma_2^2) + (1 + C_2 L_2) \sigma_2^2\right), \end{aligned}$$

where in (a) we used A7, and in (b) we used Lemma 12 and $\alpha_{k-1} L_2 \leq 1$. For $k > i$, applying A7 and A1, we have

$$\begin{aligned} \mathbb{E}[\|\theta_k^{(i)} - \theta_k\|^2 | \mathcal{F}_{k-1}] &\leq \|\theta_{k-1}^{(i)} - \theta_{k-1}\|^2 - 2\alpha_k \langle \theta_{k-1}^{(i)} - \theta_{k-1}, \nabla f(\theta_{k-1}^{(i)}) - \nabla f(\theta_{k-1}) \rangle \\ &\quad + 2\alpha_k^2 (L_1 + L_2)^2 \|\theta_{k-1}^{(i)} - \theta_{k-1}\|^2. \end{aligned}$$

Taking expectation from both sides and applying A1 with Lemma 15(a), we obtain

$$\begin{aligned} \mathbb{E}[\|\theta_k^{(i)} - \theta_k\|^2] &\leq (1 - 2\alpha_k \mu + 2\alpha_k^2 (L_1 + L_2)^2) \mathbb{E}[\|\theta_{k-1}^{(i)} - \theta_{k-1}\|^2] \\ &\leq \exp\left\{\frac{2c_0^2(L_1 + L_2)^2}{2\gamma - 1}\right\} \exp\left\{-2\mu \sum_{j=i+1}^k \alpha_j\right\} \mathbb{E}[\|\theta_i^{(i)} - \theta_i\|^2]. \end{aligned}$$

Combining the above inequalities completes the proof. \square

Lemma 7. *Assume A1, A3, A7(4) and A8. Then for any $k \in \mathbb{N}$ and $1 \leq i \leq n-1$ it holds*

$$\mathbb{E}[\|\theta_k^{(i)} - \theta_k\|^4] \leq \alpha_i^4 R_{4,1} \exp\left\{-4\mu \sum_{j=i+1}^k \alpha_j\right\} \left(R_{4,2} \exp\left\{-\frac{2\mu C_0}{4}(i+k_0-1)^{1-\gamma}\right\} (\|\theta_0 - \theta^*\|^4 + \sigma_4^4) + R_{4,3} \sigma_4^4\right)$$

where we have set

$$R_{4,1} = 64 \exp\left\{\frac{4(L_1 + L_2)^2(1 + 3c_0(L_1 + L_2))^2}{2\gamma - 1}\right\}, \quad R_{4,2} = L_2^4 C_{4,1}, \quad R_{4,3} = 1 + L_2^2 C_{4,2}. \quad (35)$$

And constant $C_{4,1}$, $C_{4,2}$ are defined in Corollary 4.

Proof. Repeating the proof of the Lemma 6 for $k = i$, we get

$$\begin{aligned} \mathbb{E}[\|\theta_i^{(i)} - \theta_i\|^4] &\leq 64\alpha_i^4(L_2^4\mathbb{E}[\|\theta_{i-1} - \theta^*\|^4] + \sigma_4^4) \\ &\leq 64\alpha_i^4\left(L_2^4C_{4,1}\exp\left\{-\frac{2\mu c_0}{4}(i+k_0-1)^{1-\gamma}\right\}(\|\theta_0 - \theta^*\|^4 + \sigma_4^4) + (1 + L_2^2C_{4,2})\sigma_4^4\right). \end{aligned}$$

For $k > i$ we denote $\delta_k^{(i)} = \|\theta_k^{(i)} - \theta_k\|$, similar to (60), we obtain

$$E[\{\delta_k^{(i)}\}^4|\mathcal{F}_{k-1}] \leq (1 - 4\mu\alpha_k + 4\alpha_k^2(L_1 + L_2)^2(1 + 3c_0(L_1 + L_2))^2)\{\delta_{k-1}^{(i)}\}^4.$$

Using Lemma 15(a), we obtain

$$E[\{\delta_k^{(i)}\}^4] \leq \exp\left\{\frac{4(L_1 + L_2)^2(1 + 3c_0(L_1 + L_2))^2}{2\gamma - 1}\right\} \exp\left\{-4\mu \sum_{j=i+1}^k \alpha_j\right\} \mathbb{E}[\|\theta_i^{(i)} - \theta_i\|^4].$$

Combining the above inequalities completes the proof. \square

C Proof of Theorem 3

Since the matrix Σ_n^b concentrates around Σ_n due to Lemma 2, there is a set Ω_1 such that $\mathbb{P}(\Omega_1) \geq 1 - 1/n$ and $\lambda_{\min}(\Sigma_n^b) > 0$ on Ω_1 . Moreover, on this set Applying Lemma 1 with

$$X = \{\Sigma_n^b\}^{-1/2}W^b, \quad Y = \{\Sigma_n^b\}^{-1/2}D^b,$$

we get

$$\begin{aligned} &\sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\sqrt{n}\{\Sigma_n^b\}^{-1/2}(\bar{\theta}_n^b - \bar{\theta}_n) \in B) - \mathbb{P}^b(Y^b \in B)| \\ &\leq \sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\{\Sigma_n^b\}^{-1/2}W^b \in B) - \mathbb{P}^b(Y^b \in B)| + 2c_d(\mathbb{E}^b[\|\{\Sigma_n^b\}^{-1/2}D^b\|^p])^{1/(1+p)}. \end{aligned}$$

By (Shao and Zhang, 2022) (with $D = 0$) we may estimate

$$\begin{aligned} &\sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\{\Sigma_n^b\}^{-1/2}W^b \in B) - \mathbb{P}^b(Y^b \in B)| \\ &\leq \frac{259d^{1/2}}{n^{3/2}} \sum_{i=1}^n \mathbb{E}^b[|w_i - 1|^3] \|(\{\Sigma_n^b\}^{-1/2}Q_i\eta(\xi_i))\|^3. \end{aligned} \tag{36}$$

Applying Lemma 3 and Corollary 1 we get

$$\sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}^b(\{\Sigma_n^b\}^{-1/2}W^b \in B) - \mathbb{P}^b(Y^b \in B)| \leq \frac{259d^{1/2}(\sqrt{2}C_\Sigma C_Q C_{1,\xi})^3 W_{\max}}{n^{1/2}}.$$

From Proposition 2 and Corollary 1 it follows that on the set $\Omega_0 \cap \Omega_1$ the following bound is satisfied

$$(\mathbb{E}^b[\|\{\Sigma_n^b\}^{-1/2}D^b\|^p])^{1/(p+1)} \leq \sqrt{2}C_\Sigma(M_{1,1}^b e^{1/p} p^{3/2} n^{1/p-\gamma/2} + M_{2,1}^b e^{2/p} p n^{1/2+1/p-\gamma} p/(p+1)).$$

Since $p \geq 2$, $M_{1,1}^b, M_{2,1}^b \geq 1$, we obtain

$$(\mathbb{E}^b[\|\{\Sigma_n^b\}^{-1/2}D^b\|^p])^{1/(p+1)} \leq \sqrt{2}C_\Sigma(e^{1/2}M_{1,1}^b p^{3/2} n^{\frac{1}{p+1}} n^{-\gamma/2} n^{\frac{\gamma/2}{(p+1)}} + eM_{2,1}^b p n^{\frac{1}{p+1}} n^{1/2-\gamma} n^{-\frac{1/2-\gamma}{p+1}}).$$

Setting $p = \log n - 1$, we get

$$(\mathbb{E}^b[\|\{\Sigma_n^b\}^{-1/2}D^b\|^p])^{1/(p+1)} \leq \sqrt{2}C_\Sigma(M_{1,1}^b (\log n)^{3/2} e^{3/2+\gamma/2} n^{-\gamma/2} + M_{2,1}^b (\log n) e^{3/2+\gamma} n^{1/2-\gamma}).$$

Setting

$$\begin{aligned}
 M_{3,1}^b &= 259(\sqrt{2}C_\Sigma C_Q C_{1,\xi})^3 W_{\max} \sqrt{d}, \\
 M_{3,2}^b &= 2^{3/2} c_d C_\Sigma M_{2,1}^b e^{3/2+\gamma}, \\
 M_{3,2}^b &= 2^{3/2} c_d C_\Sigma M_{1,1}^b e^{3/2+\gamma/2},
 \end{aligned} \tag{37}$$

and $M_{1,1}^b, M_{2,1}^b$ are defined in (41) and combining the above inequalities, we complete the proof.

Remark 7. We use (Shao and Zhang, 2022) with $D = 0$ to prove (36) since we are not aware of Berry-Esseen results for non i.i.d. random vectors in dimension d with precise constants and dependence on d . The result Bentkus (2003) may be applied for i.i.d. vectors only.

C.1 From non-linear to linear statistics

In this section we prove (20). We start from the definition of an isoperimetric constant. Define

$$A^\varepsilon = \{x \in \mathbb{R}^d : \rho_A(x) \leq \varepsilon\} \quad \text{and} \quad A^{-\varepsilon} = \{x \in A : B_\varepsilon(x) \subset A\},$$

where $\rho_A(x) = \inf_{y \in A} \|x - y\|$ is the distance between $A \subset \mathbb{R}^d$ and $x \in \mathbb{R}^d$, and

$$B_\varepsilon(x) = \{y \in \mathbb{R}^d : \|x - y\| \leq \varepsilon\}.$$

For some class \mathcal{A} of subsets of \mathbb{R}^d we define its isoperimetric constant $a_d(\mathcal{A})$ (depending only on d and \mathcal{A}) as follows: for all $A \in \mathcal{A}$ and $\varepsilon > 0$,

$$\mathbb{P}\{Y \in A^\varepsilon \setminus A\} \leq a_d \varepsilon, \quad \mathbb{P}\{Y \in A \setminus A^{-\varepsilon}\} \leq a_d \varepsilon$$

where Y follows the standard Gaussian distribution on \mathbb{R}^d . Ball (1993) has proved that

$$e^{-1} \sqrt{\ln d} \leq \sup_{A \in \mathcal{C}} \int_{\partial A} p(x) ds \leq 4d^{1/4}, \tag{38}$$

where $p(x)$ is the standard normal d -dimensional density and ds is the surface measure on the boundary ∂A of A . Using (38) one can show that for the class of convex sets

$$e^{-1} \sqrt{\ln d} \leq a_d(\mathcal{C}(\mathbb{R}^d)) \leq 4d^{1/4}.$$

We denote $c_d = a_d(\mathcal{C}(\mathbb{R}^d))$.

Proposition 1. Let ν be a standard Gaussian measure in \mathbb{R}^d . Then for any random vectors X, Y taking values in \mathbb{R}^d , and any $p \geq 1$,

$$\sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}(X + Y \in B) - \nu(B)| \leq \sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}(X \in B) - \nu(B)| + 2c_d^{p/(p+1)} \mathbb{E}^{1/(p+1)}[\|Y\|^p],$$

where c_d is the isoperimetric constant of class $\mathcal{C}(\mathbb{R}^d)$.

Proof. Let $\varepsilon \geq 0$. Define $\rho(B) = \mathbb{P}(X + Y \in B) - \nu(B)$. Let B be such that $\rho(B) \geq 0$. By Markov's inequality

$$\begin{aligned}
 \rho(B) &\leq \mathbb{P}(X + Y \in B, |Y| \leq \varepsilon) + \frac{1}{\varepsilon^p} \mathbb{E}[\|Y\|^p] - \nu(B) \\
 &\leq \sup_A |\mathbb{P}(X \in A) - \nu(A)| + \mathbb{P}(Y \in B^\varepsilon \setminus B) + \frac{1}{\varepsilon^p} \mathbb{E}[\|Y\|^p].
 \end{aligned}$$

Choosing

$$\varepsilon = \frac{1}{c_d^{1/(p+1)}} \mathbb{E}^{1/(p+1)}[\|Y\|^p] \tag{39}$$

we obtain

$$\sup_B |\mathbb{P}(X + Y \in B) - \nu(B)| \leq \sup_B |\mathbb{P}(X \in B) - \nu(B)| + 2c_d^{p/(p+1)} \mathbb{E}^{1/(p+1)}[\|Y\|^p].$$

Assume now that $\rho(B) < 0$. We distinguish between $B^{-\varepsilon} = \emptyset$ or $B^{-\varepsilon} \neq \emptyset$. In the first case, $\mathbb{P}(Y \in B^{-\varepsilon}) = 0$ and

$$-\rho(B) \leq \gamma(B) = \mathbb{P}(Y \in B) - \mathbb{P}(Y \in B^{-\varepsilon}) = \mathbb{P}(Y \in B \setminus B^{-\varepsilon}) \leq c_d \varepsilon.$$

Finally, in the case $B^{-\varepsilon} \neq \emptyset$,

$$-\rho(B) \leq \sup_A |\mathbb{P}(X \in A) - \nu(A)| + \mathbb{P}(Y \in B \setminus B^{-\varepsilon}) + \frac{1}{\varepsilon^p} \mathbb{E}[\|Y\|^p].$$

Taking ε as in (39) we conclude the proof. \square

C.2 Bounds for D^b

Recall that the term D^b defined in (17), has a form:

$$\begin{aligned} D^b &= -\frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} (w_i - 1) Q_i \left(G(\theta_{i-1}^b - \theta^*) + g(\theta_{i-1}^b, \xi_i) + H(\theta_{i-1}^b) \right) \\ &\quad - \frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i \left(H(\theta_{i-1}^b) + g(\theta_{i-1}^b, \xi_i) - H(\theta_{i-1}) - g(\theta_{i-1}, \xi_i) \right). \end{aligned}$$

To prove Theorem 3, we need to obtain a high-probability bound for the non-linear statistic D^b . To this end, we first derive a bound on $\mathbb{E}^{1/p}[\|D^b\|^p]$, where the expectation is taken with respect to the joint distribution of the bootstrap weights and the data Ξ^{n-1} . We then apply Markov's inequality to convert this moment bound into a high-probability bound.

