

# Bridging Stochastic Gradient Descent and Markov Chains: Constant Step-Size Convergence and Richardson-Romberg Extrapolation

Chris Junchi Li<sup>◊</sup>

Department of Electrical Engineering and Computer Sciences<sup>◊</sup>  
University of California, Berkeley

October 4, 2024

## Abstract

Stochastic Gradient Descent (SGD) is a fundamental optimization method widely used in machine learning and related fields. This work revisits constant step-size SGD and its connection to Markov chains. We provide an asymptotic expansion of the averaged SGD iterates, highlighting the key effects of noise, initial conditions, and step-size choices on convergence behavior. Furthermore, we introduce Richardson-Romberg extrapolation as a technique to accelerate convergence towards the optimum, even in the presence of persistent oscillations around the optimal point. Our theoretical analysis is complemented by empirical results demonstrating the improvements obtained through this method. These insights allow for better understanding of both the bias-variance trade-off and the convergence behavior of SGD in high-dimensional, noisy environments.

**Keywords:** Stochastic Gradient Descent, Markov Chains, Constant Step-Size, Richardson-Romberg Extrapolation, Bias-Variance Trade-off

## 1 Introduction

Stochastic Gradient Descent (SGD) has become one of the most essential tools in large-scale optimization, with wide applications in machine learning, deep learning, and related fields. Its simplicity and scalability have made it indispensable for minimizing objective functions where only stochastic estimates of gradients are available. Despite its popularity, the behavior of SGD with constant step-size remains underexplored in many settings, particularly when considering its non-convergence to an exact optimum in general, due to the persistent noise inherent in stochastic approximations.

In this paper, we focus on constant step-size SGD, a practical variant used extensively due to its ease of implementation and reasonable empirical performance. A significant challenge in constant step-size SGD is understanding how the iterates behave, oscillating around the optimum rather than converging. Our approach utilizes a Markov chain framework to model SGD iterates, shedding light on their long-term behavior, including the rapid forgetting of initial conditions and the trade-offs between noise and step-size.

One of the primary contributions of this work is the introduction of Richardson-Romberg extrapolation as a technique to accelerate the convergence of constant step-size SGD. By combining iterates with different step-sizes, this method provably reduces the bias introduced by constant step-size while maintaining computational efficiency.

**Settings** We consider the minimization of an objective function given access to unbiased estimates of the function gradients. This key methodological problem has raised interest in different

communities: in large-scale machine learning [9, 51, 52], optimization [41, 44], and stochastic approximation [27, 46, 50]. The most widely used algorithms are stochastic gradient descent (SGD), a.k.a. Robbins-Monro algorithm [49], and some of its modifications based on averaging of the iterates [46, 48, 53].

While the choice of the step-size may be done robustly in the deterministic case (see *e.g.* [8]), this remains a traditional theoretical and practical issue in the stochastic case. Indeed, early work suggested to use step-size decaying with the number  $k$  of iterations as  $O(1/k)$  [49], but it appeared to be non-robust to ill-conditioning and slower decays such as  $O(1/\sqrt{k})$  together with averaging lead to both good practical and theoretical performance [3].

We consider in this paper constant step-size SGD, which is often used in practice. Although the algorithm is not converging in general to the global optimum of the objective function, constant step-sizes come with benefits: (a) there is a single parameter value to set as opposed to the several choices of parameters to deal with decaying step-sizes, *e.g.* as  $1/(\square k + \triangle)^\circ$ ; the initial conditions are forgotten exponentially fast for well-conditioned (*e.g.* strongly convex) problems [39, 40], and the performance, although not optimal, is sufficient in practice (in a machine learning set-up, being only 0.1% away from the optimal prediction often does not matter).

The main goals of this paper are (a) to gain a complete understanding of the properties of constant-step-size SGD in the strongly convex case, and (b) to propose provable improvements to get closer to the optimum when precision matters or in high-dimensional settings. We consider the iterates of the SGD recursion on  $\mathbb{R}^d$  defined starting from  $\theta_0 \in \mathbb{R}^d$ , for  $k \geq 0$ , and a step-size  $\gamma > 0$  by

$$\theta_{k+1}^{(\gamma)} = \theta_k^{(\gamma)} - \gamma[f'(\theta_k^{(\gamma)}) + \varepsilon_{k+1}(\theta_k^{(\gamma)})] , \quad (1)$$

where  $f$  is the objective function to minimize (in machine learning the generalization performance),  $\varepsilon_{k+1}(\theta_k^{(\gamma)})$  the zero-mean statistically independent noise (in machine learning, obtained from a single observation). Following [5], we leverage the property that the sequence of iterates  $(\theta_k^{(\gamma)})_{k \geq 0}$  is an *homogeneous Markov chain*.

This interpretation allows us to capture the general behavior of the algorithm. In the strongly convex case, this Markov chain converges exponentially fast to a unique stationary distribution  $\pi_\gamma$  (see Proposition 1) highlighting the facts that (a) initial conditions of the algorithms are forgotten quickly and (b) the algorithm does not converge to a point but oscillates around the mean of  $\pi_\gamma$ . See an illustration in Figure 1 (left). It is known that the oscillations of the non-averaged iterates have an average magnitude of  $\gamma^{1/2}$  [45].

Consider the process  $(\bar{\theta}_k^{(\gamma)})_{k \geq 0}$  given for all  $k \geq 0$  by

$$\bar{\theta}_k^{(\gamma)} = \frac{1}{k+1} \sum_{j=0}^k \theta_j^{(\gamma)} . \quad (2)$$

Then under appropriate conditions on the Markov chain  $(\theta_k^{(\gamma)})_{k \geq 0}$ , a central limit theorem on  $(\bar{\theta}_k^{(\gamma)})_{k \geq 0}$  holds which implies that  $\bar{\theta}_k^{(\gamma)}$  converges at rate  $O(1/\sqrt{k})$  to

$$\bar{\theta}_\gamma = \int_{\mathbb{R}^d} \vartheta \, d\pi_\gamma(\vartheta) . \quad (3)$$

The deviation between  $\bar{\theta}_k^{(\gamma)}$  and the global optimum  $\theta^*$  is thus composed of a stochastic part  $\bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma$  and a deterministic part  $\bar{\theta}_\gamma - \theta^*$ .

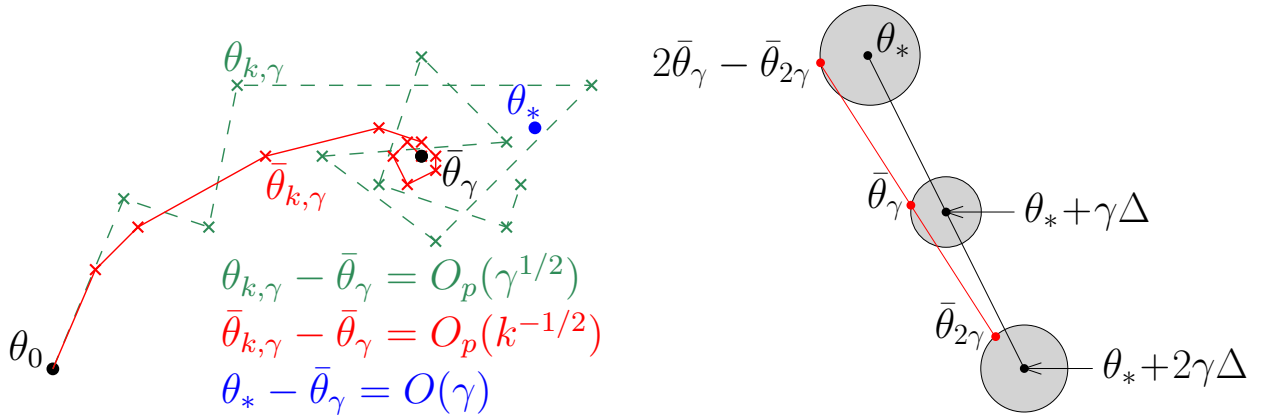
For quadratic functions, it turns out that the deterministic part vanishes [5], that is,  $\bar{\theta}_\gamma = \theta^*$  and thus averaged SGD with a constant step-size does converge. However, it is not true for general objective functions where we can only show that  $\bar{\theta}_\gamma - \theta^* = O(\gamma)$ , and this deviation is the reason why constant step-size SGD is not convergent.

The first main contribution of the paper is to provide an explicit asymptotic expansion in the step-size  $\gamma$  of  $\bar{\theta}_\gamma - \theta^*$ . Second, a quantitative version of a central limit theorem is established which gives a bound on  $\mathbb{E}[\|\bar{\theta}_\gamma - \bar{\theta}_k^{(\gamma)}\|^2]$  that highlights all dependencies on initial conditions and noise variance, as achieved for least-squares by [14], with an explicit decomposition into “bias” and “variance” terms: the bias term characterizes how fast initial conditions are forgotten and is proportional to  $N(\theta_0 - \theta^*)$ , for a suitable norm  $N : \mathbb{R}^d \rightarrow \mathbb{R}_+$ ; while the variance term characterizes the effect of the noise in the gradient, independently of the starting point, and increases with the covariance of the noise.

Moreover, akin to weak error results for ergodic diffusions, we achieve a non-asymptotic weak error expansion in the step-size between  $\pi_\gamma$  and the Dirac measure on  $\mathbb{R}^d$  concentrated at  $\theta^*$ . Namely, we prove that for all functions  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ , regular enough,  $\int_{\mathbb{R}^d} g(\theta) d\pi_\gamma(\theta) = g(\theta^*) + \gamma C_1^g + r_\gamma^g$ ,  $r_\gamma^g \in \mathbb{R}^d$ ,  $\|r_\gamma^g\| \leq C_2^g \gamma^2$ , for some  $C_1^g, C_2^g \geq 0$  independent of  $\gamma$ . Given this expansion, we can now use a very simple trick from numerical analysis, namely Richardson-Romberg extrapolation [54]: if we run two SGD recursions  $(\theta_k^{(\gamma)})_{k \geq 0}$  and  $(\theta_k^{(2\gamma)})_{k \geq 0}$  with the two different step-sizes  $\gamma$  and  $2\gamma$ , then the average processes  $(\bar{\theta}_k^{(\gamma)})_{k \geq 0}$  and  $(\bar{\theta}_k^{(2\gamma)})_{k \geq 0}$  will converge to  $\bar{\theta}_\gamma$  and  $\bar{\theta}_{2\gamma}$  respectively. Since  $\bar{\theta}_\gamma = \theta^* + \gamma \Delta_1^{\text{Id}} + r_\gamma^{\text{Id}}$  and  $\bar{\theta}_{2\gamma} = \theta^* + 2\gamma \Delta_1^{\text{Id}} + r_{2\gamma}^{\text{Id}}$ , for  $r_\gamma^{\text{Id}}, r_{2\gamma}^{\text{Id}} \in \mathbb{R}^d$ ,  $\max(\|2r_\gamma^{\text{Id}}\|, \|r_{2\gamma}^{\text{Id}}\|) \leq 2C\gamma^2$ , for  $C \geq 0$  and  $\Delta \in \mathbb{R}^d$  independent of  $\gamma$ , the combined iterates  $2\bar{\theta}_k^{(\gamma)} - \bar{\theta}_k^{(2\gamma)}$  will converge to  $\theta^* + 2r_\gamma^{\text{Id}} - r_{2\gamma}^{\text{Id}}$  which is closer to  $\theta^*$  by a factor  $\gamma$ . See illustration in Figure 1(right).

In summary, we make the following contributions:

- We provide in Section 2 an asymptotic expansion in  $\gamma$  of  $\bar{\theta}_\gamma - \theta^*$  and an explicit version of a central limit theorem is given which bounds  $\mathbb{E}[\|\bar{\theta}_\gamma - \bar{\theta}_k^{(\gamma)}\|^2]$ . These two results outlines the dependence on initial conditions, the effect of noise and the step-size.
- We show in Section 2 that Richardson-Romberg extrapolation may be used to get closer to the global optimum.
- We bring and adapt in Section 3 tools from analysis of discretization of diffusion processes



**Figure 1.** (Left) Convergence of iterates  $\theta_k^{(\gamma)}$  and averaged iterates  $\bar{\theta}_k^{(\gamma)}$  to the mean  $\bar{\theta}_\gamma$  under the stationary distribution  $\pi_\gamma$ . (Right) Richardson-Romberg extrapolation, the disks are of radius  $O(\gamma^2)$ .

into the one of SGD and create new ones. We believe that this analogy and the associated ideas are interesting in their own right.

- We show in Section 4 empirical improvements of the extrapolation schemes.