Proposition 2. *Assume A1- A5. Then it holds for any $p \geq 2$ that*

$$\mathbb{E}^{1/p}[\|D^b\|^p] \leq M_{1,1}^b e^{1/p} p^{3/2} n^{-\gamma/2} + M_{2,1}^b e^{2/p} p n^{1/2-\gamma}, \quad (40)$$

where the constants are given by

$$\begin{aligned} M_{1,1}^b &= 4C_Q \max(L_1, L_2) \frac{\max(\sqrt{K_2}, \sqrt{K_1}) \sqrt{c_0 k_0^{1-\gamma}} (W_{\max} + 1)}{\sqrt{2}(1-\gamma)}, \\ M_{2,1}^b &= 3C_Q L_H \frac{c_0 k_0^{1-\gamma} \max(K_2, K_1) (W_{\max} + 1)}{2(1-\gamma)}, \end{aligned} \quad (41)$$

and K_1, K_2 are defined in (62), (66), respectively. Moreover, there is a set $\Omega_0 \in \mathcal{F}_{n-1} = \sigma(\xi_1, \dots, \xi_{n-1})$, such that $\mathbb{P}(\Omega_0) \geq 1 - 1/n$, and on Ω_0 it holds that

$$\{\mathbb{E}^b[\|D^b\|^p]\}^{1/p} \leq M_{1,1}^b e^{1/p} p^{3/2} n^{1/p-\gamma/2} + M_{2,1}^b e^{2/p} p n^{1/2+1/p-\gamma}. \quad (42)$$

Proof. We first show (40). We split

$$D^b = D_1^b + D_2^b,$$

where

$$\begin{aligned} D_1^b &= -\frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} (w_i - 1) Q_i (G(\theta_{i-1}^b - \theta^*) + g(\theta_{i-1}^b, \xi_i)) - \frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i (g(\theta_{i-1}^b, \xi_i) - g(\theta_{i-1}, \xi_i)), \\ D_2^b &= -\frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} (w_i - 1) Q_i H(\theta_{i-1}^b) - \frac{1}{\sqrt{n}} \sum_{i=1}^{n-1} Q_i (H(\theta_{i-1}^b) - H(\theta_{i-1})). \end{aligned}$$

Applying Minkowski's inequality together with Lemma 8 and Lemma 9 we get (40).

To proof (42) we consider

$$\Omega_0 = \{ \{\mathbb{E}^b[\|D^b\|^p]\}^{1/p} \leq M_{1,1}^b e^{1/p} p^{3/2} n^{1/p-\gamma/2} + M_{2,1}^b e^{2/p} p n^{1/2+1/p-\gamma} \}.$$

Note that by Markov's inequality

$$\begin{aligned} \mathbb{P}(\Omega_0^c) &\leq \frac{\mathbb{E}\{\{\mathbb{E}^b[\|D^b\|^p]\}\}}{n(M_{1,1}^b e^{2/p} p^{3/2} n^{-\gamma/2} + M_{2,1}^b e^{1/p} p n^{1/2-\gamma})^p} \\ &= \frac{\mathbb{E}[\|D^b\|^p]}{n(M_{1,1}^b e^{1/p} p^{3/2} n^{-\gamma/2} + M_{2,1}^b e^{2/p} p n^{1/2-\gamma})^p} \leq \frac{1}{n}. \end{aligned}$$

□

Lemma 8. *Assume A1-A5. Then for any $p \geq 2$ it holds*

$$\mathbb{E}^{1/p}[\|D_1^b\|^p] \leq M_{1,1}^b e^{1/p} p^{3/2} n^{-\gamma/2},$$

where

$$M_{1,1}^b = 4C_Q \max(L_1, L_2) \frac{\max(\sqrt{K_2}, \sqrt{K_1}) \sqrt{c_0} (W_{\max} + 1)}{\sqrt{2}(1-\gamma)}, \quad (43)$$

and K_1, K_2 are defined in (62), (66), respectively.

Proof. We split D_1^b into four parts, where each part is a sum of martingale differences. Note that $\{Q_i(g(\theta_{i-1}, \xi_i) - g(\theta^*, \xi_i))\}_{i=1}^n$ is a martingale difference with respect to \mathcal{F}_{i-1} . Then applying Burholder's inequality (Osekowski, 2012, Theorem 8.6) together with Minkowski's inequality and Lemma 3, we obtain that

$$\begin{aligned} \mathbb{E}^{1/p} \left[\left\| \sum_{i=1}^{n-1} Q_i(g(\theta_{i-1}, \xi_i) - g(\theta^*, \xi_i)) \right\|^p \right] &\leq p \left(\mathbb{E}^{2/p} \left[\left(\sum_{i=1}^{n-1} \|Q_i(g(\theta_{i-1}, \xi_i) - g(\theta^*, \xi_i))\|^2 \right)^{p/2} \right] \right)^{1/2} \\ &\leq C_Q p \left(\mathbb{E}^{2/p} \left[\left(\sum_{i=1}^{n-1} \|g(\theta_{i-1}, \xi_i) - g(\theta^*, \xi_i)\|^2 \right)^{p/2} \right] \right)^{1/2} \\ &\leq C_Q p \left(\sum_{i=1}^{n-1} \mathbb{E}^{2/p} [\|g(\theta_{i-1}, \xi_i) - g(\theta^*, \xi_i)\|^p] \right)^{1/2}. \end{aligned}$$

Finally, using A7 and Lemma 14, we obtain

$$\begin{aligned} \mathbb{E}^{1/p} \left[\left\| \sum_{i=1}^{n-1} Q_i(g(\theta_{i-1}, \xi_i) - g(\theta^*, \xi_i)) \right\|^p \right] &\leq p C_Q L_2 \left(\sum_{i=1}^{n-1} \mathbb{E}^{2/p} [\|\theta_{i-1} - \theta^*\|^p] \right)^{1/2} \\ &\leq C_Q L_2 (en)^{1/p} p^{3/2} \frac{\sqrt{K_2}}{\sqrt{2}} \left(\sum_{i=0}^{n-2} \alpha_i \right)^{1/2} \\ &\leq C_Q L_2 (en)^{1/p} p^{3/2} \frac{\sqrt{K_2}}{\sqrt{2}} \left(c_0 \frac{(k_0 + n - 2)^{1-\gamma} - (k_0 - 1)^{1-\gamma}}{1 - \gamma} \right)^{1/2}. \end{aligned}$$

Since $k_0 \geq 1$ and $(k_0 + n - 2)^{1-\gamma} - (k_0 - 1)^{1-\gamma} \leq n^{1-\gamma}$ we complete the proof for

$$\mathbb{E}^{1/p} \left[\left\| \sum_{i=1}^{n-1} Q_i(g(\theta_{i-1}, \xi_i) - g(\theta^*, \xi_i)) \right\|^p \right].$$

The proof for other three terms is analogous, since each of the terms

$$\{Q_i(g(\theta_{i-1}^b, \xi_i) - g(\theta^*, \xi_i))\}_{i=1}^{n-1}, \{(w_i - 1)Q_i(g(\theta_{i-1}^b, \xi_i) - g(\theta^*, \xi_i))\}_{i=1}^{n-1}, \{(w_i - 1)Q_i G(\theta_{i-1}^b - \theta^*)\}_{i=1}^{n-1},$$

are martingale differences with respect to $\tilde{\mathcal{F}}_{i-1}$ (see definition in (19)). We finish the proof applying Minkowski's inequality. □

Lemma 9. *Assume A1- A5. Then for any $p \geq 2$ it holds*

$$\mathbb{E}^{1/p}[\|D_2^b\|^p] \leq M_{2,1}^b e^{2/p} p n^{1/2-\gamma},$$

$$M_{2,1}^b = 3C_Q L_H \frac{c_0 \max(K_2, K_1)(W_{\max} + 1)}{2(1-\gamma)}, \quad (44)$$

and K_1, K_2 are defined in (62), (66), respectively.

Proof. Using Minkowski's inequality, we get

$$\begin{aligned} \mathbb{E}^{1/p}[\|D_2^b\|^p] &\leq \frac{1}{\sqrt{n}} \mathbb{E}^{1/p}[\|\sum_{i=1}^{n-1} Q_i H(\theta_{i-1})\|^p] \\ &\quad + \frac{1}{\sqrt{n}} \mathbb{E}^{1/p}[\|\sum_{i=1}^{n-1} (w_i - 1) Q_i \left(H(\theta_{i-1}^b) \right)\|^p] \\ &\quad + \frac{1}{\sqrt{n}} \mathbb{E}^{1/p}[\|\sum_{i=1}^{n-1} Q_i H(\theta_{i-1}^b)\|^p]. \end{aligned} \quad (45)$$

We will now consider each term separately. Using Minkowski's inequality together with Lemma 18, we obtain

$$\begin{aligned} \frac{1}{\sqrt{n}} \mathbb{E}^{1/p} \left[\left\| \sum_{i=1}^{n-1} Q_i H(\theta_{i-1}) \right\|^p \right] &\leq \frac{C_Q L_H}{\sqrt{n}} \sum_{i=0}^{n-2} \mathbb{E}^{1/p} \left[\|\theta_i - \theta^*\|^{2p} \right] \\ &\leq \frac{C_Q L_H p}{\sqrt{n}} (en)^{2/p} (K_2/2) \sum_{i=0}^{n-1} \alpha_i \\ &\leq \frac{C_Q L_H p}{\sqrt{n}} (en)^{2/p} (K_2/2) \left(c_0 \frac{(k_0 + n - 2)^{1-\gamma} - (k_0 - 1)^{1-\gamma}}{1 - \gamma} \right). \end{aligned}$$

Since $k_0 \geq 1$ and $(k_0 + n - 2)^{1-\gamma} - (k_0 - 1)^{1-\gamma} \leq n^{1-\gamma}$ we complete the proof for the first term in the r.h.s. of (45). The proof for other two terms is analogous. \square

D Proof of Theorem 4

The result follows immediately from the triangle inequality together with Theorem 2 and Lemma 1. Therefore, it remains to prove Lemma 1 to complete the proof.

D.1 Proof of Lemma 1

By definition of Σ_n and Σ_∞ we may write

$$\begin{aligned} \Sigma_n - \Sigma_\infty &= \underbrace{\frac{1}{n} \sum_{t=1}^{n-1} (Q_t - G^{-1}) \Sigma_\xi G^{-\top} + \frac{1}{n} \sum_{t=1}^{n-1} G^{-1} \Sigma_\xi (Q_t - G^{-1})^\top}_{D_1} \\ &\quad + \underbrace{\frac{1}{n} \sum_{t=1}^{n-1} (Q_t - G^{-1}) \Sigma_\xi (Q_t - G^{-1})^\top - \frac{1}{n} \Sigma_\infty}_{D_2}. \end{aligned}$$

The following lemma is an analogue of (Wu et al., 2024, pp. 26-30).

Lemma 10. *The following identities hold*

$$Q_i - G^{-1} = S_i - G^{-1} G_{i:n-1}^{(\alpha)}, \quad S_i = \sum_{j=i+1}^{n-1} (\alpha_i - \alpha_j) G_{i+1:j-1}^{(\alpha)}, \quad (46)$$

and

$$\sum_{i=1}^{n-1} (Q_i - G^{-1}) = -G^{-1} \sum_{j=1}^{n-1} G_{1:j}^{(\alpha)}, \quad (47)$$

where

$$G_{i:j}^{(\alpha)} = \prod_{k=i}^j (I - \alpha_k G)$$

Proof. To prove (47) we first change the order of summation and then use the properties of the telescopic sums we get

$$\begin{aligned} \sum_{i=1}^{n-1} Q_i &= \sum_{i=1}^{n-1} \alpha_i \sum_{j=i}^{n-1} \prod_{k=i+1}^j (I - \alpha_k G) = \sum_{j=1}^{n-1} \sum_{i=1}^j \alpha_i \prod_{k=i+1}^j (I - \alpha_k G) \\ &= \sum_{j=1}^{n-1} \sum_{i=1}^j G^{-1} \left(\prod_{k=i+1}^j - \prod_{k=i}^j \right) (I - \alpha_k G) = G^{-1} \sum_{j=1}^{n-1} \left(I - \prod_{k=1}^j (I - \alpha_k G) \right). \end{aligned}$$

The proof of (46) could be obtained by the following arguments. Note that

$$\begin{aligned} \alpha_i G Q_i &= Q_i - (I - \alpha_i G) Q_i = \\ &= \alpha_i I + \alpha_i \sum_{j=i+1}^{n-1} \prod_{k=i+1}^j (I - \alpha_k G) - \alpha_i \sum_{j=i+1}^{n-1} \prod_{k=i}^{j-1} (I - \alpha_k G) - \alpha_i \prod_{k=i}^{n-1} (I - \alpha_k G). \end{aligned}$$

It remains to note that

$$\prod_{k=i+1}^j (I - \alpha_k G) - \prod_{k=i}^{j-1} (I - \alpha_k G) = (\alpha_i - \alpha_j) G \prod_{k=i+1}^{j-1} (I - \alpha_k G).$$

The last two equations imply (46). \square

Lemma 11. *It holds that*

(a)

$$\|S_i\| \leq C_S (i + k_0)^{\gamma-1},$$

where

$$C_S = 2c_0 \exp\left\{ \frac{\mu c_0}{k_0^\gamma} \right\} \left(2^{\gamma/(1-\gamma)} \frac{1}{\mu c_0} + \left(\frac{1}{\mu c_0} \right)^{1/(1-\gamma)} \Gamma\left(\frac{1}{1-\gamma} \right) \right).$$

(b)

$$\sum_{i=1}^{n-1} \|G_{i:n-1}^{(\alpha)}\|^2 \leq \frac{1}{1 - (1 - c_0 \mu (n + k_0 - 2))^{-\gamma}}$$

(c)

$$\left\| \sum_{i=1}^{n-1} G_{i:n-1}^{(\alpha)} \right\| \leq \frac{k_0^\gamma n^\gamma}{c_0 \mu}$$

Proof. For simplicity we define $m_i^j = \sum_{k=i}^j (k + k_0)^{-\gamma}$. Note that

$$\left\| \sum_{j=i+1}^{n-1} (\alpha_i - \alpha_j) G_{i+1:j-1}^{(\alpha)} \right\| \leq \sum_{j=i}^{n-2} \frac{c_0}{(j + k_0 + 1)^\gamma} \left(\left(\frac{j + k_0 + 1}{i + k_0} \right)^\gamma - 1 \right) \exp\{-\mu c_0 m_{i+1}^j\}$$

Following the proof of (Wu et al., 2024, Lemma A.5), we have

$$\left(\frac{j+k_0+1}{i+k_0}\right)^\gamma - 1 \leq (i+k_0)^{\gamma-1} \left(1 + (1-\gamma)m_i^j\right)^{\gamma/(1-\gamma)}$$