**Notations** We first introduce several notations. We consider the finite dimensional euclidean space  $\mathbb{R}^d$  embedded with its canonical inner product  $\langle \cdot, \cdot \rangle$ . Denote by  $\{\mathbf{e}_1, \dots, \mathbf{e}_d\}$  the canonical basis of  $\mathbb{R}^d$ . Let  $E$  and  $F$  be two real vector spaces, denote by  $E \otimes F$  the tensor product of  $E$  and  $F$ . For all  $x \in E$  and  $y \in F$  denote by  $x \otimes y \in E \otimes F$  the tensor product of  $x$  and  $y$ . Denote by  $E^{\otimes k}$  the  $k^{\text{th}}$  tensor power of  $E$  and  $x^{\otimes k} \in E^{\otimes k}$  the  $k^{\text{th}}$  tensor power of  $x$ . We let  $\mathcal{L}((\mathbb{R}^d)^{\otimes k}, \mathbb{R}^\ell)$  stand for the set of linear maps from  $(\mathbb{R}^d)^{\otimes k}$  to  $\mathbb{R}^\ell$  and for  $L \in \mathcal{L}((\mathbb{R}^d)^{\otimes k}, \mathbb{R}^\ell)$ , we denote by  $\|L\|$  the operator norm of  $L$ .

Let  $n \in \mathbb{N}^*$ , denote by  $C^n(\mathbb{R}^d, \mathbb{R}^m)$  the set of  $n$  times continuously differentiable functions from  $\mathbb{R}^d$  to  $\mathbb{R}^m$ . Let  $f \in C^n(\mathbb{R}^d, \mathbb{R}^m)$ , denote by  $F^{(n)}$  or  $D^n f$ , the  $n^{\text{th}}$  differential of  $f$ . Let  $f \in C^n(\mathbb{R}^d, \mathbb{R})$ . For any  $x \in \mathbb{R}^d$ ,  $f^{(n)}(x)$  is a tensor of order  $n$ . For example, for all  $x \in \mathbb{R}^d$ ,  $f^{(3)}(x)$  is a third order tensor. In addition, for any  $x \in \mathbb{R}^d$  and any matrix,  $M \in \mathbb{R}^{d \times d}$ , we define  $f^{(3)}(x)M$  as the vector in  $\mathbb{R}^d$  given by: for any  $l \in \{1, \dots, d\}$ , the  $l^{\text{th}}$  coordinate is given by  $(f^{(3)}(x)M)_l = \sum_{i,j=1}^d M_{i,j} \frac{\partial^3 f}{\partial x_i \partial x_j \partial x_l}(x)$ . By abuse of notations, for  $f \in C^1(\mathbb{R}^d)$ , we identify  $f'$  with the gradient of  $f$  and if  $f \in C^2(\mathbb{R}^d)$ , we identify  $f''$  with the Hessian matrix of  $f$ . A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^q$  is said to be locally Lipschitz if there exists  $\alpha \geq 0$  such that for all  $x, y \in \mathbb{R}^d$ ,  $\|f(x) - f(y)\| \leq (1 + \|x\|^\alpha + \|y\|^\alpha) \|x - y\|$ . For ease of notations and depending on the context, we consider  $M \in \mathbb{R}^{d \times d}$  either as a matrix or a second order tensor. More generally, any  $M \in \mathcal{L}((\mathbb{R}^d)^{\otimes k}, \mathbb{R})$  will be also consider as an element of  $\mathcal{L}((\mathbb{R}^d)^{\otimes(k-1)}, \mathbb{R}^d)$  by the canonical bijection. Besides, For any matrices  $M, N \in \mathbb{R}^{d \times d}$ ,  $M \otimes N$  is defined as the endomorphism of  $\mathbb{R}^{d \times d}$  such that  $M \otimes N : P \mapsto MPN$ . For any matrix  $M \in \mathbb{R}^{d \times d}$ ,  $\text{tr}(M)$  is the trace of  $M$ , *i.e.* the sum of diagonal elements of the matrix  $M$ .

For  $a, b \in \mathbb{R}$ , denote by  $a \vee b$  and  $a \wedge b$  the maximum and the minimum of  $a$  and  $b$  respectively. Denote by  $\lfloor \cdot \rfloor$  and  $\lceil \cdot \rceil$  the floor and ceiling function respectively.

Denote by  $\mathcal{B}(\mathbb{R}^d)$  the Borel  $\sigma$ -field of  $\mathbb{R}^d$ . For all  $x \in \mathbb{R}^d$ ,  $\delta_x$  stands for the Dirac measure at  $x$ .

## 2 Main results

In this section, we describe the assumptions underlying our analysis, describe our main results and their implications.

### 2.1 Setting

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be an objective function, satisfying the following assumptions:

**Assumption 1.** *The function  $f$  is strongly convex with convexity constant  $\mu > 0$ , *i.e.* for all  $\theta_1, \theta_2 \in \mathbb{R}^d$  and  $t \in [0, 1]$ ,*

$$f(t\theta_1 + (1-t)\theta_2) \leq tf(\theta_1) + (1-t)f(\theta_2) - (\mu/2)t(1-t) \|\theta_1 - \theta_2\|^2 .$$

**Assumption 2.** *The function  $f$  is five times continuously differentiable with second to fifth uniformly bounded derivatives: for all  $k \in \{2, \dots, 5\}$ ,  $\sup_{\theta \in \mathbb{R}^d} \|f^{(k)}(\theta)\| < +\infty$ . Especially  $f$  is  $L$ -smooth with  $L \geq 0$ : for all  $\theta_1, \theta_2 \in \mathbb{R}^d$*

$$\|f'(\theta_1) - f'(\theta_2)\| \leq L \|\theta_1 - \theta_2\| .$$

If there exists a positive definite matrix  $\Sigma \in \mathbb{R}^{d \times d}$ , such that the function  $f$  is the quadratic function  $\theta \mapsto \|\Sigma^{1/2}(\theta - \theta^*)\|^2/2$ , then Assumptions [Assumption 1](#), [Assumption 2](#) are satisfied.

In the definition of SGD given by [\(1\)](#),  $(\varepsilon_k)_{k \geq 1}$  is a sequence of random functions from  $\mathbb{R}^d$  to  $\mathbb{R}^d$  satisfying the following properties.

**Assumption 3.** *There exists a filtration  $(\mathcal{F}_k)_{k \geq 0}$  (i.e. for all  $k \in \mathbb{N}$ ,  $\mathcal{F}_k \subset \mathcal{F}_{k+1}$ ) on some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  such that for any  $k \in \mathbb{N}$  and  $\theta \in \mathbb{R}^d$ ,  $\varepsilon_{k+1}(\theta)$  is a  $\mathcal{F}_{k+1}$ -measurable random variable and  $\mathbb{E}[\varepsilon_{k+1}(\theta) | \mathcal{F}_k] = 0$ . In addition,  $(\varepsilon_k)_{k \in \mathbb{N}^*}$  are independent and identically distributed (i.i.d.) random fields. Moreover, we assume that  $\theta_0$  is  $\mathcal{F}_0$ -measurable.*

[Assumption 3](#) expresses that we have access to an i.i.d. sequence  $(f'_k)_{k \in \mathbb{N}^*}$  of unbiased estimator of  $f'$ , i.e. for all  $k \in \mathbb{N}$  and  $\theta \in \mathbb{R}^d$ ,

$$f'_{k+1}(\theta) = f'(\theta) + \varepsilon_{k+1}(\theta). \quad (4)$$

Note that we do not assume random vectors  $(\varepsilon_{k+1}(\theta_k^{(\gamma)}))_{k \in \mathbb{N}}$  to be i.i.d., a stronger assumption generally referred to as the semi-stochastic setting. Moreover, as  $\theta_0$  is  $\mathcal{F}_0$ -measurable, for any  $k \in \mathbb{N}$ ,  $\theta_k$  is  $\mathcal{F}_k$ -measurable.

We also consider the following conditions on the noise, for  $p \geq 2$ :

**Assumption 4** ( $p$ ). *For any  $k \in \mathbb{N}^*$ ,  $f'_k$  is almost surely  $L$ -co-coercive (with the same constant as in [Assumption 2](#)): that is, for any  $\eta, \theta \in \mathbb{R}^d$ ,  $L \langle f'_k(\theta) - f'_k(\eta), \theta - \eta \rangle \geq \|f'_k(\theta) - f'_k(\eta)\|^2$ . Moreover, there exists  $\tau_p \geq 0$ , such that for any  $k \in \mathbb{N}^*$ ,  $\mathbb{E}^{1/p}[\|\varepsilon_k(\theta^*)\|^p] \leq \tau_p$ .*

Almost sure  $L$ -co-coercivity [\[59\]](#) is for example satisfied if for any  $k \in \mathbb{N}^*$ , there exists a random function  $f_k$  such that  $f'_k = (f_k)'$  and which is a.s. convex and  $L$ -smooth. Weaker assumptions on the noise are discussed in [Section 6.1](#). Finally we emphasize that under [Assumption 3](#) then to verify that [Assumption 4](#)( $p$ ) holds,  $p \geq 2$ , it suffices to show that  $f'_1$  is almost surely  $L$ -co-coercive and  $\mathbb{E}^{1/p}[\|\varepsilon_1(\theta^*)\|^p] \leq \tau_p$ . Under [Assumption 3](#)-[Assumption 4](#)(2), consider the function  $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$  defined for all  $\theta \in \mathbb{R}^d$  by

$$\mathcal{C}(\theta) = \mathbb{E}[\varepsilon_1(\theta)^{\otimes 2}]. \quad (5)$$

**Assumption 5.** *The function  $\mathcal{C}$  is three time continuously differentiable and there exist  $M_\varepsilon, k_\varepsilon \geq 0$  such that for all  $\theta \in \mathbb{R}^d$ ,*

$$\max_{i \in \{1, 2, 3\}} \left\| \mathcal{C}^{(i)}(\theta) \right\| \leq M_\varepsilon \left\{ 1 + \|\theta - \theta^*\|^{k_\varepsilon} \right\}.$$

In other words, we assume that the covariance matrix  $\theta \mapsto \mathcal{C}(\theta)$  is a regular enough function, which is satisfied in natural settings.

**Example 1** (Learning from i.i.d. observations). Our main motivation comes from machine learning; consider two sets  $\mathcal{X}, \mathcal{Y}$  and a convex loss function  $L : \mathcal{X} \times \mathcal{Y} \times \mathbb{R}^d \rightarrow \mathbb{R}$ . The objective function is the generalization error  $f_L(\theta) = \mathbb{E}_{X,Y}[L(X,Y,\theta)]$ , where  $(X,Y)$  are some random variables. Given i.i.d. observations  $(X_k, Y_k)_{k \in \mathbb{N}^*}$  with the same distribution as  $(X,Y)$ , for any  $k \in \mathbb{N}^*$ , we define  $f_k(\cdot) = L(X_k, Y_k, \cdot)$  the loss with respect to observation  $k$ . SGD then corresponds to following gradient of the loss on a single independent observation  $(X_k, Y_k)$  at each step; [Assumption 3](#) is then satisfied with  $\mathcal{F}_k = \sigma((X_j, Y_j)_{j \in \{1, \dots, k\}})$ .

Two classical situations are worth mentioning. On the first hand, in *least-squares regression*,  $\mathcal{X} = \mathbb{R}^d$ ,  $\mathcal{Y} = \mathbb{R}$ , and the loss function is  $L(X, Y, \theta) = (\langle X, \theta \rangle - Y)^2$ . Then  $f_\Sigma$  is the quadratic function  $\theta \mapsto \|\Sigma^{1/2}(\theta - \theta^*)\|^2/2$ , with  $\Sigma = \mathbb{E}[XX^\top]$ , which satisfies Assumption 2. For any  $\theta \in \mathbb{R}^d$ ,

$$\varepsilon_k(\theta) = X_k X_k^\top \theta - X_k Y_k \quad (6)$$

Then, for any  $p \geq 2$ , Assumption 4(p) and Assumption 5 is satisfied as soon as observations are a.s. bounded, while Assumption 1 is satisfied if the second moment matrix is invertible or additional regularization is added. In this setting,  $\varepsilon_k$  can be decomposed as  $\varepsilon_k = \varrho_k + \xi_k$  where  $\varrho_k$  is the multiplicative part,  $\xi_k$  the additive part, given for  $\theta \in \mathbb{R}^d$  by  $\varrho_k(\theta) = (X_k X_k^\top - \Sigma)(\theta - \theta^*)$  and

$$\xi_k = (X_k^\top \theta^* - Y_k) X_k. \quad (7)$$

For all  $k \geq 1$ ,  $\xi_k$  does not depend on  $\theta$ . This two parts in the noise will appear in Corollary 1. Finally assume that there exists  $r \geq 0$  such that

$$\mathbb{E}[\|X_k\|^2 X_k X_k^\top] \preceq r^2 \Sigma, \quad (8)$$

then Assumption 4(4) is satisfied. This assumption is satisfied, *e.g.*, for a.s. bounded data, or for data with bounded kurtosis, see [17] for details.

On the other hand, in *logistic regression*, where  $L(X, Y, \theta) = \log(1 + \exp(-Y \langle X, \theta \rangle))$ . Assumptions Assumption 4 or Assumption 2 are similarly satisfied, while Assumption 1 needs an additional restriction to a compact set.

## 2.2 Summary and discussion of main results

Under the stated assumptions, for all  $\gamma \in (0, 2/L)$  and  $\theta_0 \in \mathbb{R}^d$ , the Markov chain  $(\theta_k^{(\gamma)})_{k \geq 0}$  converges in a certain sense specified below to a probability measure on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ ,  $\pi_\gamma$  satisfying  $\int_{\mathbb{R}^d} \|\vartheta\|^2 \pi_\gamma(d\vartheta) < +\infty$ , see Proposition 1 in Section 3. In the next section, by two different methods (Theorem 1 and Theorem 3), we show that under suitable conditions on  $f$  and the noise  $(\varepsilon_k)_{k \geq 1}$ , there exists  $\Delta \in \mathbb{R}^d$  such that for all  $\gamma \geq 0$ , small enough

$$\bar{\theta}_\gamma = \int_{\mathbb{R}^d} \vartheta \pi_\gamma(d\vartheta) = \theta^* + \gamma \Delta + r_\gamma^{(1)},$$

where  $r_\gamma^{(1)} \in \mathbb{R}^d$ ,  $\|r_\gamma^{(1)}\| \leq C\gamma^2$  for some constant  $C \geq 0$  independent of  $\gamma$ . Using Proposition 1, we get that for all  $k \geq 1$ ,

$$\mathbb{E}[\bar{\theta}_k^{(\gamma)} - \theta^*] = \frac{A(\theta_0, \gamma)}{k} + \gamma \Delta + r_\gamma^{(2)}, \quad (9)$$

where  $r_\gamma^{(2)} \in \mathbb{R}^d$ ,  $\|r_\gamma^{(2)}\| \leq C(\gamma^2 + e^{-k\mu\gamma})$  for some constant  $C \geq 0$  independent of  $\gamma$ .

This expansion in the step-size  $\gamma$  shows that a Richardson-Romberg extrapolation can be used to have better estimates of  $\theta^*$ . Consider the average iterates  $(\bar{\theta}_{2\gamma}^{(k)})_{k \geq 0}$  and  $(\bar{\theta}_k^{(\gamma)})_{k \geq 0}$  associated with SGD with step size  $2\gamma$  and  $\gamma$  respectively. Then (9) shows that  $(2\bar{\theta}_k^{(\gamma)} - \bar{\theta}_k^{(2\gamma)})_{k \geq 0}$  satisfies

$$\mathbb{E}[2\bar{\theta}_k^{(\gamma)} - \bar{\theta}_k^{(2\gamma)} - \theta^*] = \frac{2A(\theta_0, \gamma) - A(\theta_0, 2\gamma)}{k} + 2r_\gamma^{(2)} - r_{2\gamma}^{(2)},$$

and therefore is closer to the optimum  $\theta^*$ . This very simple trick improves the convergence by a factor of  $\gamma$  (at the expense of a slight increase of the variance). In practice, while the un-averaged

gradient iterate  $\theta_k^{(\gamma)}$  saturates rapidly,  $\bar{\theta}_k^{(\gamma)}$  may already perform well enough to avoid saturation on real data-sets [5]. The Richardson-Romberg extrapolated iterate  $2\bar{\theta}_k^{(\gamma)} - \bar{\theta}_k^{(2\gamma)}$  very rarely reaches saturation in practice. This appears in synthetic experiments presented in Section 4. Moreover, this procedure only requires to compute two parallel SGD recursions, either with the same inputs, or with different ones, and is naturally parallelizable.

In Section 3.2, we give a quantitative version of a central limit theorem for  $(\bar{\theta}_k^{(\gamma)})_{k \geq 0}$ , for a fixed  $\gamma > 0$  and  $k$  going to  $+\infty$ : under appropriate conditions, there exist constants  $B_1(\gamma)$  and  $B_2(\gamma)$  such that

$$\mathbb{E} \left[ \left\| \bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma \right\|^2 \right] = B_1(\gamma)/k + B_2(\gamma)/k^2 . \quad (10)$$

Combining (9) and (10) characterizes the bias/variance trade-off of SGD used to estimate  $\theta^*$ .

### 2.3 Related work

The idea to study stochastic approximation algorithms using results and techniques from the Markov chain literature is not new. It goes back to [22], which shows under appropriate conditions that solutions of stochastic differential equations (SDE)

$$dY_t = -f'(Y_t)dt + \gamma_t dB_t ,$$

where  $(B_t)_{t \geq 0}$  is a  $d$ -dimensional Brownian motion and  $(\gamma_t)_{t \geq 0}$  is a one-dimensional positive function,  $\lim_{t \rightarrow +\infty} \gamma_t = 0$ , converge in probability to some minima of  $f$ . An other example is [47] which extends the classical Foster-Lyapunov criterion from Markov chain theory (see [37]) to study the stability of the LMS algorithm. In [10], the authors are interested in the convergence of the multidimensional Kohonen algorithm. They show that the Markov chain defined by this algorithm is uniformly ergodic and derive asymptotic properties on its limiting distribution.

The techniques we use in this paper to establish our results share a lot of similarities with previous work. For example, our first results in Section 3.1 and Section 3.2 regarding an asymptotic expansion in  $\gamma$  of  $\bar{\theta}_\gamma - \theta^*$  and an explicit version of a central limit theorem is given which bounds  $\mathbb{E}[\|\bar{\theta}_\gamma - \bar{\theta}_k^{(\gamma)}\|^2]$ , can be seen as complementary results of [2]. Indeed, in [2], the authors decompose the tracking error of a general algorithm in a linear regression model. To prove their result, they develop the error using a perturbation approach, which is quite similar to what we do.

Another and significant point of view to study stochastic approximation relies on the gradient flow equation associated with the vector field  $f'$ :  $\dot{x}_t = -f'(x_t)$ . This approach was introduced by [30] and [27] and have been applied in numerous papers since then, see [35, 36, 7, 6, 55]. We use to establish our result in Section 3.3, the strong connection between SGD and the gradient flow equation as well. The combination of the relation between stochastic approximation algorithms with the gradient flow equation and the Markov chain theory have been developed in [20] and [21]. In particular, [21] establishes under certain conditions that there exists for all  $\gamma \in (0, \gamma_0)$ , with  $\gamma_0$  small enough, an invariant distribution  $\pi_\gamma$  for the Markov chain  $(\theta_k^{(\gamma)})_{k \in \mathbb{N}}$ , and  $(\pi_\gamma)_{\gamma \in (0, \gamma_0)}$  is tight. In addition, they show that any limiting distributions is invariant for the gradient flow associated with  $\nabla f$ . Note that their conditions and results are different from ours. In particular, we do not assume that  $(\theta_k^{(\gamma)})_{k \in \mathbb{N}}$  is Feller but require that  $f$  is strongly convex contrary to [21].

To the authors knowledge, the use of the Richardson-Romberg method for stochastic approximation has only been considered in [38] to recover the minimax rate for recursive estimation of time varying autoregressive process.



Several attempts have been made to improve convergence of SGD. [5] proposed an online Newton algorithm which converges in practice to the optimal point with constant step-size but has no convergence guarantees. The quadratic case was studied by [5], for the (uniform) average iterate: the variance term is upper bounded by  $\sigma^2 d/n$  and the squared bias term by  $\|\theta^*\|^2/(\gamma n)$ . This last term was improved to  $\|\Sigma^{-1/2}\theta^*\|^2/(\gamma n)^2$  by [14, 15], showing that asymptotically, the bias term is negligible, see also [28]. Analysis has been extended to “tail averaging” [25], to improve the dependence on the initial conditions. Note that this procedure can be seen as a Richardson-Romberg trick with respect to  $k$ . Other strategies were suggested to improve the speed at which initial conditions were forgotten, for example using acceleration when the noise is additive [17, 26]. A criterion to check when SGD with constant step size is close to its limit distribution was recently proposed in [11].

In the context of discretization of ergodic diffusions, weak error estimates between the stationary distribution of the discretization and the invariant distribution of the associated diffusion have been first shown by [56] and [34] in the case of the Euler-Maruyama scheme. Then, [56] suggested the use of Richardson-Romberg interpolation to improve the accuracy of estimates of integrals with respect to the invariant distribution of the diffusion. Extension of these results have been obtained for other types of discretization by [1] and [12]. We show in Section 3.3 that a weak error expansion in the step-size  $\gamma$  also holds for SGD between  $\pi_\gamma$  and  $\delta_{\theta^*}$ . Interestingly as to the Euler-Maruyama discretization, SGD has a weak error of order  $\gamma$ . In addition, [18] proposed and analyzed the use of Richardson-Romberg extrapolation applied to the stochastic gradient Langevin dynamics (SGLD) algorithm. This method introduced by [58] combines SGD and the Euler-Maruyama discretization of the Langevin diffusion associated to a target probability measure [13, 19]. Note that this method is however completely different from SGD, in part because Gaussian noise of order  $\gamma^{1/2}$  (instead of  $\gamma$ ) is injected in SGD which changes the overall dynamics.

Finally, it is worth mentioning [32, 33] which are interested in showing that the invariant measure of constant step-size SGD for an appropriate choice of the step-size  $\gamma$ , can be used as a proxy to approximate the target distribution  $\pi$  with density with respect to the Lebesgue measure  $e^{-f}$ . Note that the perspective and purpose of this paper is completely different since we are interested in optimizing the function  $f$  and not in sampling from  $\pi$ .

### 3 Detailed analysis

In this Section, we describe in detail our approach. A first step is to describe the existence of a unique stationary distribution  $\pi_\gamma$  for the Markov chain  $(\theta_k^{(\gamma)})_{k \geq 0}$  and the convergence of this Markov chain to  $\pi_\gamma$  in the Wasserstein distance of order 2.

**Limit distribution** We cast in this section SGD in the Markov chain framework and introduce basic notion related to this theory, see [37] for an introduction to this topic. Consider the Markov kernel  $R_\gamma$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  associated with SGD iterates  $(\theta_k^{(\gamma)})_{k \in \mathbb{N}}$ , *i.e.* for all  $k \in \mathbb{N}$  and  $A \in \mathcal{B}(\mathbb{R}^d)$ , almost surely  $R_\gamma(\theta_k, A) = \mathbb{P}(\theta_{k+1} \in A | \theta_k)$ , for all  $\theta_0 \in \mathbb{R}^d$  and  $A \in \mathcal{B}(\mathbb{R}^d)$ ,  $\theta \mapsto R_\gamma(\theta, A)$  is Borel measurable and  $R_\gamma(\theta_0, \cdot)$  is a probability measure on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ . For all  $k \in \mathbb{N}^*$ , we define the Markov kernel  $R_\gamma^k$  recursively by  $R_\gamma^1 = R_\gamma$  and for  $k \geq 1$ , for all  $\theta_0 \in \mathbb{R}^d$  and  $A \in \mathcal{B}(\mathbb{R}^d)$

$$R_\gamma^{k+1}(\theta_0, A) = \int_{\mathbb{R}^d} R_\gamma^k(\theta_0, d\theta) R_\gamma(\theta, A) .$$



For any probability measure  $\lambda$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ , we define the probability measure  $\lambda R_\gamma$  for all  $A \in \mathcal{B}(\mathbb{R}^d)$  by

$$\lambda R_\gamma^k(A) = \int_{\mathbb{R}^d} \lambda(d\theta) R_\gamma^k(\theta, A) .$$

By definition, for all probability measure  $\lambda$  on  $\mathcal{B}(\mathbb{R}^d)$  and  $k \in \mathbb{N}^*$ ,  $\lambda R_\gamma^k$  is the distribution of  $\theta_k^{(\gamma)}$  started from  $\theta_0$  drawn from  $\lambda$ . For any function  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}_+$  and  $k \in \mathbb{N}^*$ , define the measurable function  $R_\gamma^k \phi : \mathbb{R}^d \rightarrow \mathbb{R}$  for all  $\theta_0 \in \mathbb{R}^d$ ,

$$R_\gamma^k \phi(\theta_0) = \int_{\mathbb{R}^d} \phi(\theta) R_\gamma^k(\theta_0, d\theta) .$$

For any measure  $\lambda$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  and any measurable function  $h : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $\lambda(h)$  denotes  $\int_{\mathbb{R}^d} h(\theta) d\lambda(\theta)$  when it exists. Note that with such notations, for any  $k \in \mathbb{N}^*$ , probability measure  $\lambda$  on  $\mathcal{B}(\mathbb{R}^d)$ , measurable function  $h : \mathbb{R}^d \rightarrow \mathbb{R}_+$ , we have  $\lambda(R_\gamma^k h) = (\lambda R_\gamma^k)(h)$ . A probability measure  $\pi_\gamma$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  is said to be a invariant probability measure for  $R_\gamma$ ,  $\gamma > 0$ , if  $\pi_\gamma R_\gamma = \pi_\gamma$ . A Markov chain  $(\theta_k^{(\gamma)})_{k \in \mathbb{N}}$  satisfying the SGD recursion (1) for  $\gamma > 0$  will be said at stationarity if it admits a invariant measure  $\pi_\gamma$  and  $\theta_k^{(\gamma)}$  is distributed according to  $\pi_\gamma$ . Note that in this case for all  $k \in \mathbb{N}$ , the distribution of  $\theta_k^{(\gamma)}$  is  $\pi_\gamma$ .