Hence, we obtain

$$\begin{aligned} \|S_i\| &\leq c_0(i+k_0)^{\gamma-1} \sum_{j=i}^{n-2} \frac{1}{(j+k_0+1)^\gamma} \left(1 + (1-\gamma)m_i^j\right)^{\gamma/(1-\gamma)} \exp\{-\mu c_0 m_{i+1}^j\} \\ &\leq c_0(i+k_0)^{\gamma-1} \sum_{j=i}^{n-2} \frac{1}{(j+k_0)^\gamma} \left(1 + (1-\gamma)m_i^j\right)^{\gamma/(1-\gamma)} \exp\{\mu c_0(k_0+i)^{-\gamma}\} \exp\{-\mu c_0 m_i^j\} \\ &\leq c_0 \exp\left\{\frac{\mu c_0}{k_0^\gamma}\right\} (i+k_0)^{\gamma-1} \sum_{j=i}^{n-2} (m_i^j - m_i^{j-1}) \left(1 + (1-\gamma)m_i^j\right)^{\gamma/(1-\gamma)} \exp\{-\mu c_0 m_i^j\} \\ &\leq 2c_0 \exp\left\{\frac{\mu c_0}{k_0^\gamma}\right\} (i+k_0)^{\gamma-1} \int_0^{+\infty} \left(1 + (1-\gamma)m\right)^{\gamma/(1-\gamma)} \exp\{-\mu c_0 m\} dm \\ &\leq 2c_0 \exp\left\{\frac{\mu c_0}{k_0^\gamma}\right\} (i+k_0)^{\gamma-1} \left(2^{\gamma/(1-\gamma)} \frac{1}{\mu c_0} + \left(\frac{1}{\mu c_0}\right)^{1/(1-\gamma)} \Gamma\left(\frac{1}{1-\gamma}\right)\right). \end{aligned}$$

Note that

$$\begin{aligned} \left\| \sum_{i=1}^{n-1} G_{i:n-1}^{(\alpha)} \right\| &\leq \sum_{i=1}^{n-1} \prod_{k=i}^{n-1} (1 - \alpha_k \mu) = \sum_{i=1}^{n-1} \prod_{k=i}^{n-1} \alpha_{i-1}^{-1} \alpha_{i-1} (1 - \alpha_k \mu) \\ &\leq \frac{(k_0 + n - 2)^\gamma}{c_0 \mu} \sum_{i=1}^{n-1} \left(\prod_{k=i}^{n-1} (1 - \alpha_k \mu) - \prod_{k=i-1}^{n-1} (1 - \alpha_k \mu) \right) \leq \frac{k_0^\gamma n^\gamma}{\mu c_0}, \end{aligned}$$

where in the last inequality we use that $(k_0 + n - 2)^\gamma \leq (k_0 n)^\gamma$. Bound for $\sum_{i=1}^{n-1} \|G_{i:n-1}^{(\alpha)}\|^2$ is obtained similarly to $\|\sum_{i=1}^{n-1} G_{i:n-1}^{(\alpha)}\|$. \square

To finish the proof of Lemma 1 we need to bound D_1, D_2 . By (47) we obtain

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1}^{n-1} (Q_i - G^{-1}) \Sigma_\xi G^{-\top} \right\| &= \left\| -\frac{1}{n} G^{-1} \sum_{j=1}^{n-1} G_{1:j}^{(\alpha)} \Sigma_\xi G^{-\top} \right\| \\ &= \|n^{-1} \Sigma_\infty \sum_{j=1}^{n-1} G_{1:j}^{(\alpha)}\| \leq n^{-1} \|\Sigma_\infty\| \cdot \left\| \sum_{j=1}^{n-1} G_{1:j}^{(\alpha)} \right\|. \end{aligned}$$

It remains to apply Lemma 3 which gives

$$\left\| \frac{1}{n} \sum_{i=1}^{n-1} (Q_i - G^{-1}) \Sigma_\xi G^{-\top} \right\| \leq \|\Sigma_\infty\| C_Q \frac{k_0^\gamma n^{\gamma-1}}{c_0}$$

Hence,

$$\|D_1\| \leq 2 \|\Sigma_\infty\| C_Q \frac{k_0^\gamma n^{\gamma-1}}{c_0}$$

To bound D_2 we use (46) which gives

$$\begin{aligned}
 & n^{-1} \sum_{i=1}^{n-1} (Q_i - G^{-1}) \Sigma_\xi (Q_i - G^{-1})^\top \\
 &= n^{-1} \sum_{i=1}^{n-1} (S_i - G^{-1} \prod_{k=i}^{n-1} (I - \alpha_k G)) \Sigma_\xi (S_i - G^{-1} \prod_{k=i}^{n-1} (I - \alpha_k G))^\top \\
 &= n^{-1} \underbrace{\sum_{i=1}^{n-1} S_i \Sigma_\xi S_i^\top}_{D_{21}} + n^{-1} \underbrace{\sum_{i=1}^{n-1} G^{-1} \prod_{k=i}^{n-1} (I - \alpha_k G) \Sigma_\xi G^{-\top} \prod_{k=i}^{n-1} (I - \alpha_k G)^\top}_{D_{22}} \\
 &\quad - n^{-1} \underbrace{\sum_{i=1}^{n-1} G^{-1} \prod_{k=i}^{n-1} (I - \alpha_k G) \cdot \Sigma_\xi S_i^\top}_{D_{23}} - n^{-1} \underbrace{\sum_{i=1}^{n-1} S_i \Sigma_\xi G^{-\top} \prod_{k=i}^{n-1} (I - \alpha_k G)^\top}_{D_{24}}.
 \end{aligned}$$

To bound D_{21} we use Lemma 11, and obtain

$$\begin{aligned}
 \|D_{21}\| &= \left\| n^{-1} \sum_{i=1}^{n-1} S_i \Sigma_\xi S_i^\top \right\| \leq n^{-1} \sum_{i=1}^{n-1} \|\Sigma_\xi\| \|S_i\|^2 \\
 &\leq n^{-1} \|\Sigma_\xi\| C_S^2 \sum_{i=1}^{n-1} (i + k_0)^{2(\gamma-1)} \\
 &\leq n^{-1} \|\Sigma_\xi\| C_S^2 \frac{(n + k_0 - 1)^{2\gamma-1} - k_0^{2\gamma-1}}{2\gamma - 1} \\
 &\leq \|\Sigma_\xi\| C_S^2 \frac{n^{2(\gamma-1)}}{2\gamma - 1}
 \end{aligned}$$

The bound for D_{22} follows from Lemma 11

$$\begin{aligned}
 \|D_{22}\| &= \left\| n^{-1} \sum_{i=1}^{n-1} \prod_{k=i}^{n-1} (I - \alpha_k G) G^{-1} \Sigma_\xi G^{-\top} \prod_{k=i}^{n-1} (I - \alpha_k G)^\top \right\| \leq n^{-1} \|\Sigma_\infty\| \sum_{i=1}^{n-1} \|G_{i:n-1}^{(\alpha)}\|^2 \\
 &\leq n^{-1} \frac{\|\Sigma_\infty\|}{2c_0\mu(n + k_0 - 2)^{-\gamma} - c_0^2\mu^2(n + k_0 - 2)^{-2\gamma}} \leq \|\Sigma_\infty\| k_0^\gamma \frac{n^{\gamma-1}}{c_0\mu}.
 \end{aligned}$$

Since $D_{23} = D_{24}^\top$, we concentrate on $\|D_{24}\|$. Lemma 11 immediately imply

$$\begin{aligned}
 \|D_{24}\| &\leq n^{-1} \|\Sigma_\xi G^{-\top}\| \sum_{i=1}^{n-1} \|S_i\| \left\| \prod_{k=i}^{n-1} (I - \alpha_k G)^\top \right\| \\
 &\leq n^{-1} \|\Sigma_\xi\| \frac{1}{\mu} C_S \sum_{i=1}^{n-1} (i + k_0)^{\gamma-1} \prod_{k=i}^{n-1} \left(1 - \mu \frac{c_0}{(k + k_0)^\gamma}\right) \\
 &\leq n^{-1} \|\Sigma_\xi\| \frac{1}{\mu} C_S \sum_{i=1}^{n-1} (i + k_0)^{2\gamma-1} (i + k_0)^{-\gamma} \prod_{k=i+1}^{n-1} \left(1 - \mu \frac{c_0}{(k + k_0)^\gamma}\right) \\
 &\leq \|\Sigma_\xi\| C_S k_0^{2\gamma-1} \frac{n^{2(\gamma-1)}}{\mu^2 c_0}
 \end{aligned}$$

Combining all inequalities above, we obtain

$$\|\Sigma_n - \Sigma_\infty\| \leq C'_\infty n^{\gamma-1}, \tag{48}$$

where

$$C'_\infty = \left(\frac{k_0^\gamma}{c_0\mu} + 2C_Q \frac{k_0^\gamma}{c_0} + 1 \right) \|\Sigma_\infty\| + \left(C_S^2 \frac{1}{2\gamma - 1} + C_S \frac{k_0^{2\gamma-1}}{\mu^2 c_0} \right) \|\Sigma_\xi\|.$$

To finish the proof it remains to apply Lemma 20, since

$$3/2\|\Sigma_n^{-1/2}\Sigma_\infty\Sigma_n^{-1/2} - I\|_F \leq C_\infty n^{\gamma-1}, \text{ where } C_\infty = 3/2\sqrt{d}C_\Sigma^2 C'_\infty. \quad (49)$$

E Lower bounds

In the following computations we provide a lower bound on the quantity $|\frac{1}{n}\sum_{j=1}^{n-1} Q_j^2 - 1|$, provided that the number of observations n is large enough. For simplicity in this bound we consider $k_0 = 1$. We first note that

$$\frac{1}{n}\sum_{j=1}^{n-1} Q_j^2 - 1 = \frac{1}{n}\sum_{j=1}^{n-1} (Q_j - 1)(Q_j + 1) - \frac{1}{n} = \frac{T_1}{n} + \frac{T_2}{n},$$

where

$$T_1 = \sum_{j=1}^{n-1} (Q_j - 1)^2, \quad T_2 = -2\sum_{j=1}^{n-1} (Q_j - 1) - 1,$$

and treat the terms T_1 and T_2 separately. Using the identity (47), we get, since $G = 1$, that

$$\sum_{j=1}^{n-1} (Q_j - 1) = -\sum_{j=1}^{n-1} \prod_{\ell=1}^j (1 - \alpha_\ell).$$

Hence, with Lemma 17,

$$\left| \sum_{i=1}^{n-1} (Q_i - 1) \right| \leq \frac{C_Q}{c_0}.$$

Hence, we can conclude that

$$|T_2| \leq \left(\frac{2C_Q}{c_0} + 1 \right),$$

and proceed with T_1 . Here we notice that, applying (46),

$$Q_i - 1 = S_i - \prod_{\ell=i}^{n-1} (1 - \alpha_\ell), \quad S_i = \sum_{j=i+1}^{n-1} (\alpha_i - \alpha_j) \prod_{\ell=i+1}^{j-1} (1 - \alpha_\ell).$$

Thus, the term T_1 can be represented as

$$T_1 = \sum_{j=1}^{n-1} S_j^2 - 2\sum_{j=1}^{n-1} S_j \prod_{\ell=j}^{n-1} (1 - \alpha_\ell) + \sum_{j=1}^{n-1} \prod_{\ell=j}^{n-1} (1 - \alpha_\ell)^2. \quad (50)$$

Due to item (a) from Lemma 11, it holds that $|S_j| \leq C_S/(j+1)^{1-\gamma}$. Hence, similarly to the proof of Lemma 1 we can show that

$$\frac{1}{n} \left| \sum_{j=1}^{n-1} S_j^2 \right| \leq C_S^2 n^{2(\gamma-1)} / (2\gamma - 1),$$

and

$$\frac{1}{n} \left| \sum_{j=1}^{n-1} S_j \prod_{\ell=j}^{n-1} (1 - \alpha_\ell) \right| \leq C_S n^{2(\gamma-1)} / c_0.$$

Now, we proceed with the last term in (50), and provide a lower bound on the last remaining component of T_1 in (50), that is,

$$T_3 = \sum_{j=1}^{n-1} \prod_{\ell=j}^{n-1} (1 - \alpha_\ell)^2.$$

Since $\alpha_j = \frac{c_0}{(1+j)^\gamma}$, we get, using an elementary inequality $1 - x \geq \exp\{-2x\}$, valid for $0 \leq x \leq 1/2$, we get that

$$\begin{aligned} \sum_{j=1}^{n-1} \prod_{\ell=j}^{n-1} (1 - \alpha_\ell)^2 &\geq \sum_{j=1}^{n-1} \exp\left\{-\sum_{\ell=j}^{n-1} \frac{4c_0}{(1+\ell)^\gamma}\right\} \\ &\geq \sum_{j=1}^{n-1} \exp\left\{-\frac{4c_0}{1-\gamma}(n^{1-\gamma} - j^{1-\gamma})\right\} \\ &= \exp\left\{-\frac{4c_0}{1-\gamma}n^{1-\gamma}\right\} \sum_{j=1}^{n-1} \exp\left\{\frac{4c_0}{1-\gamma}j^{1-\gamma}\right\} \end{aligned}$$

Now we get that

$$\begin{aligned} \sum_{j=1}^{n-1} \exp\left\{\frac{4c_0}{1-\gamma}j^{1-\gamma}\right\} &\geq \int_0^{n-1} \exp\left\{\frac{4c_0}{1-\gamma}y^{1-\gamma}\right\} dy \\ &= (n-1) \int_0^1 \exp\left\{\frac{4c_0}{1-\gamma}((n-1)z)^{1-\gamma}\right\} dz. \end{aligned}$$

Now we proceed with Laplace approximation (see e.g. (Fedoryuk, 1977) or (Olver, 1997)) for the inner integral:

$$\int_0^1 \exp\left\{\frac{4c_0}{1-\gamma}((n-1)z)^{1-\gamma}\right\} dz = \exp\left\{\frac{4c_0}{1-\gamma}(n-1)^{1-\gamma}\right\} \frac{(n-1)^{\gamma-1}}{4c_0} [1 + \mathcal{O}(n^{\gamma-1})]$$

Since $n^{1-\gamma} - (n-1)^{1-\gamma} \leq 1$ and $\frac{n-1}{n} \geq 1/2$ for $n \geq 2$, we get

$$\frac{1}{n} \sum_{j=1}^{n-1} \prod_{\ell=j}^{n-1} (1 - \alpha_\ell)^2 \geq \frac{1}{4c_0} \exp\left\{-\frac{8c_0}{1-\gamma}\right\} \frac{1}{(n-1)^{1-\gamma}} + \mathcal{O}(n^{2(\gamma-1)}).$$

Hence, we conclude that for n large enough,

$$|\sigma_{n,\gamma}^2 - 1| > \frac{C_1(\gamma, c_0)}{n^{1-\gamma}},$$

and the statement follows. To prove the second part, it remains to apply the lower bound on the total variation distance between Gaussian random vectors given in (Devroye et al., 2018, Theorem 1.1).

E.1 Numerical demonstration

In order to illustrate numerically the tightness of bounds provided in Theorem 5, we consider the following simple experiment. We consider the statistics

$$|\sigma_{n,\gamma}^2 - 1| \cdot n^{1-\gamma}, \quad n \in \{2^{10}, \dots, 2^{27}\}.$$

We illustrate numerically the tightness of our bound in the Figure 2 below by calculating

$$n^{1-\gamma} \cdot |\sigma_{n,\gamma}^2 - 1|$$

for different values of $\gamma \in \{0.5, \dots, 0.9\}$ and n . Here we fix the values of parameter $k_0 = 1$ and $c_0 = 1$. Code to reproduce the plot is provided in https://anonymous.4open.science/r/gaussian_approximation_sgd-DBDD/.