To show that  $(\theta_k^{(\gamma)})_{k \geq 0}$  admits a unique stationary distribution  $\pi_\gamma$  and quantify the convergence of  $(\nu_0 R_\gamma^k)_{k \geq 0}$  to  $\pi_\gamma$ , we use the Wasserstein distance. A probability measure  $\lambda$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  is said to have a finite second moment if  $\int_{\mathbb{R}^d} \|\vartheta\|^2 \lambda(d\vartheta) < +\infty$ . The set of probability measure on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  having a finite second moment is denoted by  $\mathcal{P}_2(\mathbb{R}^d)$ . For all probability measures  $\nu$  and  $\lambda$  in  $\mathcal{P}_2(\mathbb{R}^d)$ , define the *Wasserstein distance* of order 2 between  $\lambda$  and  $\nu$  by

$$W_2(\lambda, \nu) = \inf_{\xi \in \Pi(\lambda, \nu)} \left( \int \|x - y\|^2 \xi(dx, dy) \right)^{1/2} ,$$

where  $\Pi(\mu, \nu)$  is the set of probability measure  $\xi$  on  $\mathcal{B}(\mathbb{R}^d \times \mathbb{R}^d)$  satisfying for all  $A \in \mathcal{B}(\mathbb{R}^d)$ ,  $\xi(A \times \mathbb{R}^d) = \nu(A)$ ,  $\xi(\mathbb{R}^d \times A) = \lambda(A)$ .

**Proposition 1.** *Assume Assumption 1-Assumption 2-Assumption 3-Assumption 4(2). For any step-size  $\gamma \in (0, 2/L)$ , the Markov chain  $(\theta_k^{(\gamma)})_{k \geq 0}$ , defined by the recursion (1), admits a unique stationary distribution  $\pi_\gamma \in \mathcal{P}_2(\mathbb{R}^d)$ . In addition*

(a) *for all  $\theta \in \mathbb{R}^d$ ,  $k \in \mathbb{N}^*$ :*

$$W_2^2(R_\gamma^k(\theta, \cdot), \pi_\gamma) \leq (1 - 2\mu\gamma(1 - \gamma L/2))^k \int_{\mathbb{R}^d} \|\theta - \vartheta\|^2 d\pi_\gamma(\vartheta) ;$$

(b) *for any Lipschitz function  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ , with Lipschitz constant  $L_\phi$ , for all  $\theta \in \mathbb{R}^d$ ,  $k \in \mathbb{N}^*$ :*

$$\left| R_\gamma^k \phi(\theta) - \pi_\gamma(\phi) \right| \leq L_\phi (1 - 2\mu\gamma(1 - \gamma L/2))^{k/2} \left( \int \|\theta - \vartheta\|^2 d\pi_\gamma(\vartheta) \right)^{1/2} .$$

*Proof.* Let  $\gamma \in (0, 2/L)$  and  $\lambda_1, \lambda_2 \in \mathcal{P}_2(\mathbb{R}^d)$ . By [57, Theorem 4.1], there exists a couple of random variables  $\theta_0^{(1)}, \theta_0^{(2)}$  such that  $W_2^2(\lambda_1, \lambda_2) = \mathbb{E}[\|\theta_0^{(1)} - \theta_0^{(2)}\|^2]$  independent of  $(\varepsilon_k)_{k \in \mathbb{N}^*}$ . Let

$(\theta_k^{(1)})_{k \geq 0}, (\theta_k^{(2)})_{k \geq 0}$  be the SGD iterates associated with the step-size  $\gamma$ , starting from  $\theta_0^{(1)}$  and  $\theta_0^{(2)}$  respectively and sharing the same noise, *i.e.* for all  $k \geq 0$ ,

$$\begin{cases} \theta_{k+1}^{(1)} &= \theta_k^{(1)} - \gamma[f'(\theta_k^{(1)}) + \varepsilon_{k+1}(\theta_k^{(1)})] \\ \theta_{k+1}^{(2)} &= \theta_k^{(2)} - \gamma[f'(\theta_k^{(2)}) + \varepsilon_{k+1}(\theta_k^{(2)})] \end{cases} . \quad (11)$$

Note that using that  $\theta_0^{(1)}, \theta_0^{(2)}$  are independent of  $\varepsilon_1$ , we have for  $i, j \in \{1, 2\}$  using Assumption 3, that

$$\mathbb{E}[\langle \theta_0^{(i)}, \varepsilon(\theta_0^{(j)}) \rangle] = 0 . \quad (12)$$

Since for all  $k \geq 0$ , the distribution of  $(\theta_k^{(1)}, \theta_k^{(2)})$  belongs to  $\Pi(\lambda_1 R_\gamma^k, \lambda_2 R_\gamma^k)$ , by definition of the Wasserstein distance we get

$$\begin{aligned} W_2^2(\lambda_1 R_\gamma, \lambda_2 R_\gamma) &\leq \mathbb{E} \left[ \|\theta_1^{(1)} - \theta_1^{(2)}\|^2 \right] \\ &\leq \mathbb{E} \left[ \|\theta_0^{(1)} - \gamma f'_1(\theta_0^{(1)}) - (\theta_0^{(2)} - \gamma f'_1(\theta_0^{(2)}))\|^2 \right] \\ &\stackrel{i)}{\leq} \mathbb{E} \left[ \left\| \theta_0^{(1)} - \theta_0^{(2)} \right\|^2 - 2\gamma \left\langle f'(\theta_0^{(1)}) - f'(\theta_0^{(2)}), \theta_0^{(1)} - \theta_0^{(2)} \right\rangle \right] \\ &\quad + \gamma^2 \mathbb{E} \left[ \left\| f'_1(\theta_0^{(1)}) - f'_1(\theta_0^{(2)}) \right\|^2 \right] \\ &\stackrel{ii)}{\leq} \mathbb{E} \left[ \left\| \theta_0^{(1)} - \theta_0^{(2)} \right\|^2 - 2\gamma(1 - \gamma L/2) \left\langle f'(\theta_0^{(1)}) - f'(\theta_0^{(2)}), \theta_0^{(1)} - \theta_0^{(2)} \right\rangle \right] \\ &\stackrel{iii)}{\leq} (1 - 2\mu\gamma(1 - \gamma L/2)) \mathbb{E} \left[ \left\| \theta_0^{(1)} - \theta_0^{(2)} \right\|^2 \right] , \end{aligned}$$

using (12) for  $i$ ), Assumption 4(2) for  $ii$ ), and finally Assumption 1 for  $iii$ ).

Thus by a straightforward induction, we get, setting  $\rho = (1 - 2\mu\gamma(1 - \gamma L/2))$

$$\begin{aligned} W_2^2(\lambda_1 R_\gamma^k, \lambda_2 R_\gamma^k) &\leq \mathbb{E} \left[ \|\theta_k^{(1)} - \theta_k^{(2)}\|^2 \right] \\ &\leq \rho \mathbb{E} \left[ \|\theta_{k-1}^{(1)} - \theta_{k-1}^{(2)}\|^2 \right] \leq \rho^k W_2^2(\lambda_1, \lambda_2) , \end{aligned} \quad (13)$$

Since by Assumption 2-Assumption 3-Assumption 4(2),  $\lambda_1 R_\gamma \in \mathcal{P}_2(\mathbb{R}^d)$ , taking  $\lambda_2 = \lambda_1 R_\gamma$  in (13), for any  $N \in \mathbb{N}^*$ , we have  $\sum_{k=1}^N W_2^2(\lambda_1 R_\gamma^k, \lambda_2 R_\gamma^k) \leq \sum_{k=1}^N \rho^k W_2^2(\lambda_1, \lambda_1 R_\gamma)$ . Therefore, we get  $\sum_{k=1}^{+\infty} W_2^2(\lambda_1 R_\gamma^k, \lambda_1 R_\gamma^{k+1}) < +\infty$ . By [57, Theorem 6.16], the space  $\mathcal{P}_2(\mathbb{R}^d)$  endowed with  $W_2$  is a Polish space. Then,  $(\lambda_1 R_\gamma^k)_{k \geq 0}$  is a Cauchy sequence and converges to a limit  $\pi_\gamma^{\lambda_1} \in \mathcal{P}_2(\mathbb{R}^d)$ :

$$\lim_{k \rightarrow +\infty} W_2(\lambda_1 R_\gamma^k, \pi_\gamma^{\lambda_1}) = 0 . \quad (14)$$

We show that the limit  $\pi_\gamma^{\lambda_1}$  does not depend on  $\lambda_1$ . Assume that there exists  $\pi_\gamma^{\lambda_2}$  such that  $\lim_{k \rightarrow +\infty} W_2(\lambda_2 R_\gamma^k, \pi_\gamma^{\lambda_2}) = 0$ . By the triangle inequality

$$W_2(\pi_\gamma^{\lambda_1}, \pi_\gamma^{\lambda_2}) \leq W_2(\pi_\gamma^{\lambda_1}, \lambda_1 R_\gamma^k) + W_2(\lambda_1 R_\gamma^k, \lambda_2 R_\gamma^k) + W_2(\pi_\gamma^{\lambda_2}, \lambda_2 R_\gamma^k) .$$

Thus by (13) and (14), taking the limits as  $k \rightarrow +\infty$ , we get  $W_2(\pi_\gamma^{\lambda_1}, \pi_\gamma^{\lambda_2}) = 0$  and  $\pi_\gamma^{\lambda_1} = \pi_\gamma^{\lambda_2}$ . The limit is thus the same for all initial distributions and is denoted by  $\pi_\gamma$ .

Moreover,  $\pi_\gamma$  is invariant for  $R_\gamma$ . Indeed for all  $k \in \mathbb{N}^*$ ,

$$W_2(\pi_\gamma R_\gamma, \pi_\gamma) \leq W_2(\pi_\gamma R_\gamma, \pi_\gamma R_\gamma^k) + W_2(\pi_\gamma R_\gamma^k, \pi_\gamma).$$

Using (13) and (14), we get taking  $k \rightarrow +\infty$ ,  $W_2(\pi_\gamma R_\gamma, \pi_\gamma) = 0$  and  $\pi_\gamma R_\gamma = \pi_\gamma$ . The fact that  $\pi_\gamma$  is the unique stationary distribution is straightforward by contradiction and using (13).

Taking  $\lambda_1 = \delta_\theta$ ,  $\lambda_2 = \pi_\gamma$ , using the invariance of  $\pi_\gamma$  and (13), we get (a).

Finally, if we take  $\lambda_1 = \delta_\theta$  and  $\lambda_2 = \pi_\gamma$ , using  $\pi_\gamma R_\gamma = \pi_\gamma$ , (13), and the Cauchy-Schwarz inequality, we have for any  $k \in \mathbb{N}^*$ :

$$\begin{aligned} \left| R_\gamma^k \phi(\theta) - \pi_\gamma(\phi) \right| &= \left| \mathbb{E} \left[ \phi(\theta_{k,\gamma}^{(1)}) - \phi(\theta_{k,\gamma}^{(2)}) \right] \right| \leq L_\phi \mathbb{E}^{1/2} \left[ \left\| \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\|^2 \right] \\ &\leq L_\phi (1 - 2\mu\gamma(1 - \gamma L/2))^{k/2} \left( \int \|\theta - \vartheta\|^2 d\pi_\gamma(\vartheta) \right)^{1/2}, \end{aligned}$$

which concludes the proof of (b).  $\square$

A consequence of Proposition 1 is that the expectation of  $\bar{\theta}_k^{(\gamma)}$  defined by (2) converges to  $\int_{\mathbb{R}^d} \vartheta d\pi_\gamma(\vartheta)$  as  $k$  goes to infinity at a rate of order  $O(k^{-1})$ , see Proposition 3 in Section 6.2.