F Results for the last iterate

F.1 Last iterate bound

To prove Theorem 2, we need a bound on the $2p$ -th moment for $p \geq 1$ of the last iterate. Our approach is based on induction: first, we establish the result for $p = 2$, and then we show how to proceed from $2(p-1)$ to $2p$.

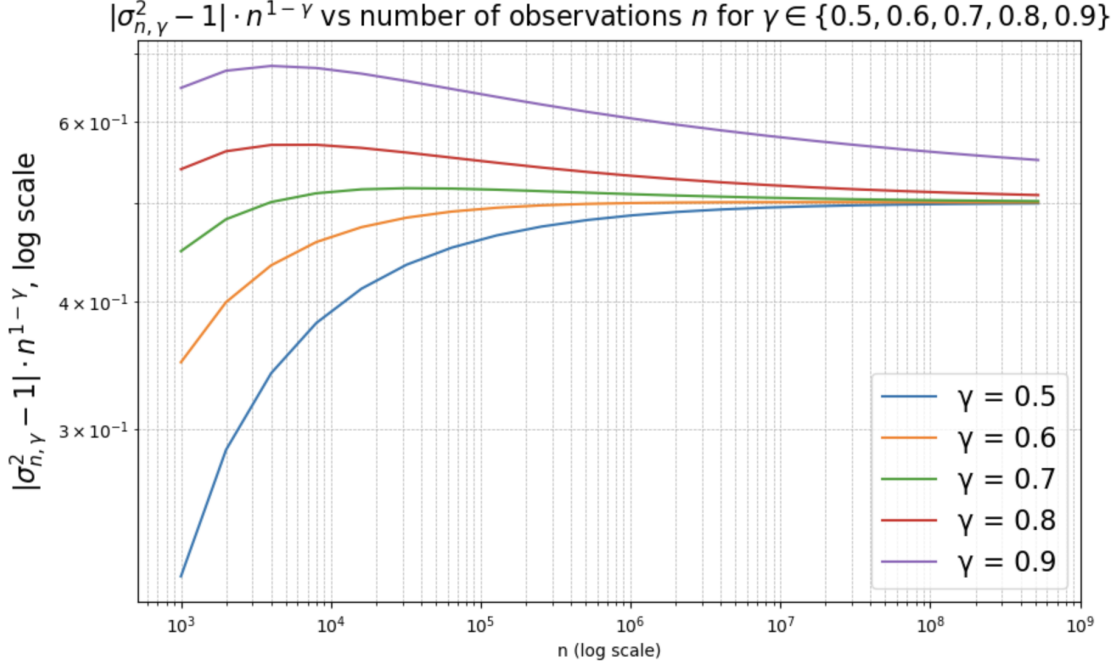


Figure 2: Numerical verification of the lower bound given in Theorem 5

Lemma 12. *Assume A1, A3, A7(2), and A8. Then for any $k \in \mathbb{N}$ it holds that*

$$\mathbb{E}[\|\theta_k - \theta^*\|^2] \leq C_1 \exp\left\{-\frac{\mu c_0}{4}(k + k_0)^{1-\gamma}\right\} [\|\theta_0 - \theta^*\|^2 + \sigma_2^2] + C_2 \sigma_2^2 \alpha_k,$$

where σ_2^2 is defined in A7(2), and the constants C_1 and C_2 are given by

$$C_1 = \exp\left\{\frac{3\mu c_0}{4(1-\gamma)} k_0^{1-\gamma}\right\} \left((1 + L_2^{-2}) \exp\left\{\frac{6c_0^2 L_2^2}{2\gamma - 1}\right\} + \frac{2c_0^2}{2\gamma - 1} \right),$$

$$C_2 = \frac{2^{1+\gamma}}{\mu}.$$

Proof. From (2) and A7 it follows that

$$\|\theta_k - \theta^*\|^2 = \|\theta_{k-1} - \theta^*\|^2 - 2\alpha_k \langle \theta_{k-1} - \theta^*, \nabla f(\theta_{k-1}) + \zeta_k \rangle + \alpha_k^2 \|\nabla f(\theta_{k-1}) + \zeta_k\|^2.$$

Using A1 and A7(2), we obtain

$$2\alpha_k \mathbb{E}[\langle \theta_{k-1} - \theta^*, \nabla f(\theta_{k-1}) + \zeta_k \rangle | \mathcal{F}_{k-1}] = 2\alpha_k \langle \theta_{k-1} - \theta^*, \nabla f(\theta_{k-1}) - \nabla f(\theta^*) \rangle.$$

Using A7(2) and A1, we get

$$\begin{aligned} \mathbb{E}[\|\nabla f(\theta_{k-1}) + \zeta_k\|^2 | \mathcal{F}_{k-1}] &= \|\nabla f(\theta_{k-1}) - \nabla f(\theta^*)\|^2 + \mathbb{E}[\|\eta(\xi_k) + g(\theta_{k-1}, \xi_k)\|^2 | \mathcal{F}_{k-1}] \\ &\leq L_1 \langle \nabla f(\theta_{k-1}) - \nabla f(\theta^*), \theta_{k-1} - \theta^* \rangle + 2L_2^2 \|\theta_{k-1} - \theta^*\|^2 + 2\sigma_2^2. \end{aligned}$$

Combining the above inequalities, we obtain

$$\mathbb{E}[\|\theta_k - \theta^*\|^2] \leq (1 - \mu\alpha_k(2 - \alpha_k L_1) + 2\alpha_k^2 L_2^2) \mathbb{E}[\|\theta_{k-1} - \theta^*\|^2] + 2\alpha_k^2 \sigma_2^2. \quad (51)$$

By applying the recurrence (51), we obtain that

$$\mathbb{E}[\|\theta_k - \theta^*\|^2] \leq A_{1,k} \|\theta_0 - \theta^*\|^2 + 2\sigma_2^2 A_{2,k},$$

where we have set

$$\begin{aligned} A_{1,k} &= \prod_{i=1}^k (1 - (3/2)\alpha_i\mu + 2\alpha_i^2L_2^2), \\ A_{2,k} &= \sum_{i=1}^k \prod_{j=i+1}^k (1 - (3/2)\alpha_j\mu + 2\alpha_j^2L_2^2)\alpha_i^2. \end{aligned} \tag{52}$$

Using the elementary bound $1 + t \leq e^t$ for any $t \in \mathbb{R}$, we get

$$A_{1,k} \leq \exp\left\{- (3/2)\mu \sum_{i=1}^k \alpha_i\right\} \exp\left\{2L_2^2 \sum_{i=1}^k \alpha_i^2\right\}.$$

Using Lemma 15, we obtain

$$A_{1,k} \leq c_1 \exp\left\{-\frac{3\mu c_0}{4(1-\gamma)}(k+k_0)^{1-\gamma}\right\},$$

where we have set

$$c_1 = \exp\left\{\frac{2c_0^2L_2^2}{2\gamma-1} + \frac{3\mu c_0}{4(1-\gamma)}k_0^{1-\gamma}\right\}. \tag{53}$$

Now we estimate $A_{2,k}$. Let k_1 be the largest index k such that $4\alpha_k^2L_2^2 \geq \alpha_k\mu$. Then, for $i > k_1$, we have that

$$1 - (3/2)\alpha_i\mu + 2\alpha_i^2L_2^2 \leq 1 - \alpha_i\mu.$$

Thus, using the definition of $A_{2,k}$ in (52), we obtain that

$$A_{2,k} \leq \sum_{i=1}^k \alpha_i^2 \prod_{j=i+1}^k (1 - \alpha_j\mu) + \sum_{i=1}^{k_1} \alpha_i^2 \left\{ \prod_{j=i+1}^{k_1} (1 + 2\alpha_j^2L_2^2) \right\} \left\{ \prod_{j=k_1+1}^k (1 - \alpha_j\mu) \right\}.$$

Note that

$$\begin{aligned} \sum_{i=1}^{k_1} \alpha_i^2 \prod_{j=i+1}^{k_1} (1 + 2\alpha_j^2L_2^2) &= \frac{1}{2L_2^2} \sum_{i=1}^{k_1} \left(\prod_{j=i}^{k_1} (1 + 2\alpha_j^2L_2^2) - \prod_{j=i+1}^{k_1} (1 + 2\alpha_j^2L_2^2) \right) \\ &\leq \frac{1}{2L_2^2} \prod_{j=1}^{k_1} (1 + 2\alpha_j^2L_2^2) \leq \frac{1}{2L_2^2} \exp\left\{2L_2^2 \sum_{j=1}^{k_1} \alpha_j^2\right\}. \end{aligned}$$

Note, that for $k \leq k_1$, $\alpha_k \geq \mu/(4L_2^2)$, hence, we have

$$\prod_{j=k_1+1}^k (1 - \alpha_j\mu) \leq \exp\left\{-\mu \sum_{i=1}^k \alpha_i\right\} \exp\left\{\mu \sum_{i=1}^{k_1} \alpha_i\right\} \leq \exp\left\{-\mu \sum_{i=1}^k \alpha_i\right\} \exp\left\{4L_2^2 \sum_{i=1}^{k_1} \alpha_i^2\right\}.$$

Moreover, for any $m \in \{1, \dots, k\}$, we obtain

$$\begin{aligned} \sum_{i=1}^k \alpha_i^2 \prod_{j=i+1}^k (1 - \alpha_j\mu) &= \sum_{i=1}^m \prod_{j=i+1}^k (1 - \alpha_j\mu)\alpha_i^2 + \sum_{i=m+1}^k \prod_{j=i+1}^k (1 - \alpha_j\mu)\alpha_i^2 \\ &\leq \prod_{j=m+1}^k (1 - \alpha_j\mu) \sum_{i=1}^m \alpha_i^2 + \alpha_m \sum_{i=m+1}^k \prod_{j=i+1}^k (1 - \alpha_j\mu)\alpha_i \\ &\leq \exp\left\{-\mu \sum_{j=m+1}^k \alpha_j\right\} \sum_{i=1}^m \alpha_i^2 + \frac{\alpha_m}{\mu} \sum_{i=m+1}^k \left(\prod_{j=i+1}^k (1 - \alpha_j\mu) - \prod_{j=i}^k (1 - \alpha_j\mu) \right) \\ &\leq \exp\left\{-\mu \sum_{j=m+1}^k \alpha_j\right\} \sum_{i=1}^m \alpha_i^2 + \frac{\alpha_m}{\mu} \left(1 - \prod_{j=m+1}^k (1 - \alpha_j\mu) \right) \\ &\leq \exp\left\{-\mu \sum_{j=m+1}^k \alpha_j\right\} \sum_{i=1}^m \alpha_i^2 + \frac{\alpha_m}{\mu}. \end{aligned}$$

Thus, setting $m = \lfloor k/2 \rfloor$, and using the definition of $A_{2,k}$ in (52), we obtain that

$$A_{2,k} \leq \exp\left\{\frac{-\mu c_0}{2(1-\gamma)}((k+k_0)^{1-\gamma} - (\lfloor k/2 \rfloor + k_0)^{1-\gamma})\right\} \frac{c_0^2}{2\gamma-1} + \frac{c_0}{\mu(k_0 + \lfloor k/2 \rfloor)^\gamma} + c_2 \exp\left\{-\frac{\mu c_0}{2(1-\gamma)}(k+k_0)^{1-\gamma}\right\},$$

where we have set

$$c_2 = \frac{1}{2L_2^2} \exp\left\{\frac{6c_0^2 L_2^2}{2\gamma-1} + \frac{\mu c_0}{2(1-\gamma)} k_0^{1-\gamma}\right\}. \quad (54)$$

Using that $\lfloor k/2 \rfloor \leq k/2$ together with the elementary inequality

$$\frac{x^\beta}{\beta} - \frac{(x/2)^\beta}{\beta} \geq \frac{x^\beta}{2},$$

which is valid for $\beta \in (0, 1]$, and $\frac{c_0}{\mu(k_0 + \lfloor k/2 \rfloor)^\gamma} \leq \frac{2^\gamma c_0}{\mu(k+k_0)^\gamma}$, we obtain that

$$A_{2,k} \leq \exp\left\{-\frac{\mu c_0}{4}(k+k_0)^{1-\gamma}\right\} \exp\left\{\frac{\mu c_0}{2(1-\gamma)} k_0^{1-\gamma}\right\} \frac{c_0^2}{2\gamma-1} + \frac{2^\gamma c_0}{\mu(k+k_0)^\gamma} + c_2 \exp\left\{\frac{-\mu c_0}{2(1-\gamma)}(k+k_0)^{1-\gamma}\right\}.$$

Combining the bounds for $A_{1,k}$ and $A_{2,k}$, we obtain that

$$\begin{aligned} \mathbb{E}[\|\theta_k - \theta^*\|^2] &\leq c_1 \exp\left\{-\frac{\mu c_0}{(1-\gamma)}(k+k_0)^{1-\gamma}\right\} \|\theta_0 - \theta^*\|^2 \\ &\quad + \exp\left\{-\frac{\mu c_0}{4}(k+k_0)^{1-\gamma}\right\} \frac{2c_0^2 \sigma_2^2}{2\gamma-1} \exp\left\{\frac{\mu c_0}{2(1-\gamma)} k_0^{1-\gamma}\right\} + \frac{2^{1+\gamma} c_0 \sigma_2^2}{\mu(k+k_0)^\gamma} \\ &\quad + 2c_2 \sigma_2^2 \exp\left\{\frac{-\mu c_0}{2(1-\gamma)}(k+k_0)^{1-\gamma}\right\} \\ &\leq C_1 \exp\left\{-\frac{\mu c_0}{4}(k+k_0)^{1-\gamma}\right\} [\|\theta_0 - \theta^*\|^2 + \sigma_2^2] + C_2 \alpha_k, \end{aligned}$$

where we have set constants C_1 and C_2 using the definitions of c_1 and c_2 from (53) and (54). \square

The first term in the bound from Lemma 12 decays exponentially with k , which implies that $\mathbb{E}^{1/2}[\|\theta_k - \theta^*\|^2] \lesssim \alpha_k$. Below, we state this result with an explicit constant, as this specific form of the bound for the last iteration is needed for the induction step.

Corollary 2. *Under the assumptions of Lemma 12, it holds that*

$$\mathbb{E}[\|\theta_k - \theta^*\|^2] \leq D_1 (\|\theta_0 - \theta^*\|^2 + \sigma_2^2) \alpha_k,$$

where

$$D_1 = C_1(1/c_0 + C_2) \left(\frac{4\gamma}{(1-\gamma)\mu c_0 e}\right)^{\gamma/(1-\gamma)}.$$

Proof. Define $C_3 = (\frac{4\gamma}{(1-\gamma)\mu c_0 e})^{\gamma/(1-\gamma)} > 1$, then $\exp\{-\mu c_0(k+k_0)^{1-\gamma}/4\} \leq C_3(k+k_0)^{-\gamma}$, and the statement follows. \square

Now we provide bound for p-moment of last iterate.

Proposition 3. *Assume A1, A3, A7(2p), and A8. Then for any $k \in \mathbb{N}$ it holds that*

$$\mathbb{E}[\|\theta_k - \theta^*\|^{2p}] \leq C_{2p,1} \exp\left\{-\frac{p\mu c_0}{4}(k+k_0)^{1-\gamma}\right\} (\|\theta_0 - \theta^*\|^{2p} + \sigma_{2p}^{2p}) + C_{2p,2} \sigma_{2p}^{2p} \alpha_k^p,$$

where

$$C_{2p,1} = 2^{2p-1}(D_{2(p-1)}C_4^p c_0^p + 1)c_4 ,$$

$$C_{2p,2} = 2^{2p-1}D_{2(p-1)}C_4^p \frac{2^{1+\gamma p}}{\mu p c_0} ,$$

constants $D_{2(p-1)}$ are defined in (61), and

$$C_4^p = (4c_0^{1/2}2^{\gamma/2} + 2^\gamma + 4c_0)^p$$

$$c_4 = \left(\exp \left\{ \exp \left\{ 5pc_0(L_1 + L_2) \right\} \frac{4p^2(L_1 + L_2)^2}{2\gamma - 1} \right\} + 1 \right) \exp \left\{ \frac{p\mu c_0}{1 - \gamma} k_0^{1-\gamma} \right\} \frac{1}{\gamma(p+1) - 1}$$

Proof. We prove the statement by induction in p . We first assume that $\theta_0 = \theta^*$ and then provide a result for arbitrary initial condition. The result for $p = 1$ is provided in Corollary 2. Assume that for any $t \leq p - 1$ and all $k \in \mathbb{N}$ we proved that

$$\mathbb{E}[\|\theta_k - \theta^*\|^{2t}] \leq D_{2t}\sigma_{2t}^{2t}\alpha_k^t , \quad (55)$$

and the sequence of constants $\{D_{2t}\}$ is non-decreasing in t . Inequality (55) implies that, since $\sigma_{2t} \leq \sigma_{2p}$ for $t \leq p - 1$,

$$\mathbb{E}[\|\theta_k - \theta^*\|^{2t}] \leq D_{2t}\sigma_{2p}^{2t}\alpha_k^t .$$

For any $k \in \mathbb{N}$ we denote $\delta_k = \|\theta_k - \theta^*\|$. Using (2), we get

$$\begin{aligned} \delta_k^{2p} &= (\delta_{k-1}^2 - 2\alpha_k \langle \theta_{k-1} - \theta^*, \nabla f(\theta_{k-1}) + \zeta_k \rangle + \alpha_k^2 \|\nabla f(\theta_{k-1}) + \zeta_k\|^2)^p \\ &= \sum_{\substack{i+j+l=p \\ i,j,l \in \{0, \dots, p\}}} \frac{p!}{i!j!l!} \delta_{k-1}^{2i} (-2\alpha_k \langle \theta_{k-1} - \theta^*, \nabla f(\theta_{k-1}) + \zeta_k \rangle)^j \alpha_k^{2l} \|\nabla f(\theta_{k-1}) + \zeta_k\|^{2l} . \end{aligned}$$

Now we bound each term in the sum above.

1. First, for $i = p, j = 0, l = 0$, the corresponding term in the sum equals δ_{k-1}^{2p} .
2. Second, for $i = p - 1, j = 1, l = 0$, we obtain, applying A1, that

$$\begin{aligned} 2p\alpha_k \mathbb{E}[\langle \theta_{k-1} - \theta^*, \nabla f(\theta_{k-1}) + \zeta_k \rangle \delta_{k-1}^{2(p-1)} | \mathcal{F}_{k-1}] &= 2p\alpha_k \langle \theta_{k-1} - \theta^*, \nabla f(\theta_{k-1}) - \nabla f(\theta^*) \rangle \delta_{k-1}^{2(p-1)} \\ &\geq 2p\mu\alpha_k \delta_{k-1}^{2p} . \end{aligned}$$

3. Third, for $l \geq 1$ or $j \geq 2$ (that is, $2l + j \geq 2$), we use Cauchy-Schwartz inequality

$$|\langle \theta_{k-1} - \theta^*, \nabla f(\theta_{k-1}) + \zeta_k \rangle^j| \leq \|\theta_{k-1} - \theta^*\|^j \|\nabla f(\theta_{k-1}) + \zeta_k\|^j ,$$

moreover, applying A1 and A7(2p) together with the Lyapunov inequality, we get

$$\begin{aligned} \mathbb{E}[\|\nabla f(\theta_{k-1}) + \zeta_k\|^{2l+j} | \mathcal{F}_{k-1}] &= \mathbb{E}[\|\nabla f(\theta_{k-1}) + g(\theta_{k-1}, \xi_k) + \eta(\xi_k)\|^{2l+j} | \mathcal{F}_{k-1}] \\ &\leq 2^{2l+j-1}((L_1 + L_2)^{2l+j} \delta_{k-1}^{2l+j} + \sigma_{2p}^{2l+j}) . \end{aligned}$$

Combining inequalities above, we get

$$\begin{aligned} \mathbb{E}[\delta_k^{2p} | \mathcal{F}_{k-1}] &\leq \left(1 - 2p\mu\alpha_k + \sum_{\substack{i+j+l=p \\ i,j,l \in \{0, \dots, p\} \\ j+2l \geq 2}} \frac{p!}{i!j!l!} \alpha_k^{j+2l} 2^{2l+2j-1} (L_1 + L_2)^{2l+j} \right) \delta_{k-1}^{2p} \\ &\quad + \sum_{\substack{i+j+l=p \\ i,j,l \in \{0, \dots, p\} \\ j+2l \geq 2}} \frac{p!}{i!j!l!} \delta_{k-1}^{2i+j} \alpha_k^{j+2l} 2^{2l+2j-1} \sigma_{2p}^{2l+j} . \end{aligned}$$

Consider the first term above, and note that

$$\begin{aligned}
& \sum_{\substack{i+j+l=p; \\ i,j,l \in \{0, \dots, p\}; \\ j+2l \geq 2}} \frac{p!}{i!j!l!} \alpha_k^{j+2l} 2^{2l+2j-1} (L_1 + L_2)^{2l+j} \\
& \leq 2\alpha_k^2 (L_1 + L_2)^2 \left(\sum_{\substack{i+j+l=p; \\ i,j,l \in \{0, \dots, p\}; \\ l \geq 1}} \frac{p!}{i!j!l!} (4\alpha_k (L_1 + L_2))^j (4\alpha_k^2 (L_1 + L_2)^2)^{l-1} + \sum_{\substack{i+j+l=p; \\ i,j,l \in \{0, \dots, p\}; \\ l=0; j \geq 2}} \frac{p!}{i!j!} (4\alpha_k (L_1 + L_2))^{j-2} \right) \\
& \leq 2p^2 \alpha_k^2 (L_1 + L_2)^2 (1 + 5\alpha_k (L_1 + L_2))^p.
\end{aligned} \tag{56}$$

Hence,

$$\mathbb{E}[\delta_k^{2p}] \leq (1 - 2p\mu\alpha_k + 2p^2\alpha_k^2(L_1 + L_2)^2(1 + 5\alpha_k(L_1 + L_2))^p) \delta_{k-1}^{2p} + T_1,$$

where we have defined

$$T_1 = \sum_{\substack{i+j+l=p; \\ i,j,l \in \{0, \dots, p\}; \\ j+2l \geq 2; \\ i+j=p}} \frac{p!}{i!j!l!} \mathbb{E}[\delta_{k-1}^{2i+j}] \alpha_k^{j+2l} 2^{2l+2j-1} \sigma_{2p}^{2l+j}.$$

For the last term we apply Hölder's inequality together with induction assumption (55) and $(k + k_0 - 1)^{-\gamma} \leq 2^\gamma(k + k_0)^{-\gamma}$ and obtain

$$\begin{aligned}
T_1 & \leq \sum_{\substack{i+j+l=p; \\ i,j,l \in \{0, \dots, p\}; \\ j+2l \geq 2}} \frac{p!}{i!j!l!} \mathbb{E}^{1/2}[\delta_{k-1}^{2(i+\lfloor j/2 \rfloor)}] \mathbb{E}^{1/2}[\delta_{k-1}^{2(i+\lceil j/2 \rceil)}] \alpha_k^{j+2l} 2^{2l+2j-1} \sigma_{2p}^{2l+j} \\
& \leq D_{2(p-1)} \frac{(4c_0^{1/2} 2^{\gamma/2} + 2^\gamma + 4c_0)^p}{2} c_0^p \sigma_{2p}^{2p} (k + k_0)^{-\gamma(p+1)}.
\end{aligned}$$

Hence, combining the above bounds, we obtain that

$$\mathbb{E}[\delta_k^{2p}] \leq (1 - 2p\mu\alpha_k + 16\alpha_k^2(L_1 + L_2)^2 3^p) \mathbb{E}[\delta_{k-1}^{2p}] + D_{2(p-1)} C_4^p c_0^p \sigma_{2p}^{2p} k^{-\gamma(p+1)}, \tag{57}$$

where we have defined

$$C_4^p = (4c_0^{1/2} 2^{\gamma/2} + 2^\gamma + 4c_0)^p.$$

Note that

$$1 - 2p\mu\alpha_k + 2p^2\alpha_k^2(L_1 + L_2)^2(1 + 5\alpha_k(L_1 + L_2))^p > 1 - 2p\mu\alpha_k + \alpha_k^2\mu^2p^2 \geq 0.$$

Unrolling the recurrence (57), we get

$$\mathbb{E}[\delta_k^{2p}] \leq A'_{2,k} D_{2(p-1)} C_4^p c_0^p \sigma_{2p}^{2p},$$

where we have set

$$A'_{2,k} = \sum_{t=1}^k \prod_{i=t+1}^k (1 - 2p\mu\alpha_i + 2p^2\alpha_i^2(L_1 + L_2)^2(1 + 5\alpha_i(L_1 + L_2))^p) (t + k_0)^{-\gamma(p+1)}. \tag{58}$$

For simplicity, we define $C_5 = 2p^2(L_1 + L_2)^2$. Let k_1 is the largest k such that $\alpha_k^2 C_5 (1 + 5\alpha_i(L_1 + L_2))^p \geq p\mu\alpha_k$. Then, for $i > k_1$, we have

$$1 - 2p\mu\alpha_i + C_5\alpha_i^2(1 + 5\alpha_i(L_1 + L_2))^p \leq 1 - p\mu\alpha_i.$$

Hence, using the definition of $A'_{2,k}$ in (58), we get

$$\begin{aligned}
 A'_{2,k} &= \sum_{t=k_1+1}^k \prod_{i=t+1}^k \exp\{-p\mu\alpha_i\} (t+k_0)^{-\gamma(p+1)} \\
 &\quad + \prod_{t=k_1+1}^k \exp\{-p\mu\alpha_t\} \sum_{t=1}^{k_1} \prod_{i=t+1}^{k_1} \exp\left\{C_5\alpha_i^2(1+5\alpha_i(L_1+L_2))^p\right\} (t+k_0)^{-\gamma(p+1)} \\
 &\leq \sum_{t=1}^k \prod_{i=t+1}^k \exp\{-p\mu\alpha_i\} (t+k_0)^{-\gamma(p+1)} \\
 &\quad + \prod_{t=1}^k \exp\{-p\mu\alpha_t\} \prod_{t=1}^{k_1} \exp\{p\mu\alpha_t\} \prod_{i=1}^{k_1} \exp\left\{C_5\alpha_i^2(1+5\alpha_i(L_1+L_2))^p\right\} \sum_{t=1}^{k_1} (t+k_0)^{-\gamma(p+1)} \\
 &\leq \sum_{t=1}^k \prod_{i=t+1}^k \exp\{-p\mu\alpha_i\} (t+k_0)^{-\gamma(p+1)} \\
 &\quad + \prod_{i=1}^{k_1} \exp\left\{2C_5\alpha_i^2(1+5\alpha_i(L_1+L_2))^p\right\} \prod_{t=1}^k \exp\{-p\mu\alpha_t\} \sum_{t=1}^{k_1} (t+k_0)^{-\gamma(p+1)}
 \end{aligned}$$

For any $m \in \{1, \dots, k\}$ we have

$$\begin{aligned}
 &\sum_{t=1}^k \prod_{i=t+1}^k \exp\{-p\mu\alpha_i\} (t+k_0)^{-\gamma(p+1)} \\
 &= \sum_{t=1}^m \prod_{i=t+1}^k \exp\{-p\mu\alpha_i\} (t+k_0)^{-\gamma(p+1)} + \sum_{t=m+1}^k \prod_{i=t+1}^k \exp\{-p\mu\alpha_i\} (t+k_0)^{-\gamma(p+1)} \\
 &\leq \prod_{i=m+1}^k \exp\{-p\mu\alpha_i\} \sum_{t=1}^m (t+k_0)^{-\gamma(p+1)} + \sum_{t=m+1}^k \prod_{i=t+1}^k \exp\{-p\mu\alpha_i\} (m+k_0)^{-\gamma p} (t+k_0)^{-\gamma} \\
 &\leq \prod_{i=m+1}^k \exp\{-p\mu\alpha_i\} \sum_{t=1}^k (t+k_0)^{-\gamma(p+1)} + (m+k_0)^{-\gamma p} \sum_{t=1}^k \prod_{i=t+1}^k \exp\{-p\mu\alpha_i\} (t+k_0)^{-\gamma}
 \end{aligned}$$

Applying Lemma 15(b), we have

$$\begin{aligned}
 &\sum_{t=1}^k \prod_{i=t+1}^k \exp\{-p\mu\alpha_i\} t^{-\gamma} \leq \sum_{t=1}^k \exp\left\{\frac{-p\mu c_0}{2(1-\gamma)}((k+k_0)^{1-\gamma} - (t+k_0)^{1-\gamma})\right\} (t+k_0)^{-\gamma} \\
 &\leq \exp\left\{\frac{-p\mu c_0}{2(1-\gamma)}(k+k_0)^{1-\gamma}\right\} \frac{2}{p\mu c_0} \int_0^{\frac{p\mu c_0}{2(1-\gamma)}(k+k_0)^{1-\gamma}} e^u du \leq \frac{2}{p\mu c_0}.
 \end{aligned}$$