### 3.1 Expansion of moments of $\pi_\gamma$ when $\gamma$ is in a neighborhood of 0

In this sub-section, we analyze the properties of the chain starting at  $\theta_0$  distributed according to  $\pi_\gamma$ . As a result, we prove that the mean of the stationary distribution  $\bar{\theta}_\gamma = \int_{\mathbb{R}^d} \vartheta d\pi_\gamma(\vartheta)$  is such that  $\bar{\theta}_\gamma = \theta^* + \gamma\Delta + O(\gamma^2)$ . Simple developments of Equation (1) at the equilibrium, result in expansions of the first two moments of the chain. It extends [45, 31] which showed that  $(\gamma^{-1/2}(\pi_\gamma - \delta_{\theta^*}))_{\gamma>0}$  converges in distribution to a normal law as  $\gamma \rightarrow 0$ .

**Quadratic case** When  $f$  is a quadratic function, *i.e.*  $f'$  is affine, we have the following result.

**Proposition 2.** Assume  $f = f_\Sigma$ ,  $f_\Sigma : \theta \mapsto \|\Sigma^{1/2}(\theta - \theta^*)\|^2/2$ , where  $\Sigma$  is a positive definite matrix, and Assumption 2-Assumption 3-Assumption 4(4). Let  $\gamma \in (0, 2/L)$ . Then, it holds  $\bar{\theta}_\gamma = \theta^*$ ,  $\Sigma \otimes I + I \otimes \Sigma - \gamma\Sigma \otimes \Sigma$  is invertible and

$$\int_{\mathbb{R}^d} (\theta - \theta^*)^{\otimes 2} \pi_\gamma(d\theta) = \gamma(\Sigma \otimes I + I \otimes \Sigma - \gamma\Sigma \otimes \Sigma)^{-1} \left[ \int_{\mathbb{R}^d} \mathcal{C}(\theta) \pi_\gamma(d\theta) \right],$$

where  $\bar{\theta}_\gamma$  and  $\mathcal{C}$  are given by (3) and (5) respectively, and  $\pi_\gamma$  is the invariant probability measure of  $R_\gamma$  given by Proposition 1.

The first part of the result, which highlights the crucial fact that for a quadratic function, the mean under the limit distribution is the optimal point, is easy to prove. Indeed, since  $\pi_\gamma$  is invariant for  $(\theta_k^{(\gamma)})_{k \geq 0}$ , if  $\theta_0^{(\gamma)}$  is distributed according to  $\pi_\gamma$ , then  $\theta_1^{(\gamma)}$  is distributed according to  $\pi_\gamma$  as well. Thus as  $\theta_1^{(\gamma)} = \theta_0^{(\gamma)} - \gamma f'(\theta_0^{(\gamma)}) + \gamma \varepsilon_1(\theta_0^{(\gamma)})$  taking expectations on both sides, we get  $\int_{\mathbb{R}^d} f'(\vartheta) d\pi_\gamma(\vartheta) = 0$ . For a quadratic function, whose gradient is linear:  $\int_{\mathbb{R}^d} f'(\vartheta) d\pi_\gamma(\vartheta) = f'(\bar{\theta}_\gamma) = 0$  and thus  $\bar{\theta}_\gamma = \theta^*$ . This implies that the averaged iterate converges to  $\theta^*$ , see *e.g.* [5]. The proof for the second expression is given in Section 6.3.

**General case** While the quadratic case led to particularly simple expressions, in general, we can only get a first order development of these expectations as  $\gamma \rightarrow 0$ . Note that it improved on [45], which shows a similar expansion but an error of order of  $O(\gamma^{3/2})$ .

**Theorem 1.** *Assume Assumption 1-Assumption 2-Assumption 3-Assumption 4( $6 \vee [2(k_\varepsilon + 1)]$ )-Assumption 5 and let  $\gamma \in (0, 2/L)$ . Then  $f''(\theta^*) \otimes I + I \otimes f''(\theta^*)$  is invertible and*

$$\bar{\theta}_\gamma - \theta^* = \gamma f''(\theta^*)^{-1} f'''(\theta^*) \mathbf{AC}(\theta^*) + O(\gamma^2) \quad (15)$$

$$\int_{\mathbb{R}^d} (\theta - \theta^*)^{\otimes 2} \pi_\gamma(d\theta) = \gamma \mathbf{AC}(\theta^*) + O(\gamma^2), \quad (16)$$

where

$$\mathbf{A} = (f''(\theta^*) \otimes I + I \otimes f''(\theta^*))^{-1}, \quad (17)$$

$\bar{\theta}_\gamma$  and  $\mathcal{C}$  are given by (3) and (5) respectively, and  $\pi_\gamma$  is the invariant probability measure of  $R_\gamma$  given by Proposition 1.

*Proof.* The proof is postponed to Section 6.4.  $\square$

This shows that  $\gamma \mapsto \bar{\theta}_\gamma$  is a differentiable function at  $\gamma = 0$ . The “drift”  $\bar{\theta}_\gamma - \theta^*$  can be understood as an additional error occurring because the function is non quadratic ( $f'''(\theta^*) \neq 0$ ) and the step-sizes are not decaying to zero. The mean under the limit distribution is at distance  $\gamma$  from  $\theta^*$ . In comparison, the final iterate oscillates in a sphere of radius proportional to  $\sqrt{\gamma}$ .

### 3.2 Expansion for a given $\gamma > 0$ when $k$ tends to $+\infty$

In this sub-section, we analyze the convergence of  $\bar{\theta}_k^{(\gamma)}$  to  $\bar{\theta}_\gamma$ , when  $k \rightarrow \infty$ , and the convergence of  $\mathbb{E}[\|\bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma\|^2]$  to 0. Under suitable conditions [23],  $\bar{\theta}_k^{(\gamma)}$  satisfies a central limit theorem:  $\{\sqrt{k}(\bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma)\}_{k \in \mathbb{N}^*}$  converges in law to a  $d$ -dimensional Gaussian distribution with zero-mean. However, this result is purely asymptotic and we propose a new tighter development that describes how the initial conditions are forgotten. We show that the convergence behaves similarly to the convergence in the quadratic case, where the expected squared distance decomposes as a sum of a bias term, that scales as  $k^{-2}$ , and a variance term, that scales as  $k^{-1}$ , plus linearly decaying residual terms. We also describe how the asymptotic bias and variance can be easily expressed as moments of solutions associated to several *Poisson equations*.

For any Lipschitz function  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^q$ , by Lemma 1 in Section 6.2, the function  $\psi_\gamma = \sum_{i=0}^{+\infty} \{R_\gamma^i \varphi - \pi_\gamma(\varphi)\}$  is well-defined, Lipschitz and satisfies  $\pi_\gamma(\psi_\gamma) = 0$ ,  $(\text{Id} - R_\gamma)\psi_\gamma = \varphi$ .  $\psi_\gamma$  will be referred to as the *Poisson solution* associated with  $\varphi$ . Consider the three following functions:

- $\psi_\gamma$  the Poisson solution associated to  $\varphi : \theta \mapsto \theta - \theta^*$ ,
- $\varpi_\gamma$  the Poisson solution associated to  $\theta \mapsto \psi_\gamma(\theta)$ ,
- $\chi_\gamma^1$  the Poisson solution associated to  $\theta \mapsto (\psi_\gamma(\theta))^{\otimes 2}$ ,
- $\chi_\gamma^2$  the Poisson solution associated to  $\theta \mapsto ((\psi_\gamma - \varphi)(\theta))^{\otimes 2}$ .

**Theorem 2.** *Assume Assumption 1-Assumption 2-Assumption 3-Assumption 4(4) and let  $\gamma \in (0, 1/(2L))$ . Then setting  $\rho = (1 - \gamma\mu)^{1/2}$ , for any starting point  $\theta_0 \in \mathbb{R}^d$ ,  $k \in \mathbb{N}^*$ ,*

$$\mathbb{E}[\bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma] = k^{-1}(\psi_\gamma(\theta_0) + O(\rho^k)),$$

$$\begin{aligned} \mathbb{E} \left[ \left( \bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma \right)^{\otimes 2} \right] &= k^{-1} \pi_\gamma \left( \psi_\gamma^{\otimes 2} - (\psi_\gamma - \varphi)^{\otimes 2} \right) \\ &\quad - k^{-2} \left[ \pi_\gamma \left( \varpi_\gamma \varphi^\top + \varphi \varpi_\gamma^\top \right) + \chi_\gamma^2(\theta_0) - \chi_\gamma^1(\theta_0) \right] + O(k^{-3}), \end{aligned}$$

where  $\bar{\theta}_k^{(\gamma)}$ ,  $\bar{\theta}_\gamma$  are given by (2) and (3) respectively, and  $\pi_\gamma$  is the invariant probability measure of  $R_\gamma$  given by Proposition 1.

Equation (2) is a sum of three terms: (i) a variance term, that scales as  $1/k$ , and does not depend on the initial distribution (but only on the asymptotic distribution  $\pi_\gamma$ ), and (ii) a bias term, which scales as  $1/k^2$ , and depends on the initial point  $\theta_0 \in \mathbb{R}^d$ , (iii) a non-positive residual term, which scales as  $1/k^2$ .

*Proof.* In order to give the intuition of the proof and to underline how the associated Poisson solutions are introduced, we here sketch the proof of the first result. By definition of  $\varphi : \theta \mapsto \theta - \theta^*$  and since  $\psi_\gamma$  satisfies  $(\text{Id} - R_\gamma)\psi_\gamma = \varphi$ , we have

$$\mathbb{E} \left[ \bar{\theta}_{k+1}^{(\gamma)} \right] - \theta^* = (k+1)^{-1} \sum_{i=0}^k (R_\gamma^i \varphi)(\theta_0) = \pi_\gamma(\varphi) + (k+1)^{-1} \psi_\gamma(\theta_0) + R_\gamma^{k+1} \psi_\gamma(\theta_0),$$

where we have used that

$$\sum_{i=0}^{\infty} R_\gamma^i (\varphi - \pi_\gamma(\varphi)) - R_\gamma^{k+1} \sum_{i=0}^{\infty} R_\gamma^i (\varphi - \pi_\gamma(\varphi)) = \psi_\gamma - R_\gamma^{k+1} \psi_\gamma.$$

Finally, we have that  $R_\gamma^k \psi_\gamma(\theta_0)$  converges to 0 at linear speed, using Proposition 1 and  $\pi_\gamma(\psi_\gamma) = 0$ .

The formal and complete proof of this result is postponed to Section 6.5.  $\square$

This result gives an exact closed form for the asymptotic bias and variance, for a fixed  $\gamma$ , as  $k \rightarrow \infty$ . Unfortunately, in the general case, it is neither possible to compute the Poisson solutions exactly, nor is it possible to prove a first order development of the limits as  $\gamma \rightarrow 0$ .

When  $f_\Sigma$  is a quadratic function, it is possible, for any  $\gamma > 0$ , to compute  $\psi_\gamma$  and  $\chi_\gamma^{1,2}$  explicitly; we get the following decomposition of the error, which exactly recovers the result of [2] or [14].

**Corollary 1.** *Assume that  $f$  is an objective function of a least-square regression problem, i.e. with the notations of Example 1,  $f = f_\Sigma$ ,  $\Sigma = \mathbb{E}[XX^\top]$ ,  $\varepsilon_k$  are defined by (6), and step-size  $\gamma \leq 1/r^2$ , with  $r$  defined by (8). Assume Assumption 1-Assumption 2-Assumption 3-Assumption 4(4). For any starting point  $\theta_0 \in \mathbb{R}^d$  :*

$$\begin{aligned} \mathbb{E} \bar{\theta}_k^{(\gamma)} - \theta^* &= (1/(k\gamma)) \Sigma^{-1} (\theta_0 - \theta^*) + O(\rho^k) \\ \mathbb{E} \left[ \left( \bar{\theta}_k^{(\gamma)} - \theta^* \right)^{\otimes 2} \right] &= (1/k) \Sigma^{-1} \left\{ \int_{\mathbb{R}^d} \mathcal{C}(\theta) d\pi_\gamma(\theta) \right\} \Sigma^{-1} \\ &\quad + (1/(k^2 \gamma^2)) \Sigma^{-1} \Omega [\varphi(\theta_0)^{\otimes 2} - \pi_\gamma(\varphi^{\otimes 2})] \Sigma^{-1} \\ &\quad - (1/(k^2 \gamma^2)) (\Sigma^{-2} \otimes \text{Id} + \text{Id} \otimes \Sigma^{-2}) \pi_\gamma(\varphi^{\otimes 2}) + O(k^{-3}). \end{aligned}$$

With  $\Omega = (\Sigma \otimes I + I \otimes \Sigma - \gamma \Sigma \otimes \Sigma)(\Sigma \otimes I + I \otimes \Sigma - \gamma \mathbf{T})^{-1}$ , and

$$\mathbf{T} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}, A \mapsto \mathbb{E} \left[ (X^\top A X) X X^\top \right]. \quad (18)$$

*Proof.* The proof is postponed to the supplementary paper [16], Section S3.  $\square$

The bound on the second order moment is composed of a variance term  $k^{-1}\Sigma^{-1}\pi_\gamma(\mathcal{C})\Sigma^{-1}$ , a bias term which decays as  $k^{-2}$ , and a non-positive residual term. Interestingly, the bias is 0 if we start under the limit distribution.