Applying Lemma 15(a), we get

$$\sum_{i=1}^k (i+k_0)^{-\gamma(p+1)} \leq \frac{1}{(p+1)\gamma-1},$$

and

$$\sum_{i=1}^k 2C_5\alpha_i^2(1+5\alpha_i(L_1+L_2))^p \leq 2C_5(1+5c_0(L_1+L_2))^p \sum_{i=1}^{+\infty} \alpha_k^2 \leq \exp\left\{5pc_0(L_1+L_2)\right\} \frac{2C_5}{2\gamma-1}$$

Substituting $m = \lfloor k/2 \rfloor$ and applying ((b)), we get

$$\begin{aligned}
 A'_{2,k} &\leq \exp\left\{-\frac{p\mu c_0}{2(1-\gamma)}((k+k_0)^{1-\gamma} - (\lfloor k/2 \rfloor + k_0)^{1-\gamma})\right\} \frac{1}{\gamma(p+1)-1} + \frac{2(\lfloor k/2 \rfloor + k_0)^{-\gamma p}}{p\mu c_0} \\
 &\quad + c_3 \exp\left\{-\frac{p\mu c_0}{2(1-\gamma)}(k+k_0)^{1-\gamma}\right\},
 \end{aligned}$$

where we have set

$$c_3 = \exp\left\{\exp\left\{5pc_0(L_1 + L_2)\right\}\frac{2C_5}{2\gamma - 1} + \frac{p\mu c_0}{2(1 - \gamma)}k_0^{1-\gamma}\right\}\frac{1}{\gamma(p + 1) - 1}.$$

Using that $\lfloor k/2 \rfloor \leq k/2$ together with the elementary inequality

$$\frac{x^\beta}{\beta} - \frac{(x/2)^\beta}{\beta} \geq \frac{x^\beta}{2},$$

which is valid for $\beta \in (0, 1]$, and $\frac{2}{\mu pc_0(\lfloor k/2 \rfloor + k_0)^{\gamma p}} \leq \frac{2^{1+\gamma p}}{\mu pc_0(k + k_0)^{\gamma p}}$, we obtain that

$$\begin{aligned} A'_{2,k} &\leq \exp\left\{-\frac{p\mu c_0}{4}(k + k_0)^{1-\gamma}\right\} \exp\left\{\frac{p\mu c_0}{2(1 - \gamma)}k_0^{1-\gamma}\right\} \frac{1}{\gamma(p + 1) - 1} \\ &\quad + \frac{2^{1+\gamma p}}{\mu pc_0(k + k_0)^{\gamma p}} + c_3 \exp\left\{-\frac{p\mu c_0}{2(1 - \gamma)}(k + k_0)^{1-\gamma}\right\} \\ &\leq c_4 \exp\left\{-\frac{p\mu c_0}{4}(k + k_0)^{1-\gamma}\right\} + c_5(k + k_0)^{-\gamma p}, \end{aligned}$$

where we have set

$$\begin{aligned} c_4 &= \left(\exp\left\{\exp\left\{5pc_0(L_1 + L_2)\right\}\frac{4p^2(L_1 + L_2)^2}{2\gamma - 1}\right\} + 1\right) \exp\left\{\frac{p\mu c_0}{1 - \gamma}k_0^{1-\gamma}\right\} \frac{1}{\gamma(p + 1) - 1} \\ c_5 &= \frac{2^{1+\gamma p}}{\mu pc_0} \end{aligned}$$

Finally, we get

$$\mathbb{E}[\delta_k^{2p}] \leq C'_{2p,1} \exp\left\{-\frac{p\mu c_0}{4}(k + k_0)^{1-\gamma}\right\} \sigma_{2p}^{2p} + C'_{2p,2} \sigma_{2p}^{2p} \alpha_k^p,$$

where

$$\begin{aligned} C'_{2p,1} &= D_{2(p-1)} C_4^p c_0^p c_4 \\ C'_{2p,2} &= D_{2(p-1)} C_4^p c_5. \end{aligned}$$

To provide the result for arbitrary start point $\theta_0 = \theta$ we consider the synchronous coupling construction defined by the recursions

$$\begin{aligned} \theta_k &= \theta_{k-1} - \alpha_k(\nabla f(\theta_{k-1}) + g(\theta_{k-1}, \xi_k) + \eta(\xi_k)), & \theta_0 &= \theta \\ \theta'_k &= \theta'_{k-1} - \alpha_k(\nabla f(\theta'_{k-1}) + g(\theta'_{k-1}, \xi_k) + \eta(\xi_k)), & \theta'_0 &= \theta^* \end{aligned} \tag{59}$$

For any $k \in \mathbb{N}$ we denote $\delta'_k = \|\theta_k - \theta'_k\|$. Using (59) together with A1 and A7(2p), we get

$$\begin{aligned} \delta_k'^{2p} &= (\delta_{k-1}'^2 - 2\alpha_k \langle \theta_{k-1} - \theta'_{k-1}, \nabla f(\theta_{k-1}) - \nabla f(\theta'_{k-1}) + g(\theta_{k-1}, \xi_k) - g(\theta'_{k-1}, \xi_k) \rangle + \alpha_k^2(L_1 + L_2)^2 \delta_{k-1}'^2)^p \\ &\leq \sum_{\substack{i+j+l=p \\ i,j,l \in \{0, \dots, p\}}} \frac{p!}{i!j!l!} \delta_{k-1}'^{2i} (-2\alpha_k \langle \theta_{k-1} - \theta'_{k-1}, \nabla f(\theta_{k-1}) - \nabla f(\theta'_{k-1}) + g(\theta_{k-1}, \xi_k) - g(\theta'_{k-1}, \xi_k) \rangle)^j (\alpha_k(L_1 + L_2) \delta_{k-1}'^2)^{2l} \end{aligned}$$

Now we bound each term in the sum above.

1. First, for $i = p, j = 0, l = 0$, the corresponding term in the sum equals $\delta_{k-1}'^{2p}$.
2. Second, for $i = p - 1, j = 1, l = 0$, we obtain, applying A1, that

$$\begin{aligned} &2p\alpha_k \mathbb{E}[\langle \theta_{k-1} - \theta'_{k-1}, \nabla f(\theta_{k-1}) - \nabla f(\theta'_{k-1}) + g(\theta_{k-1}, \xi_k) - g(\theta'_{k-1}, \xi_k) \rangle \delta_{k-1}'^{2(p-1)} | \mathcal{F}_{k-1}] \\ &= 2p\alpha_k \langle \theta_{k-1} - \theta'_{k-1}, \nabla f(\theta_{k-1}) - \nabla f(\theta'_{k-1}) \rangle \delta_{k-1}'^{2(p-1)} \geq 2p\mu\alpha_k \delta_{k-1}'^{2p}. \end{aligned}$$

3. Third, for $l \geq 1$ or $j \geq 2$ (that is, $2l + j \geq 2$), we use Cauchy-Schwartz inequality together with A7 and A1

$$|(\theta_{k-1} - \theta'_{k-1}, \nabla f(\theta_{k-1}) - \nabla f(\theta'_{k-1}) + g(\theta_{k-1}, \xi_k) - g(\theta'_{k-1}, \xi_k))^j| \leq \|\theta_{k-1} - \theta'_{k-1}\|^{2j} (L_1 + L_2)^j,$$

Combining inequalities above, we obtain

$$\mathbb{E}[\delta_k'^{2p} | \mathcal{F}_{k-1}] \leq (1 - 2p\mu\alpha_k + \sum_{\substack{i+j+l=p; \\ i,j,l \in \{0,\dots,p\}; \\ j+2l \geq 2}} \frac{p!}{i!j!l!}) 2^j \alpha_k^{j+2l} (L_1 + L_2)^{j+2l} \delta_{k-1}'^{2p} \quad (60)$$

Similar to (56), we have

$$\sum_{\substack{i+j+l=p; \\ i,j,l \in \{0,\dots,p\}; \\ j+2l \geq 2}} \frac{p!}{i!j!l!} 2^j \alpha_k^{j+2l} (L_1 + L_2)^{j+2l} \delta_{k-1}'^{2p} \leq \alpha_k^2 p^2 (L_1 + L_2)^2 (1 + 3\alpha_k (L_1 + L_2))^p$$

Enrolling recurrence (60), we get

$$\begin{aligned} \mathbb{E}[\delta_k'^{2p}] &\leq \exp\left\{-2p\mu \sum_{i=1}^k \alpha_i\right\} \exp\left\{p^2 (L_1 + L_2)^2 \sum_{i=1}^k \alpha_i^2 (1 + 3\alpha_i (L_1 + L_2))^p\right\} \|\theta_0 - \theta^*\|^{2p} \\ &\leq c_6 \exp\left\{-\frac{p\mu c_0}{1-\gamma} (k + k_0)^{1-\gamma}\right\} \|\theta_0 - \theta^*\|^{2p}, \end{aligned}$$

where we have set

$$c_6 = \exp\left\{\exp\left\{3pc_0(L_1 + L_2)\right\} \frac{p^2(L_1 + L_2)^2}{2\gamma - 1} + \frac{p\mu c_0}{1-\gamma} k_0^{1-\gamma}\right\}.$$

It remains to note that

$$\begin{aligned} \mathbb{E}[\|\theta_k - \theta^*\|^{2p}] &\leq 2^{2p-1} \mathbb{E}[\|\theta'_k - \theta^*\|^{2p}] + 2^{2p-1} \mathbb{E}[\|\theta_k - \theta'_k\|^{2p}] \\ &\leq C_{2p,1} \exp\left\{-\frac{p\mu c_0}{4} (k + k_0)^{1-\gamma}\right\} (\|\theta_0 - \theta^*\|^{2p} + \sigma_{2p}^{2p}) + C_{2p,2} \sigma_{2p}^{2p} \alpha_k^p. \end{aligned}$$

□

For validity of induction in Proposition 3, we need the following corollary.

Corollary 3. *Under the assumptions of Proposition 3, it holds that*

$$\mathbb{E}[\|\theta_k - \theta^*\|^{2p}] \leq D_{2p} (\|\theta_0 - \theta^*\|^{2p} + \sigma_{2p}^{2p}) \alpha_k^p,$$

where

$$D_{2p} = C_{2p,1} (1/c_0^p + C_{2p,2}) \left(\frac{4\gamma}{(1-\gamma)\mu p c_0 e} \right)^{\gamma p / (1-\gamma)}. \quad (61)$$

Proof. Define $C_5 = \left(\frac{4\gamma}{(1-\gamma)\mu p c_0 e}\right)^{\gamma p / (1-\gamma)} > 1$, then $\exp\{-\mu p c_0 (k + k_0)^{1-\gamma} / 4\} \leq C_5 (k + k_0)^{-p\gamma}$, and the statement follows. □

Corollary 4. *Assume A1, A3, A7(4) and A8. Then for any $k \in \mathbb{N}$ it holds that*

$$\mathbb{E}[\|\theta_k - \theta^*\|^4] \leq C_{4,1} \exp\left\{-\frac{2\mu c_0}{4} k^{1-\gamma}\right\} (\|\theta_0 - \theta^*\|^4 + \sigma_4^4) + C_{4,2} \sigma_4^4 \alpha_k^2,$$

with

$$C_{4,1} = 2^3 \left(C_1 (1/c_0 + C_2) \left(\frac{4\gamma}{(1-\gamma)\mu c_0 e} \right)^{\gamma / (1-\gamma)} (4c_0^{1/2} 2^{\gamma/2} + 2^\gamma + 4c_0)^2 c_0^2 + 1 \right) c_{2,4}$$

and

$$C_{4,2} = 2^3 C_1 (1/c_0 + C_2) \left(\frac{4\gamma}{(1-\gamma)\mu c_0 e} \right)^{\gamma/(1-\gamma)} (4c_0^{1/2} 2^{\gamma/2} + 2^\gamma + 4c_0)^2 c_{2,5}.$$

Here C_1 and C_2 are defined in Lemma 12 and

$$c_{2,4} = \left(\exp \left\{ \exp \left\{ 10c_0(L_1 + L_2) \right\} \frac{16(L_1 + L_2)^2}{2\gamma - 1} \right\} + 1 \right) \exp \left\{ \frac{2\mu c_0}{1-\gamma} k_0^{1-\gamma} \right\} \frac{1}{3\gamma - 1},$$

$$c_{2,5} = \frac{2^{1+2\gamma}}{2\mu c_0}.$$

Proof. The proof follows directly from Proposition 3 and Corollary 2. \square

F.2 High probability bounds on the last iterate

In this section, we establish a high-probability bound for the last iterate, which is instrumental in controlling the non-linear statistic D^b . Our analysis adapts the approach of (Madden et al., 2024, Theorem 9), which relies on the assumption that the noise variables ζ_k are sub-Gaussian. This result, in turn, generalizes the previous results of (Harvey et al., 2019), where the authors assumed that both the additive noise component $\eta(\xi)$ and state-dependent component $g(\theta, \xi)$ are uniformly bounded.