### 3.3 Continuous interpretation of SGD and weak error expansion

Under the stated assumptions on  $f$  and  $(\varepsilon_k)_{k \in \mathbb{N}^*}$ , we have analyzed the convergence of the stochastic gradient recursion (1). We here describe how this recursion can be seen as a noisy discretization of the following *gradient flow* equation, for  $t \in \mathbb{R}_+$ :

$$\dot{\theta}_t = -f'(\theta_t). \quad (19)$$

Note that since  $f'(\theta^*) = 0$  by definition of  $\theta^*$  and Assumption 1, then  $\theta^*$  is an equilibrium point of (19), i.e.  $\theta_t = \theta^*$  for all  $t \geq 0$  if  $\theta_0 = \theta^*$ . Under Assumption 2, (19) admits a unique solution on  $\mathbb{R}_+$  for any starting point  $\theta \in \mathbb{R}^d$ . Denote by  $(\varphi_t)_{t \geq 0}$  the flow of (19), defined for all  $\theta \in \mathbb{R}^d$  by  $(\varphi_t(\theta))_{t \geq 0}$  as the solution of (19) starting at  $\theta$ .

Denote by  $(\mathcal{A}, D(\mathcal{A}))$ , the *infinitesimal generator* associated with the flow  $(\varphi_t)_{t \geq 0}$  defined by

$$\begin{aligned} D(\mathcal{A}) &= \left\{ h : \mathbb{R}^d \rightarrow \mathbb{R} : \text{for all } \theta \in \mathbb{R}^d, \lim_{t \rightarrow 0} \frac{h(\varphi_t(\theta)) - h(\theta)}{t} \text{ exists} \right\} \\ \mathcal{A}h(\theta) &= \lim_{t \rightarrow 0} \frac{h(\varphi_t(\theta)) - h(\theta)}{t} \text{ for all } h \in D(\mathcal{A}), \theta \in \mathbb{R}^d. \end{aligned} \quad (20)$$

Note that for any  $h \in C^1(\mathbb{R}^d)$ ,  $h \in D(\mathcal{A})$ ,  $\mathcal{A}h = -\langle f', h' \rangle$ .

Under Assumption 1 and Assumption 2, for any locally Lipschitz function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  (extension to a function  $g : \mathbb{R}^d \rightarrow \mathbb{R}^q$  can easily be done considering all assumptions and results coordinatewise), denote by  $h_g$  the solution of the continuous Poisson equation defined for all  $\theta \in \mathbb{R}^d$  by  $h_g(\theta) = \int_0^\infty (g(\varphi_s(\theta)) - g(\theta^*)) ds$ . Note that  $h_g$  is well-defined by Lemma 10-b) in Section 6.7.1, since  $g$  is assumed to be locally Lipschitz. By (20), we have for all  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ , locally Lipschitz,

$$\mathcal{A}h_g(\theta) = g(\theta^*) - g(\theta). \quad (21)$$

Under regularity assumptions on  $g$  (see Theorem 5),  $h_g$  is continuously differentiable and therefore satisfies  $\langle f', h'_g \rangle = g - g(\theta^*)$ . The idea is then to make a Taylor expansion of  $h_g(\theta_{k+1}^{(\gamma)})$  around  $\theta_k^{(\gamma)}$  to express  $k^{-1} \sum_{i=1}^k g(\theta_i^{(\gamma)}) - g(\theta^*)$  as convergent terms involving the derivatives of  $h_g$ . For  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\ell, p \in \mathbb{N}$ ,  $\ell \geq 1$  consider the following assumptions.

**Assumption 6** ( $\ell, p$ ). *There exist  $a_g, b_g \in \mathbb{R}_+$  such that  $g \in C^\ell(\mathbb{R}^d)$  and for all  $\theta \in \mathbb{R}^d$  and  $i \in \{1, \dots, \ell\}$ ,  $\|g^{(i)}(\theta)\| \leq a_g \{\|\theta - \theta^*\|^p + b_g\}$ .*

**Theorem 3.** *Let  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfying Assumption 6(5, p) for  $p \in \mathbb{N}$ . Assume Assumption 1-Assumption 2-Assumption 3-Assumption 5. Furthermore, suppose that there exists  $q \in \mathbb{N}$  and  $C \geq 0$  such that for all  $\theta \in \mathbb{R}^d$ ,*

$$\mathbb{E} \left[ \|\varepsilon_1(\theta)\|^{p+k_\varepsilon+3} \right] \leq C(1 + \|\theta - \theta^*\|^q),$$

*and Assumption 4(2 $\tilde{p}$ ) holds for  $\tilde{p} = p + 3 + q \vee k_\varepsilon$ . Then there exists a constant  $\varsigma > 0$  only depending on  $\tilde{p}$  such that for all  $\gamma \in (0, 1/(\varsigma L))$ ,  $k \in \mathbb{N}^*$  and any starting point  $\theta_0 \in \mathbb{R}^d$  it holds that:*

$$\mathbb{E} \left[ k^{-1} \sum_{i=1}^k \left\{ g(\theta_i^{(\gamma)}) - g(\theta^*) \right\} \right] = (1/(k\gamma)) \left\{ h_g(\theta_0) - \mathbb{E} \left[ h_g(\theta_{k+1}^{(\gamma)}) \right] \right\} \\ + (\gamma/2) \text{tr} \left( h_g''(\theta^*) \mathcal{C}(\theta^*) \right) - (\gamma/k) A_1(\theta_0) - \gamma^2 A_2(\theta_0, k), \quad (22)$$

where  $\theta_k^{(\gamma)}$  is the Markov chain starting from  $\theta_0$  and defined by the recursion (1) and  $\mathcal{C}$  is given by (5). In addition for some constant  $C \geq 0$  independent of  $\gamma$  and  $k$ , we have

$$A_1(\theta_0) \leq C \left\{ 1 + \|\theta_0 - \theta^*\|^{\bar{p}} \right\}, \quad A_2(\theta_0, k) \leq C \left\{ 1 + \|\theta_0 - \theta^*\|^{\bar{p}} / k \right\}.$$

*Proof.* The proof is postponed to Section 6.7.  $\square$

First in the case where  $f'$  is linear, choosing for  $g$  the identity function, then  $h_{\text{Id}} = \int_0^{+\infty} \{\varphi_s - \theta^*\} ds = \Sigma^{-1}$ , and we get that the first term in (22) vanishes which is expected since in that case  $\bar{\theta}_\gamma = \theta^*$ . Second by Lemma 11-b), we recover the first expansion of Theorem 1 for arbitrary objective functions  $f$ . Finally note that for all  $q \in \mathbb{N}$ , under appropriate conditions, Theorem 3 implies that there exist constants  $C_1, C_2(\theta_0) \geq 0$  such that  $\mathbb{E} \left[ k^{-1} \sum_{i=1}^k \|\theta_i^{(\gamma)} - \theta^*\|^{2q} \right] = C_1 \gamma + C_2(\theta_0)/k + O(\gamma^2)$ .

### 3.4 Discussion

Classical proofs of convergence rely on another decomposition, originally proposed by [42] and used in recent papers analyzing the averaged iterate [4]. We here sketch the arguments of these decompositions, in order to highlight the main difference, namely the fact that the residual term is not well controlled when  $\gamma$  goes to zero in the classical proof.

**Classical decomposition** The starting point of this decomposition is to consider a Taylor expansion of  $f'(\theta_{k+1}^{(\gamma)})$  around  $\theta^*$ . For any  $k \in \mathbb{N}$ ,

$$f'(\theta_k^{(\gamma)}) = f''(\theta^*)(\theta_k^{(\gamma)} - \theta^*) + O\left(\|\theta_k^{(\gamma)} - \theta^*\|^2\right).$$

As a consequence, using the definition of the SGD recursion (1),

$$\begin{aligned} \theta_{k+1}^{(\gamma)} - \theta_k^{(\gamma)} &= -\gamma f'(\theta_k^{(\gamma)}) - \gamma \varepsilon_{k+1}(\theta_k^{(\gamma)}) \\ &= -\gamma f''(\theta^*)(\theta_k^{(\gamma)} - \theta^*) - \gamma \varepsilon_{k+1}(\theta_k^{(\gamma)}) + \gamma O\left(\|\theta_k^{(\gamma)} - \theta^*\|^2\right). \end{aligned}$$

Thus

$$f''(\theta^*)(\theta_k^{(\gamma)} - \theta^*) = \gamma^{-1}(-\theta_{k+1}^{(\gamma)} + \theta_k^{(\gamma)}) - \varepsilon_{k+1}(\theta_k^{(\gamma)}) + O\left(\|\theta_k^{(\gamma)} - \theta^*\|^2\right).$$

Averaging over the first  $k$  iterates yields:

$$\begin{aligned} (k+1) \left( \bar{\theta}_k^{(\gamma)} - \theta^* \right) &= \gamma^{-1} f''(\theta^*)^{-1} \left( \theta_0^{(\gamma)} - \theta_{k+1}^{(\gamma)} \right) - \sum_{i=0}^k f''(\theta^*)^{-1} \varepsilon_{i+1} \left( \theta_i^{(\gamma)} \right) \\ &\quad + \sum_{i=0}^k O\left(\|\theta_i^{(\gamma)} - \theta^*\|^2\right). \end{aligned} \quad (23)$$



The term on the right-hand part of Equation (23) is composed of a bias term (depending on the initial condition), a variance term, and a residual term. This residual term differentiates the general setting from the quadratic one (in which it does not appear, as the first order Taylor expansion of  $f'$  is exact). This decomposition has been used in [4] to prove upper bound on the error, but does not allow for a tight decomposition in powers of  $\gamma$  when  $\gamma \rightarrow 0$ . Indeed, the residual  $\theta_i^{(\gamma)} - \theta^*$  simply does not go to 0 when  $\gamma \rightarrow 0$ : on the contrary, the chain becomes ill-conditioned when  $\gamma = 0$ .

**New decomposition** Here, we use the fact that for a function  $g : \mathbb{R}^d \rightarrow \mathbb{R}^q$  regular enough, there exists  $h_g : \mathbb{R}^d \rightarrow \mathbb{R}^q$  satisfying, for any  $\theta \in \mathbb{R}^d$ :

$$h'_g(\theta)f'(\theta) = g(\theta) - g(\theta^*),$$

where  $h'_g(\theta) \in \mathbb{R}^{q \times d}$ , and  $f'(\theta) \in \mathbb{R}^d$ . The starting point is then a first order Taylor development of  $h_g(\theta_{k+1}^{(\gamma)})$  around  $\theta_k^{(\gamma)}$ . For any  $k \in \mathbb{N}^*$ , we have

$$\begin{aligned} h_g(\theta_{k+1}^{(\gamma)}) &= h_g(\theta_k^{(\gamma)}) + h'_g(\theta_k^{(\gamma)})(\theta_{k+1}^{(\gamma)} - \theta_k^{(\gamma)}) + O\left(\left\|\theta_{k+1}^{(\gamma)} - \theta_k^{(\gamma)}\right\|^2\right) \\ &= h_g(\theta_k^{(\gamma)}) - \gamma h'_g(\theta_k^{(\gamma)})f'(\theta_k^{(\gamma)}) - \gamma h'_g(\theta_k^{(\gamma)})\varepsilon_{k+1}(\theta_k^{(\gamma)}) + O\left(\left\|\theta_{k+1}^{(\gamma)} - \theta_k^{(\gamma)}\right\|^2\right) \\ &= h_g(\theta_k^{(\gamma)}) - \gamma(g(\theta_k^{(\gamma)}) - g(\theta^*)) - \gamma h'_g(\theta_k^{(\gamma)})\varepsilon_{k+1}(\theta_k^{(\gamma)}) + O\left(\left\|\theta_{k+1}^{(\gamma)} - \theta_k^{(\gamma)}\right\|^2\right). \end{aligned}$$

Thus reorganizing terms,

$$\begin{aligned} g(\theta_k^{(\gamma)}) - g(\theta^*) &= \gamma^{-1} \left\{ h_g(\theta_k^{(\gamma)}) - h_g(\theta_{k+1}^{(\gamma)}) \right\} \\ &\quad + h'_g(\theta_k^{(\gamma)})\varepsilon_{k+1}(\theta_k^{(\gamma)}) + \gamma^{-1} O\left(\left\|\theta_{k+1}^{(\gamma)} - \theta_k^{(\gamma)}\right\|^2\right). \end{aligned}$$

Finally, averaging over the first  $k$  iterations and taking  $g = \text{Id}$  give

$$\begin{aligned} (k+1) \left( \bar{\theta}_k^{(\gamma)} - \theta^* \right) &= \gamma^{-1} \left( h_{\text{Id}}(\theta_0^{(\gamma)}) - h_{\text{Id}}(\theta_{k+1}^{(\gamma)}) \right) + \sum_{i=0}^k h'_{\text{Id}}(\theta_i^{(\gamma)})\varepsilon_{i+1} \left( \theta_i^{(\gamma)} \right) \\ &\quad + \gamma^{-1} \sum_{i=0}^k O\left(\left\|\theta_{i+1}^{(\gamma)} - \theta_i^{(\gamma)}\right\|^2\right). \end{aligned} \tag{24}$$

This expansion is the root of the proof of Theorem 3, which formalizes the expansion as powers of  $\gamma$ . The key difference between decomposition (23) and (24) is that in the latter, when  $\gamma \rightarrow 0$ , the expectation of the residual term tends to 0 and can naturally be controlled.