Lemma 13. *Assume A1, A2, A4, A5. Then for any $\delta \in (0, 1)$, it holds with probability at least $1 - \delta$ that for any $k \in \{1, \dots, n\}$,*

$$\|\theta_k^b - \theta^*\|^2 \leq \alpha_k K_1 \log \left(\frac{e}{\delta} \right),$$

where

$$K_1 = \frac{k_0^\gamma}{c_0} \|\theta_0 - \theta^*\|^2 + \frac{16W_{\max}^2 C_\xi^2}{\mu W_{\min}} (2d + 1) \quad (62)$$

Proof. Unrolling the recurrence (6), we have

$$\begin{aligned} \|\theta_k^b - \theta^*\|^2 &= \|\theta_{k-1}^b - \theta^*\|^2 - 2\alpha_k w_k \langle F(\theta_{k-1}^b, \xi_k), \theta_{k-1}^b - \theta^* \rangle + \alpha_k^2 w_k^2 \|F(\theta_{k-1}^b, \xi_k)\|^2 \\ &\leq \|\theta_{k-1}^b - \theta^*\|^2 - 2\alpha_k w_k \langle \nabla f(\theta_{k-1}^b), \theta_{k-1}^b - \theta^* \rangle - 2\alpha_k w_k \langle g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k), \theta_{k-1}^b - \theta^* \rangle \\ &\quad + \alpha_k^2 w_k^2 \|\nabla f(\theta_{k-1}^b) + g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k)\|^2. \end{aligned}$$

Using A1, A2, and A4, we obtain

$$\begin{aligned} \|\theta_k^b - \theta^*\|^2 &\leq (1 - 2\alpha_k \mu W_{\min} + 2\alpha_k^2 L_1^2 W_{\max}^2) \|\theta_{k-1}^b - \theta^*\|^2 \\ &\quad - 2\alpha_k w_k \langle g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k), \theta_{k-1}^b - \theta^* \rangle + 2\alpha_k^2 w_k^2 \|g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k)\|^2 \\ &\leq (1 - 2\alpha_k \mu W_{\min} + 2\alpha_k^2 L_1^2 W_{\max}^2) \|\theta_{k-1}^b - \theta^*\|^2 \\ &\quad - 2\alpha_k w_k \langle g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k), \theta_{k-1}^b - \theta^* \rangle + 2\alpha_k^2 W_{\max}^2 \|g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k)\|^2 \end{aligned}$$

Using A5, we have

$$\begin{aligned} \|\theta_k^b - \theta^*\|^2 &\leq (1 - \mu \alpha_k W_{\min}) \|\theta_{k-1}^b - \theta^*\|^2 - 2\alpha_k w_k \langle g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k), \theta_{k-1}^b - \theta^* \rangle \\ &\quad + 2\alpha_k^2 W_{\max}^2 \|g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k)\|^2. \quad (63) \end{aligned}$$

Now we introduce the quantities

$$X_k = \alpha_k^{-1} \|\theta_k^b - \theta^*\|^2, \quad Y_{k-1} = -2w_k \langle g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k), \theta_{k-1}^b - \theta^* \rangle, \quad Z_{k-1} = 2\alpha_k W_{\max}^2 \|g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k)\|^2.$$

Using (63), we obtain

$$X_k \leq \alpha_k^{-1} \alpha_{k-1} (1 - \mu W_{\min} \alpha_k) X_{k-1} + Y_{k-1} + Z_{k-1}. \quad (64)$$

Note that

$$\begin{aligned} \frac{\alpha_{k-1}}{\alpha_k} (1 - \mu W_{\min} \alpha_k) &= \left(\frac{k_0 + k}{k_0 + k - 1} \right)^\gamma - \frac{\mu W_{\min} c_0}{(k_0 + k - 1)^\gamma} \\ &\leq 1 + \frac{c_0(\gamma/c_0)}{k_0 + k - 1} - \frac{\mu W_{\min} c_0}{(k_0 + k - 1)^\gamma} \\ &= 1 - \alpha_{k-1} \left(\mu W_{\min} - \frac{(\gamma/c_0)}{(k_0 + k - 1)^{1-\gamma}} \right). \end{aligned}$$

Since $k_0 \geq \left(\frac{2\gamma}{c_0 \mu W_{\min}} \right)^{1/(1-\gamma)}$, recurrence (64) and the above identities yield

$$X_k \leq (1 - \mu W_{\min} \alpha_{k-1}/2) X_{k-1} + Y_{k-1} + Z_{k-1}.$$

Using A2 and A4, we have for any $\lambda \in \mathbb{R}$

$$\mathbb{E}[\exp\{\lambda Y_{k-1}\} | \tilde{\mathcal{F}}_{k-1}] = \mathbb{E}[\exp\{-2\lambda w_k \langle \theta_{k-1}^b - \theta^*, g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k) \rangle\} | \tilde{\mathcal{F}}_{k-1}] \leq \exp\{\lambda^2 \beta_{k-1}^2 X_{k-1}/2\},$$

where we set $\beta_{k-1}^2 = 4W_{\max}^2 C_\xi^2 \alpha_{k-1}$. Using A2 and Lemma 19, we get for any $\lambda \in [0, r_{k-1}^{-1}]$

$$\mathbb{E}[\exp\{\lambda Z_{k-1}\} | \tilde{\mathcal{F}}_{k-1}] = \mathbb{E}[\exp\{2\lambda \alpha_k W_{\max}^2 \|g(\theta_{k-1}^b, \xi_k) + \eta(\xi_k)\|^2\} | \tilde{\mathcal{F}}_{k-1}] \leq \exp\{\lambda r_{k-1}\},$$

where $r_{k-1} = 8d\alpha_k W_{\max}^2 C_\xi^2$ and $\tilde{\mathcal{F}}_{k-1}$ is defined in (19). Then, applying (Madden et al., 2024, Theorem 9), for any $k \geq 0$ it holds with probability at least $1 - \delta$ for $\delta \in (0, 1)$ that

$$X_k \leq P_k \log\left(\frac{e}{\delta}\right),$$

where $\{P_\ell\}_{\ell \in \mathbb{N}}$ is any sequence of positive real numbers, satisfying

$$P_{i+1}^2 \geq \left(1 - \frac{\mu W_{\min}}{2} \alpha_i\right) P_i + 2r_i P_{i+1} + \beta_i^2 P_i, \quad P_0 \geq X_0. \quad (65)$$

Note in particular, that the sequence $\{P_\ell\}_{\ell \in \mathbb{N}}$ given by the recurrence

$$P_{i+1} = \left(1 - \frac{\mu W_{\min}}{2} \alpha_i\right) P_i + \tau_i, \quad \text{where } \tau_i = 2r_i + \frac{\beta_i^2}{1 - \frac{\mu W_{\min}}{2} \alpha_i}, \text{ and } P_0 = X_0,$$

satisfies (65). Hence, unraveling the recursion, we have

$$P_{k+1} = \prod_{i=0}^k \left(1 - \frac{\mu W_{\min}}{2} \alpha_i\right) P_0 + \sum_{i=0}^k \prod_{j=i+1}^k \left(1 - \frac{\mu W_{\min}}{2} \alpha_j\right) \tau_i.$$

Using Lemma 16 and $\alpha_{i+1} \leq \alpha_i$, we have

$$P_{k+1} \leq \frac{k_0^\gamma}{c_0} \|\theta_0 - \theta^*\|^2 + \frac{32dW_{\max}^2 C_\xi^2}{\mu W_{\min}} + \frac{8W_{\max}^2 C_\xi^2}{\left(1 - \frac{\mu W_{\min} c_0}{2k_0^\gamma}\right) \mu W_{\min}}.$$

To complete the proof, it remains to apply the bound on k_0 given by A5. □

Corollary 5. *Under the assumptions of Lemma 13 for any $k \in \{1, \dots, n\}$ and any $p \geq 2$ it holds*

$$\mathbb{E}^{2/p}[\|\theta_k^b - \theta^*\|^p] \leq p\alpha_k (e)^{2/p} K_1/2,$$

where K_1 is defined in (62).

Proof. Note that from Lemma 13 for $\forall k \in \{1, \dots, n\}$ and for any $t \geq 0$ it holds

$$\mathbb{P}[\|\theta_k^b - \theta^*\|^2 \geq t] \leq f(t),$$

where

$$f(t) = e \exp\left\{-\frac{t}{K_1 \alpha_k}\right\}.$$

Then, we have

$$\begin{aligned} \mathbb{E}[\|\theta_k^b - \theta^*\|^p] &= \int_0^{+\infty} \mathbb{P}[\|\theta_k^b - \theta^*\|^p > u] du \leq \int_0^{+\infty} e \exp\left\{-\frac{u^{2/p}}{K_1 \alpha_k}\right\} du \\ &= e(p/2) \left(K_1 \alpha_k\right)^{p/2} \int_0^{+\infty} e^{-x} x^{p/2-1} dx \leq e \left((p/2) K_1 \alpha_k\right)^{p/2}, \end{aligned}$$

where in the last inequality we use that $\Gamma(p/2) \leq (p/2)^{p/2-1}$ (see (Anderson and Qiu, 1997, Theorem 1.5)). \square

Lemma 14. *Assume A1, A2, A5. Then for any $\delta \in (0, 1)$ with probability at least $1 - \delta$ for any $k \in \{1, \dots, n\}$ it holds*

$$\|\theta_k - \theta^*\|^2 \leq \alpha_k K_2 \log\left(\frac{en}{\delta}\right),$$

where

$$K_2 = \frac{k_0^\gamma}{c_0} \|\theta_0 - \theta^*\|^2 + \frac{16C_\xi^2}{\mu} (2d + 1) \quad (66)$$

Moreover, it holds for any $k \in \{1, \dots, n\}$ and any $p \geq 2$ that

$$\mathbb{E}^{2/p}[\|\theta_k - \theta^*\|^p] \leq p \alpha_k (e)^{2/p} K_2 / 2.$$

Proof. The proof is similar to the proof of Lemma 13 and Corollary 5. \square

G Technical bounds

We begin this section with useful technical bounds on sums of coefficients

$$\sum_{i=m}^k \alpha_i^p,$$

where the step sizes α_i have a form

$$\alpha_i = \frac{c_0}{(k_0 + i)^\gamma}, \quad 1/2 < \gamma < 1, \quad k_0 \geq 1.$$

We also bound other related quantities, which are instrumental to our further analysis.

Lemma 15. *Assume A8. Then*

(a) *for any $p \geq 2$, it holds that*

$$\sum_{i=1}^k \alpha_i^p \leq \frac{c_0^p}{p\gamma - 1},$$

(b) *for any $m \in \{0, \dots, k\}$, it holds that*

$$\sum_{i=m+1}^k \alpha_i \geq \frac{c_0}{2(1-\gamma)} \left((k+k_0)^{1-\gamma} - (m+k_0)^{1-\gamma} \right),$$

Proof. To prove (a), observe that

$$\sum_{i=1}^k \alpha_i^p \leq c_0^p \int_1^{+\infty} \frac{dx}{x^{p\gamma}} \leq \frac{c_0^p}{p\gamma - 1},$$

To prove (b), note that for any $i \geq 1$ and $k_0 \geq 1$, we have $2(i + k_0)^{-\gamma} \geq (i + k_0 - 1)^{-\gamma}$. Hence,

$$\sum_{i=m+1}^k \alpha_i \geq \frac{1}{2} \sum_{i=m}^{k-1} \alpha_i \geq \frac{c_0}{2} \int_{m+k_0}^{k+k_0} \frac{dx}{x^\gamma} = \frac{c_0}{2(1-\gamma)} ((k+k_0)^{1-\gamma} - (m+k_0)^{1-\gamma}).$$

□

Lemma 16. Lemma 24 in (Durmus et al., 2021) Let $b > 0$ and $\{\alpha_k\}_{k \geq 0}$ be a non-increasing sequence such that $\alpha_0 \leq 1/b$. Then

$$\sum_{j=0}^k \alpha_j \prod_{l=j+1}^k (1 - \alpha_l b) = \frac{1}{b} \left\{ 1 - \prod_{l=0}^k (1 - \alpha_l b) \right\}$$

Proof. Proof of this statement is given in (Durmus et al., 2021).

□

Lemma 17. For any $A > 0$, any $0 \leq i \leq n-1$, and any $\gamma \in (1/2, 1)$ it holds

$$\sum_{j=i}^{n-1} \exp\left\{-A(j^{1-\gamma} - i^{1-\gamma})\right\} \leq \begin{cases} 1 + \exp\left\{\frac{1}{1-\gamma}\right\} \frac{1}{A^{1/(1-\gamma)}(1-\gamma)} \Gamma\left(\frac{1}{1-\gamma}\right), & \text{if } Ai^{1-\gamma} \leq \frac{1}{1-\gamma} \text{ and } i \geq 1; \\ 1 + \frac{1}{A(1-\gamma)^2} i^\gamma, & \text{if } Ai^{1-\gamma} > \frac{1}{1-\gamma} \text{ and } i \geq 1; \\ 1 + \frac{1}{A^{1/(1-\gamma)}(1-\gamma)} \Gamma\left(\frac{1}{1-\gamma}\right), & \text{if } i = 0. \end{cases}$$

Proof. Note that

$$\begin{aligned} \sum_{j=i}^{n-1} \exp\left\{-A(j^{1-\gamma} - i^{1-\gamma})\right\} &\leq 1 + \exp\left\{Ai^{1-\gamma}\right\} \int_i^{+\infty} \exp\left\{-Ax^{1-\gamma}\right\} dx \\ &= 1 + \exp\left\{Ai^{1-\gamma}\right\} \frac{1}{A^{1/(1-\gamma)}(1-\gamma)} \int_{Ai^{1-\gamma}}^{+\infty} e^{-u} u^{\frac{1}{1-\gamma}-1} du \end{aligned}$$

Applying (Gabcke, 2015, Theorem 4.4.3), we get

$$\int_{Ai^{1-\gamma}}^{+\infty} e^{-u} u^{\frac{1}{1-\gamma}-1} du \leq \begin{cases} \Gamma\left(\frac{1}{1-\gamma}\right), & \text{if } Ai^{1-\gamma} < \frac{1}{1-\gamma}; \\ \frac{1}{1-\gamma} \exp\{-Ai^{1-\gamma}\} A^{\gamma/(1-\gamma)} i^\gamma, & \text{otherwise.} \end{cases}$$

Combining inequities above, for $i \geq 1$ we obtain

$$\sum_{j=i}^{n-1} \exp\left\{-A(j^{1-\gamma} - i^{1-\gamma})\right\} \leq \begin{cases} 1 + \exp\left\{\frac{1}{1-\gamma}\right\} \frac{1}{A^{1/(1-\gamma)}(1-\gamma)} \Gamma\left(\frac{1}{1-\gamma}\right), & \text{if } Ai^{1-\gamma} < \frac{1}{1-\gamma}; \\ 1 + \frac{1}{A(1-\gamma)^2} i^\gamma, & \text{otherwise.} \end{cases},$$

and for $i = 0$, we have

$$\sum_{j=0}^{n-1} \exp\left\{-Aj^{1-\gamma}\right\} \leq 1 + \frac{1}{A^{1/(1-\gamma)}(1-\gamma)} \Gamma\left(\frac{1}{1-\gamma}\right).$$

□

G.1 Proof of Lemma 3

Version of Lemma 3 with explicit constants. Assume A1 and A8. Then for any $i \in \{0, \dots, n-1\}$ it holds that

$$\lambda_{\max}(Q_i) \leq C_Q,$$

where the constant C_Q is given by

$$C_Q = \left[1 + \max \left(\exp \left\{ \frac{1}{1-\gamma} \right\} \left(\frac{2(1-\gamma)}{\mu c_0} \right)^{1/(1-\gamma)} \frac{1}{1-\gamma} \Gamma \left(\frac{1}{1-\gamma} \right), \frac{2}{\mu c_0(1-\gamma)} \right) \right] c_0. \quad (67)$$

Moreover,

$$\lambda_{\min}(Q_i) \geq C_Q^{\min}, \text{ and } \|\Sigma_n^{-1/2}\| \leq C_\Sigma, \quad (68)$$

where the matrix Σ_n is defined in (14), and

$$C_Q^{\min} = \frac{1}{L_1} (1 - (1 - \alpha_i L_1)^{n-i}), \quad (69)$$

$$C_\Sigma = \frac{\sqrt{2}L_1}{(1 - \exp\{-\frac{\mu c_0 L_1}{2(k_0+1)}\}) \sqrt{\lambda_{\min}(\Sigma_\xi)}}. \quad (70)$$