## 4 Experiments

We performed experiments on simulated data, for logistic regression, with  $n = 10^7$  observations, for  $d = 12$  and 4. Results are presented in Figure 2. The data are a.s. bounded by  $R \geq 0$ , therefore  $R^2 = L$ . We consider SGD with constant step-sizes  $1/R^2$ ,  $1/2R^2$  (and  $1/4R^2$ ) with or without averaging, with  $R^2 = L$ . Without averaging, the chain saturates with an error proportional to  $\gamma$

(since  $\|\theta_k^{(\gamma)} - \theta^*\| = O(\sqrt{\gamma})$  as  $k \rightarrow +\infty$ ). Note that the ratio between the convergence limits of the two sequences is roughly 2 in the un-averaged case, and 4 in the averaged case, which confirms the predicted limits. We consider Richardson Romberg iterates, which saturate at a much lower level, and performs much better than decaying step-sizes (as  $1/\sqrt{n}$ ) on the first iterations, as it forgets the initial conditions faster. Finally, we run the online-Newton [5], which performs very well but has no convergence guarantee. On the Right plot, we also propose an estimator that uses 3 different step-sizes to perform a higher order interpolation. More precisely, for all  $k \in \mathbb{N}^*$ , we compute  $\tilde{\theta}_k^3 = \frac{8}{3}\bar{\theta}_k^{(\gamma)} - 2\bar{\theta}_k^{(2\gamma)} + \frac{1}{3}\bar{\theta}_k^{(4\gamma)}$ . With such an estimator, the *first 2* terms in the expansion, scaling as  $\gamma$  and  $\gamma^2$ , should vanish, which explains that it does not saturate.

## 5 Conclusion

In this paper, we have used and developed Markov chain tools to analyze the behavior of constant step-size SGD, with a complete analysis of its convergence, outlining the effect of initial conditions, noise and step-sizes. For machine learning problems, this allows us to extend known results from least-squares to all loss functions. This analysis leads naturally to using Romberg-Richardson extrapolation, that provably improves the convergence behavior of the averaged SGD iterates. Our work opens up several avenues for future work: (a) show that Richardson-Romberg trick can be applied to the decreasing step-sizes setting, (b) study the extension of our results under self-concordance condition [3].

## References

- [1] A. Abdulle, G. Vilmart, and K. C. Zygalakis. High order numerical approximation of the invariant measure of ergodic SDEs. *SIAM J. Numer. Anal.*, 52(4):1600–1622, 2014.
- [2] R. Aguech, E. Moulines, and P. Priouret. On a perturbation approach for the analysis of stochastic tracking algorithms. *SIAM J. Control Optim.*, 39(3):872–899, 2000.
- [3] F. Bach. Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic regression. *J. Mach. Learn. Res.*, 15(1):595–627, January 2014.
- [4] F. Bach and E. Moulines. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Proceedings of the 24th International Conference on Neural Information Processing Systems*, NIPS’11, pages 451–459, USA, 2011. Curran Associates Inc.
- [5] F. Bach and E. Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate  $O(1/n)$ . *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- [6] M. Benaïm. A dynamical system approach to stochastic approximations. *SIAM J. Control Optim.*, 34(2):437–472, 1996.
- [7] A. Benveniste, M. Métivier, and P. Priouret. *Adaptive algorithms and stochastic approximations*, volume 22 of *Applications of Mathematics (New York)*. Springer-Verlag, Berlin, 1990. Translated from the French by Stephen S. Wilson.
- [8] D. Bertsekas. *Nonlinear programming*. Athena Scientific, 1995.
- [9] L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2008.
- [10] C. Bouton and G. Pagès. About the multidimensional competitive learning vector quantization algorithm with constant gain. *Ann. Appl. Probab.*, 7(3):679–710, 1997.
- [11] J. Chee and P. Toulis. Convergence diagnostics for stochastic gradient descent with constant step size. *arXiv preprint arXiv:1710.06382*, 2017.

- [12] C. Chen, N. Ding, and L. Carin. On the convergence of stochastic gradient MCMC algorithms with high-order integrators. In *NIPS*, pages 2269–2277, 2015.
- [13] A. S. Dalalyan. Theoretical guarantees for approximate sampling from smooth and log-concave densities. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3):651–676, 2017.
- [14] A. Défossez and F. Bach. Averaged least-mean-squares: bias-variance trade-offs and optimal sampling distributions. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2015.
- [15] A. Dieuleveut and F. Bach. Nonparametric stochastic approximation with large step-sizes. *Ann. Statist.*, 44(4):1363–1399, 08 2016.
- [16] A. Dieuleveut, A. Durmus, and F. Bach. Supplement to bridging the gap between constant step size stochastic gradient descent and Markov chains, 2018.
- [17] A. Dieuleveut, N. Flammarion, and F. Bach. Harder, better, faster, stronger convergence rates for least-squares regression. *ArXiv e-prints*, February 2016.
- [18] A. Durmus, U. Şimşekli, E. Moulines, R. Badeau, and G. Richard. Stochastic gradient Richardson-Romberg Markov Chain Monte Carlo. In *Advances in Neural Information Processing Systems*, pages 2047–2055, 2016.
- [19] A. Durmus and É. Moulines. Nonasymptotic convergence analysis for the unadjusted Langevin algorithm. *Ann. Appl. Probab.*, 27(3):1551–1587, 06 2017.
- [20] J.-C. Fort and G. Pagès. Convergence of stochastic algorithms: from the Kushner-Clark theorem to the Lyapounov functional method. *Adv. in Appl. Probab.*, 28(4):1072–1094, 1996.
- [21] J.-C. Fort and G. Pagès. Asymptotic behavior of a Markovian stochastic algorithm with constant step. *SIAM J. Control Optim.*, 37(5):1456–1482, 1999.
- [22] M. I. Freidlin and A. D. Wentzell. *Random Perturbations of Dynamical Systems*. Number vol. 260 in Grundlehren der mathematischen Wissenschaften. Springer, 1998.
- [23] P. Glynn and S. Meyn. A Liapunov bound for solutions of the Poisson equation. *Ann. Probab.*, 24(2):916–931, 04 1996.
- [24] P. Hartman. *Ordinary Differential Equations: Second Edition*. Classics in Applied Mathematics. Society for Industrial and Applied Mathematics (SIAM, 3600 Market Street, Floor 6, Philadelphia, PA 19104), 1982.
- [25] P. Jain, S. M. Kakade, R. Kidambi, P. Netrapalli, and A. Sidford. Parallelizing stochastic approximation through mini-batching and tail-averaging. *ArXiv e-prints*, October 2016.
- [26] P. Jain, S. M. Kakade, R. Kidambi, P. Netrapalli, and A. Sidford. Accelerating Stochastic Gradient Descent. *arXiv preprint arXiv:1704.08227*, 2017.
- [27] H. J. Kushner and D. S. Clark. *Stochastic approximation methods for constrained and unconstrained systems*, volume 26 of *Applied Mathematical Sciences*. Springer-Verlag, New York-Berlin, 1978.
- [28] G. Lan. An optimal method for stochastic composite optimization. *Math. Program.*, 133(1-2, Ser. A):365–397, 2012.
- [29] E. Levy. Why do partitions occur in Faa di Bruno’s chain rule for higher derivatives? Technical Report 0602183, arXiv, February 2006.
- [30] L. Ljung. Analysis of recursive stochastic algorithms. *IEEE Trans. Automatic Control*, AC-22(4):551–575, 1977.
- [31] L. Ljung, G. Pflug, and H. Walk. *Stochastic approximation and optimization of random systems*. DMV Seminar. Birkhauser Verlag, Basel, Boston, 1992.
- [32] S. Mandt, M. Hoffman, and D. M. Blei. A variational analysis of stochastic gradient algorithms. In *International Conference on Machine Learning*, pages 354–363, 2016.
- [33] S. Mandt, M. D Hoffman, and D. M. Blei. Stochastic gradient descent as approximate Bayesian inference. *arXiv preprint arXiv:1704.04289*, 2017.

- [34] J. C. Mattingly, A. M. Stuart, and D. J. Higham. Ergodicity for SDEs and approximations: locally Lipschitz vector fields and degenerate noise. *Stochastic Process. Appl.*, 101(2):185–232, 2002.
- [35] M. Métivier and P. Priouret. Applications of a Kushner and Clark lemma to general classes of stochastic algorithms. *IEEE Trans. Inform. Theory*, 30(2, part 1):140–151, 1984.
- [36] M. Métivier and P. Priouret. Théorèmes de convergence presque sure pour une classe d’algorithmes stochastiques à pas décroissant. *Probab. Theory Related Fields*, 74(3):403–428, 1987.
- [37] S. Meyn and R. Tweedie. *Markov Chains and Stochastic Stability*. Cambridge University Press, New York, NY, USA, 2nd edition, 2009.
- [38] E. Moulines, P. Priouret, and F. Roueff. On recursive estimation for time varying autoregressive processes. *Ann. Statist.*, 33(6):2610–2654, 2005.
- [39] A. Nedić and D. Bertsekas. Convergence rate of incremental subgradient algorithms. In *Stochastic optimization: algorithms and applications*, pages 223–264. Springer, 2001.
- [40] D. Needell, R. Ward, and N. Srebro. Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm. In *Advances in Neural Information Processing Systems 27*, pages 1017–1025. Curran Associates, Inc., 2014.
- [41] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust Stochastic Approximation Approach to Stochastic Programming. *SIAM J. on Optimization*, 19(4):1574–1609, January 2009.
- [42] A. S. Nemirovsky and D. B. Yudin. *Problem complexity and method efficiency in optimization*. A Wiley-Interscience Publication. John Wiley & Sons, Inc., New York, 1983. Translated from the Russian and with a preface by E. R. Dawson, Wiley-Interscience Series in Discrete Mathematics.
- [43] Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Applied Optimization. Springer, 2004.
- [44] Y. Nesterov and J. Ph. Vial. Confidence level solutions for stochastic programming. *Automatica*, 44(6):1559–1568, June 2008.
- [45] G. Pflug. Stochastic minimization with constant step-size: asymptotic laws. *SIAM Journal on Control and Optimization*, 24(4):655–666, 1986.
- [46] B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM J. Control Optim.*, 30(4):838–855, 1992.
- [47] P. Priouret and A. Veretenikov. A remark on the stability of the l.m.s. tracking algorithm. *Stochastic Analysis and Applications*, 16(1):119–129, 1998.
- [48] A. Rakhlin, O. Shamir, and K. Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. *ArXiv e-prints*, September 2011.
- [49] H. Robbins and S. Monro. A stochastic approximation method. *The Annals of mathematical Statistics*, 22(3):400–407, 1951.
- [50] D. Ruppert. Efficient estimations from a slowly convergent Robbins-Monro process. Technical report, Cornell University Operations Research and Industrial Engineering, 1988.
- [51] S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Stochastic convex optimization. In *Proceedings of the International Conference on Learning Theory (COLT)*, 2009.
- [52] S. Shalev-Shwartz, Y. Singer, and N. Srebro. Pegasos: Primal Estimated sub-GrAdient SOLver for SVM. In *Proceedings of the 24th International Conference on Machine Learning, ICML ’07*, pages 807–814, New York, NY, USA, 2007. ACM.
- [53] O. Shamir and T. Zhang. Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- [54] J. Stoer and R. Bulirsch. *Introduction to numerical analysis*, volume 12. Springer Science & Business Media, 2013.