Proof. Note that using Lemma 15(b), for $i \geq 0$, it holds that

$$\begin{aligned} \lambda_{\max}(Q_i) &\leq \alpha_i \sum_{j=i}^{n-1} \prod_{k=i+1}^j (1 - \alpha_k \mu) \leq \alpha_i \sum_{j=i}^{n-1} \exp \left\{ -\mu \sum_{k=i+1}^j \alpha_k \right\} \\ &\leq \alpha_i \sum_{j=i+k_0}^{n-1+k_0} \exp \left\{ -\frac{\mu c_0}{2(1-\gamma)} (j^{1-\gamma} - (i+k_0)^{1-\gamma}) \right\}. \end{aligned}$$

Using Lemma 17, we complete the first part with the constant C_Q defined in (67). In order to prove (68), we note that

$$\lambda_{\min}(Q_i) \geq \alpha_i \sum_{j=i}^{n-1} (1 - \alpha_i L_1)^{j-i} = \frac{1}{L_1} (1 - (1 - \alpha_i L_1)^{n-i}).$$

Then for $i \leq n/2$, we have

$$\lambda_{\min}(Q_i) \geq \frac{1}{L_1} (1 - (1 - \alpha_i L_1)^{n/2}) \geq \frac{1}{L_1} (1 - \exp\{-\mu \alpha_i L_1 n/2\}) \geq \frac{1}{L_1} (1 - \exp\{-\frac{\mu c_0 L_1}{2(k_0+1)}\}),$$

where the last inequality holds, since $\alpha_i n \geq \alpha_n n \geq \frac{c_0 n}{k_0+n} \geq \frac{c_0}{1+k_0}$. Combining previous inequalities, we get

$$\lambda_{\min}(\Sigma_n) \geq \lambda_{\min} \left(n^{-1} \sum_{i=1}^{n/2} Q_i \Sigma_\xi Q_i^\top \right) \geq \frac{\lambda_{\min}(\Sigma_\xi)}{2L_1^2} (1 - \exp\{-\frac{\mu c_0 L_1}{2(k_0+1)}\})^2,$$

and (68) follows. \square

Lemma 18. Assume A1 and A3. Then

$$\|H(\theta)\| \leq L_H \|\theta - \theta^*\|^2,$$

where $L_H = \max(L_3, 2L_1/\beta)$.

Proof. Using A3 and the definition of $H(\theta)$ in (10), we get

$$\|H(\theta)\| \mathbf{1}(\|\theta - \theta^*\| \leq \beta) \leq L_3 \|\theta - \theta^*\|^2.$$

Since $\mu \mathbf{I} \preceq \nabla^2 f(\theta) \preceq L_1 \mathbf{I}$, we also obtain

$$\|H(\theta)\| \mathbf{1}(\|\theta - \theta^*\| > \beta) \leq 2L_1 \mathbf{1}(\|\theta - \theta^*\| > \beta) \|\theta - \theta^*\| \leq \frac{2L_1}{\beta} \|\theta - \theta^*\|^2.$$

This concludes the proof. \square

G.2 Properties of sub-Gaussian random vectors

In this section we derive some auxiliary results on sub-Gaussian random vectors. Following (Vershynin, 2018) and (Jin et al., 2019), we rely on the following definition.

Definition 1. A random vector $X \in \mathbb{R}^d$ with $\mathbb{E}[X] = 0$ is called sub-Gaussian with variance proxy $\sigma^2 > 0$, if for any vector $v \in \mathbb{R}^d$,

$$\mathbb{E}[\exp\{\langle v, X \rangle\}] \leq \exp\{\|v\|^2 \sigma^2 / 2\}.$$

In this case, we write $X \in \text{SG}(\sigma^2)$.

Lemma 19. Let $X \in \text{SG}(\sigma^2)$ be a d -dimensional sub-Gaussian vector. Then for any $\lambda \in [0, 1/C_X]$,

$$\mathbb{E}[\exp\{\lambda \|X\|^2\}] \leq \exp\{\lambda C_X\}, \quad (71)$$

where $C_X = 4d\sigma^2$.

Proof. Let $Y \sim \mathcal{N}(0, I_d)$ be a random vector independent of X . Then for fixed $x \in \mathbb{R}^d$, we have

$$\exp\{\lambda \|x\|^2\} = \mathbb{E}[\exp\{\langle x, \sqrt{2\lambda}Y \rangle\}].$$

Hence, we have for $\lambda \in [0, 1/(2\sigma^2)]$:

$$\mathbb{E}[\exp\{\lambda \|X\|^2\}] = \mathbb{E}[\exp\{\langle X, \sqrt{2\lambda}Y \rangle\}] \leq \mathbb{E}[\exp\{\lambda \sigma^2 \|Y\|^2\}] \leq (1 - 2\lambda \sigma^2)^{-d/2}.$$

Note that

$$-1/2 \log(1 - a) \leq a \text{ for } a \in [0, 1/2].$$

Then, we have for $\lambda \in [0, 1/(4\sigma^2)]$:

$$\mathbb{E}[\exp\{\lambda \|X\|^2\}] \leq \exp\{2d\lambda \sigma^2\}.$$

Hence, (71) holds with $C_X = 4d\sigma^2$. \square

G.3 Gaussian comparison lemma

There are quite a few works devoted to the comparison of Gaussian measures with different covariance matrices and means. Among others, we note (Barsov and Ulyanov, 1986), (Götze et al., 2019), and (Devroye et al., 2018). We use the result from (Devroye et al., 2018, Theorem 1.1), which provides a comparison in terms of the total variation distance. Recall that the total variation distance between probability measures μ and ν on a measurable space (X, \mathcal{X}) is defined as

$$d_{\text{TV}}(\mu, \nu) = \sup_{B \in \mathcal{X}} |\mu(B) - \nu(B)|.$$

With a slight abuse of notation, when X and Y are random vectors with distributions μ and ν , respectively, we write $d_{\text{TV}}(X, Y)$ instead of $d_{\text{TV}}(\mu, \nu)$. The following lemma holds:

Lemma 20. Let Σ_1 and Σ_2 be positive definite covariance matrices in $\mathbb{R}^{d \times d}$. Let $X \sim \mathcal{N}(0, \Sigma_1)$ and $Y \sim \mathcal{N}(0, \Sigma_2)$. Then

$$d_{\text{TV}}(X, Y) \leq \frac{3}{2} \|\Sigma_2^{-1/2} \Sigma_1 \Sigma_2^{-1/2} - I_d\|_{\text{F}}.$$

Recall that our primary aim in this paper is to obtain convergence bounds in the convex distance

$$d_{\text{C}}(X, Y) = \sup_{B \in \mathcal{C}(\mathbb{R}^d)} |\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)|,$$

where $\mathcal{C}(\mathbb{R}^d)$ is a collection of convex sets on \mathbb{R}^d . We can immediately obtain the result for convex distance from Lemma 20, since

$$d_{\text{C}}(X, Y) \leq d_{\text{TV}}(X, Y).$$

H Numerical Experiments

To evaluate the behavior of the multiplier bootstrap procedure for constructing confidence sets, we conducted a series of numerical experiments on synthetic linear and logistic regression problems. The overall experimental procedure is identical for both models, and we highlight only the model-specific differences below. Our experimental pipeline follows the main steps outlined in (Fang et al., 2018). Our experiments were conducted on a single Intel Xeon Gold 6248R CPU (48 cores, 3.0–4.0 GHz), 768 GB RAM, and 240 GB SSD storage, without GPU accelerators. Code to reproduce the experiments is provided in <https://github.com/marina-shesha/Gaussian-Approximation-and-Multiplier-Bootstrap-for-Stochastic-Gradient-Descent>.

H.1 Experimental Setup

For each of 1024 trajectories, we generate a dataset (X, y) with sample sizes $N \in \{10000, 20000, 30000, 40000, 50000\}$ and $d = 5$ features. For both models, the true parameter vector is fixed as

$$\theta_{\text{true}} = [1.0, 1.0, -0.1, -0.1, -0.1].$$

For each trajectory we run the SGD algorithm with step sizes

$$\alpha_n = \frac{200}{(20000 + n)^{0.85}},$$

and compute the Polyak–Ruppert averaged estimator $\bar{\theta}_n$.

To assess coverage probabilities, we employ a multiplier bootstrap procedure. For each trajectory, we generate $N_{\text{boot}} = 256$ bootstrap trajectories. In each bootstrap run, the step sizes α_n are multiplicatively perturbed by independent random variables drawn from a Beta distribution:

$$\tilde{\alpha}_n = \alpha_n \left(1 + \frac{w_n - \mathbb{E}[w_n]}{\sqrt{\text{Var } w_n}} \right), \quad \text{where } w_n \sim \text{Beta}(0.5, 2).$$

For each trajectory length N , we construct confidence intervals for the one-dimensional functional given by a random projection of the parameter target vector. Specifically, for each trajectory we draw a unit vector $u \in \mathbb{S}^{d-1}$ (fixed with all trajectories and its bootstrap replicates) and form confidence intervals for the scalar target parameter. The coverage probabilities for this scalar target parameter are then estimated using three approaches:

1. **Empirical Quantiles:** We compute the empirical quantiles of the N_{boot} bootstrap replicates and check whether the target parameter belongs to the determined interval.
2. **Standard Deviation–Based Confidence Intervals:** We construct confidence intervals based on the sample standard deviation of the bootstrap replicates and verify whether target parameter belongs to the determined interval.
3. **Overlapping Batch Mean (OBM) Estimator (Meketon and Schmeiser, 1984; Flegal and Jones, 2010):** For the trajectory with $N = 50000$, we estimate the asymptotic variance as follows:

$$\hat{\sigma}_{\bar{\theta}}^2(u) = \frac{b_n}{n - b_n + 1} \sum_{t=0}^{n-b_n} ((\bar{\theta}_{b_n,t} - \bar{\theta}_n)^\top u)^2, \quad \bar{\theta}_{b_n,t} = \frac{1}{b_n} \sum_{\ell=t}^{t+b_n-1} \theta_\ell.$$

Using this estimated asymptotic variance, we construct confidence intervals and check whether the target parameter lies within them. The procedure is repeated for several values of the batch size b_n , and we select the one that provides the best coverage performance for the 0.95 confidence interval.

The coverage probabilities are obtained by averaging the indicator of interval inclusion over all 1024 trajectories.

H.2 Linear Regression

The feature vectors $X_i \in \mathbb{R}^p$ are sampled uniformly from $[-1, 1]$, and the response is generated according to

$$Y_i = X_i^\top \theta_{\text{true}} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 0.02^2).$$

Then the target parameter is equal to θ_{true} . The results are summarized in Table 1.

Table 1: Comparison of Different Estimation Methods for Linear Regression Problem

| (a) Empirical Quantiles | | | |
|-------------------------|----------|----------|----------|
| Trajectory length | 0.95 | 0.90 | 0.80 |
| 10000 | 0.964844 | 0.936523 | 0.855469 |
| 20000 | 0.958008 | 0.919922 | 0.852539 |
| 30000 | 0.961914 | 0.927734 | 0.839844 |
| 40000 | 0.963867 | 0.924805 | 0.834961 |
| 50000 | 0.957031 | 0.920898 | 0.829102 |

| (b) Standard Deviation-Based Confidence Intervals | | | |
|---|----------|----------|----------|
| Trajectory length | 0.95 | 0.90 | 0.80 |
| 10000 | 0.967773 | 0.943359 | 0.859375 |
| 20000 | 0.963867 | 0.927734 | 0.852539 |
| 30000 | 0.964844 | 0.932617 | 0.844727 |
| 40000 | 0.967773 | 0.934570 | 0.837891 |
| 50000 | 0.959961 | 0.924805 | 0.825195 |

| (c) Overlapping Batch Mean Estimator | | | |
|---|----------|----------|----------|
| Batch size (b_n), Trajectory length | 0.95 | 0.90 | 0.80 |
| 1700, 50000 | 0.890625 | 0.818359 | 0.712891 |

H.3 Logistic Regression

Here, $X_i \in \mathbb{R}^p$ are sampled uniformly from $[-1, 3]$, and responses follow the distribution

$$\mathbb{P}(Y_i = 1) = \sigma(X_i^\top \theta_{\text{true}}), \quad \mathbb{P}(Y_i = -1) = 1 - \mathbb{P}(Y_i = 1).$$

where $\sigma(\cdot)$ denotes the sigmoid function. We minimize the L_2 -regularized logistic loss with regularization parameter $\lambda = 10^{-4}$.

$$\theta_{\text{star}} = \arg \min_{\theta \in \mathbb{R}^d} \left\{ \mathbb{E} \left[-\log \left(\frac{1}{1 + \exp\{-Y X^\top \theta\}} \right) \right] + \lambda \|\theta\|^2 \right\}$$

Since there is no closed-form solution, we estimate the target parameter θ_{star} by running one long SGD trajectory of length 10^6 and computing Polyak-Ruppert averaged estimator along this trajectory. The results are summarized in Table 2.

H.4 Discussion

Both empirical quantiles and confidence intervals based on standard deviation are variants of the bootstrap procedure, and in our experiments they provide coverage very close to nominal levels across the entire trajectory length. In contrast, the Overlapping Batch Mean estimator systematically underestimates the coverage even when the trajectory length is rather large ($N = 50000$). Moreover, the procedure itself is very sensitive to the choice of batch size b_n , which moderate heuristics for this problem available (Flegal and Jones, 2010).

Table 2: Comparison of Different Estimation Methods for Logistic Regression Problem

(a) Empirical Quantiles

| Trajectory length | 0.95 | 0.90 | 0.80 |
|-------------------|----------|----------|----------|
| 10000 | 0.947266 | 0.901367 | 0.806641 |
| 20000 | 0.940430 | 0.909180 | 0.800781 |
| 30000 | 0.943359 | 0.903320 | 0.798828 |
| 40000 | 0.943359 | 0.894531 | 0.814453 |
| 50000 | 0.945312 | 0.888672 | 0.785156 |

(b) Standard Deviation-Based Confidence Intervals

| Trajectory length | 0.95 | 0.90 | 0.80 |
|-------------------|----------|----------|----------|
| 10000 | 0.961914 | 0.915039 | 0.833984 |
| 20000 | 0.955078 | 0.921875 | 0.826172 |
| 30000 | 0.958984 | 0.916016 | 0.822266 |
| 40000 | 0.952148 | 0.906250 | 0.825195 |
| 50000 | 0.955078 | 0.913086 | 0.815430 |

(c) Overlapping Batch Mean Estimator

| Batch size (b_n), Trajectory length | 0.95 | 0.90 | 0.80 |
|---|----------|----------|----------|
| 2500, 50000 | 0.920898 | 0.856445 | 0.743164 |