- [55] V. B. Tadić and A. Doucet. Asymptotic bias of stochastic gradient search. *Ann. Appl. Probab.*, 27(6):3255–3304, 2017.
- [56] D. Talay and L. Tubaro. Expansion of the global error for numerical schemes solving stochastic differential equations. *Stochastic Anal. Appl.*, 8(4):483–509 (1991), 1990.
- [57] C. Villani. *Optimal transport : old and new*. Grundlehren der mathematischen Wissenschaften. Springer, Berlin, 2009.
- [58] M. Welling and Y. W. Teh. Bayesian learning via Stochastic Gradient Langevin Dynamics. In *ICML*, pages 681–688, 2011.
- [59] D. L. Zhu and P. Marcotte. Co-coercivity and its role in the convergence of iterative schemes for solving variational inequalities. *SIAM Journal on Optimization*, 6(3):714–726, 1996.

## 6 Postponed proofs

### 6.1 Discussion on assumptions on the noise

Assumption [Assumption 4](#), made in the text, can be weakened in order to apply to settings where input observations are un-bounded (typically, Gaussian inputs would not satisfy Assumption [Assumption 4](#)). Especially, in many cases, we only need Assumption [Assumption 7](#) below. Let  $p \geq 2$ .

**Assumption 7 (p).** (i) *There exists  $\tilde{\tau}_p \geq 0$  such that  $\{\mathbb{E}^{1/p}[\|\varepsilon_1(\theta^*)\|^p]\} \leq \tilde{\tau}_p$ .*

(ii) *For all  $x, y \in \mathbb{R}^d$ , there exists  $L \geq 0$  such that, for  $q = 2, \dots, p$ ,*

$$\begin{aligned} \mathbb{E} [\|f'_1(x) - f'_1(y)\|^q] \\ \leq L^{q-1} \|x - y\|^{q-2} \langle x - y, f'(x) - f'(y) \rangle, \end{aligned} \quad (25)$$

where  $L$  is the same constant appearing in Assumption [Assumption 2](#) and  $f'_1$  is defined by [\(4\)](#).

On the other hand, we consider also the stronger assumption that the noise is independent of  $\theta$  (referred to as the “semi-stochastic” setting, see [\[17\]](#)), or more generally that the noise has a uniformly bounded fourth order moment.

**Assumption 8.** *There exists  $\tau \geq 0$  such that  $\sup_{\theta \in \mathbb{R}^d} \{\mathbb{E}^{1/4}[\|\varepsilon_1(\theta)\|^4]\} \leq \tau$ .*

Assumption [Assumption 7\(p\)](#),  $p \geq 2$ , is the weakest, as it is satisfied for random design least mean squares and logistic regression with bounded fourth moment of the inputs. Note that we do not assume that gradient or gradient estimates are a.s. bounded, to avoid the need for a constraint on the space where iterates live. It is straightforward to see that Assumption [7\(p\)](#),  $p \geq 2$ , implies Assumption [4\(p\)](#) with  $\tau_p = \tilde{\tau}_p$ , and Assumption [8](#)-Assumption [2](#) implies Assumption [4\(4\)](#).

It is important to note that assuming Assumption [3](#)—especially that  $(\varepsilon_k)_{k \in \mathbb{N}^*}$  are i.i.d. random fields—*does not* imply Assumption [8](#). On the contrary, making *the semi stochastic assumption*, i.e. that the noise functions  $(\varepsilon_k(\theta_{k-1}))_{k \in \mathbb{N}^*}$  are i.i.d. vectors (e.g. satisfied if  $\varepsilon_k$  is constant as a function of  $\theta$ ), is a very strong assumption, and implies Assumption [8](#).

## 6.2 Preliminary results

We preface the proofs of the main results by some technical lemmas.

**Lemma 1.** *Assume Assumption 1-Assumption 2-Assumption 3-Assumption 4(2). Let  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$  be a  $L_\phi$ -Lipschitz function. For any step-size  $\gamma \in (0, 2/L)$ , the function  $\psi_\gamma : \mathbb{R}^d \rightarrow \mathbb{R}$  defined for all  $\theta \in \mathbb{R}^d$  by*

$$\psi_\gamma(\theta) = \sum_{i=0}^{+\infty} R_\gamma^i \phi(\theta), \quad (26)$$

*is well-defined, Lipschitz and satisfies  $(\text{Id} - R_\gamma)\psi_\gamma = \phi$ ,  $\pi_\gamma(\psi_\gamma) = 0$ . In addition, if  $\tilde{\psi}_\gamma : \mathbb{R}^d \rightarrow \mathbb{R}$  is an other Lipschitz function satisfying  $(\text{Id} - R_\gamma)\tilde{\psi}_\gamma = \phi$ ,  $\pi_\gamma(\tilde{\psi}_\gamma) = 0$ , then  $\psi_\gamma = \tilde{\psi}_\gamma$ .*

*Proof.* Let  $\gamma \in (0, 2/L)$ . By Proposition 1-(b), for any Lipschitz continuous function  $\phi$ ,  $\{\theta \mapsto \sum_{i=1}^k (R_\gamma^i \phi(\theta) - \pi_\gamma(\phi))\}_{k \geq 0}$  converges absolutely on all compact sets of  $\mathbb{R}^d$ . Therefore  $\psi_\gamma$  given by (26) is well-defined. Let  $(\theta, \vartheta) \in \mathbb{R}^d \times \mathbb{R}^d$ . Consider now the two processes  $(\theta_k^{(1)})_{k \geq 0}, (\theta_k^{(2)})_{k \geq 0}$  defined by (11) with  $\lambda_1 = \delta_\theta$  and  $\lambda_2 = \delta_\vartheta$ . Then, for any  $k \in \mathbb{N}^*$ , using (13):

$$\begin{aligned} |R_\gamma^k \phi(\theta) - R_\gamma^k \phi(\vartheta)| &\leq L_\phi \mathbb{E}^{1/2} \left[ \left\| \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\|^2 \right] \\ &\leq L_\phi (1 - 2\mu\gamma(1 - \gamma L/2))^{k/2} \|\theta - \vartheta\|. \end{aligned} \quad (27)$$

Therefore by definition (26),  $\psi_\gamma$  is Lipschitz continuous. Finally, it is straightforward to verify that  $\psi_\gamma$  satisfies the stated properties.

If  $\tilde{\psi}_\gamma : \mathbb{R}^d \rightarrow \mathbb{R}$  is an other Lipschitz function satisfying these properties, we have for all  $\theta \in \mathbb{R}^d$ ,  $(\psi_\gamma - \tilde{\psi}_\gamma)(\theta) = R_\gamma(\psi_\gamma - \tilde{\psi}_\gamma)(\theta)$ . Therefore for all  $k \in \mathbb{N}^*$ ,  $\theta \in \mathbb{R}^d$ ,  $(\psi_\gamma - \tilde{\psi}_\gamma)(\theta) = R_\gamma^k(\psi_\gamma - \tilde{\psi}_\gamma)(\theta)$ . But by Proposition 1-(b),  $\lim_{k \rightarrow +\infty} R_\gamma^k(\psi_\gamma - \tilde{\psi}_\gamma)(\theta) = \pi_\gamma(\psi_\gamma - \tilde{\psi}_\gamma) = 0$ , which concludes the proof.  $\square$

**Lemma 2.** *Assume Assumption 1-Assumption 2-Assumption 3-Assumption 4(2). Then we have for any  $\gamma \in (0, 2/L)$ .*

$$\int_{\mathbb{R}^d} f'(\theta) \pi_\gamma(d\theta) = 0.$$

*Proof.* Let  $(\theta_k^{(\gamma)})_{k \in \mathbb{N}}$  be a Markov chain satisfying (1), with  $\theta_0^{(\gamma)}$  distributed according to  $\pi_\gamma$ . Then the proof follows from taking the expectation in (1) for  $k = 0$ , using that the distribution of  $\theta_1^{(\gamma)}$  is  $\pi_\gamma$ ,  $\mathbb{E}[\varepsilon_1(\theta)] = 0$  for all  $\theta \in \mathbb{R}^d$  and  $\varepsilon_1$  is independent of  $\theta_0^{(\gamma)}$ .  $\square$

**Lemma 3.** *Assume Assumption 1-Assumption 2-Assumption 3-Assumption 7(2). Then for any initial condition  $\theta_0^{(\gamma)} \in \mathbb{R}^d$ , we have for any  $\gamma > 0$ ,*

$$\mathbb{E} \left[ \left\| \theta_{k+1}^{(\gamma)} - \theta^* \right\|^2 \middle| \mathcal{F}_k \right] \leq (1 - 2\gamma\mu(1 - \gamma L)) \left\| \theta_k^{(\gamma)} - \theta^* \right\|^2 + 2\gamma^2 \tilde{\tau}_2^2,$$

*where  $(\theta_k^{(\gamma)})_{k \geq 0}$  is given by (1). Moreover, if  $\gamma \in (0, 1/L)$ , we have*

$$\int_{\mathbb{R}^d} \|\theta - \theta^*\|^2 \pi_\gamma(d\theta) \leq \gamma \tilde{\tau}_2^2 / (\mu(1 - \gamma L)). \quad (28)$$

*Proof.* The proof and result is very close to the ones from [40] but we extend it without a.s. Lipschitzness (Assumption 4) but with Assumption 7. Using Assumption 3-Assumption 1 and  $f'(\theta^*) = 0$ , we have

$$\begin{aligned} \mathbb{E} \left[ \left\| \theta_{k+1}^{(\gamma)} - \theta^* \right\|^2 \middle| \mathcal{F}_k \right] &\leq \left\| \theta_k^{(\gamma)} - \theta^* \right\|^2 + \gamma^2 \mathbb{E} \left[ \left\| f'_{k+1}(\theta_k^{(\gamma)}) \right\|^2 \middle| \mathcal{F}_k \right] \\ &\quad - 2\gamma \mathbb{E} \left[ \left\langle f'_{k+1}(\theta_k^{(\gamma)}) - f'_{k+1}(\theta^*), \theta_k^{(\gamma)} - \theta^* \right\rangle \middle| \mathcal{F}_k \right] \end{aligned} \quad (29)$$

$$\leq (1 - 2\mu\gamma) \left\| \theta_k^{(\gamma)} - \theta^* \right\|^2 + \gamma^2 \mathbb{E} \left[ \left\| f'_{k+1}(\theta_k^{(\gamma)}) \right\|^2 \middle| \mathcal{F}_k \right]. \quad (30)$$

In addition, under Assumption 3-Assumption 7(2) and using (4), we have:

$$\begin{aligned} \mathbb{E} \left[ \left\| f'_{k+1}(\theta_k^{(\gamma)}) \right\|^2 \middle| \mathcal{F}_k \right] &\leq 2 \left( \mathbb{E} \left[ \left\| f'_{k+1}(\theta_k^{(\gamma)}) - f'_{k+1}(\theta^*) \right\|^2 \middle| \mathcal{F}_k \right] + \mathbb{E} \left[ \left\| f'_{k+1}(\theta^*) \right\|^2 \middle| \mathcal{F}_k \right] \right) \\ &\leq 2 \left( \mathbb{E} \left[ \left\| f'_{k+1}(\theta_k^{(\gamma)}) - f'_{k+1}(\theta^*) \right\|^2 \middle| \mathcal{F}_k \right] + \tau^2 \right) \\ &\leq 2 \left( L \mathbb{E} \left[ \left\langle f'_{k+1}(\theta_k^{(\gamma)}) - f'_{k+1}(\theta^*), \theta_k^{(\gamma)} - \theta^* \right\rangle \middle| \mathcal{F}_k \right] + \tau^2 \right) \\ &\leq 2 \left( L \left\langle f'(\theta_k^{(\gamma)}) - f'(\theta^*), \theta_k^{(\gamma)} - \theta^* \right\rangle + \tau^2 \right). \end{aligned}$$

Combining this result and (30) concludes the proof of the first inequality.

Regarding the second bound, let a fixed initial point  $\theta_0^{(\gamma)} \in \mathbb{R}^d$ . By Jensen inequality and the first result we get for any  $k \in \mathbb{N}$  and  $M \geq 0$ ,

$$\begin{aligned} \mathbb{E} \left[ \left\| \theta_{k+1}^{(\gamma)} - \theta^* \right\|^2 \wedge M \right] &\leq (1 - 2\gamma\mu(1 - \gamma L))^{k+1} \left\| \theta_0^{(\gamma)} - \theta^* \right\|^2 \\ &\quad + 2\gamma^2 \tilde{\tau}_2^2 \sum_{i=0}^k (1 - 2\gamma\mu(1 - \gamma L))^i. \end{aligned}$$

Since by Proposition 1-(b),  $\lim_{k \rightarrow +\infty} \mathbb{E}[\left\| \theta_{k+1}^{(\gamma)} - \theta^* \right\|^2 \wedge M] = \int_{\mathbb{R}^d} \{\left\| \theta - \theta^* \right\|^2 \wedge M\} \pi_\gamma(d\theta)$ , we get for any  $M \geq 0$ ,

$$\int_{\mathbb{R}^d} \{\left\| \theta - \theta^* \right\|^2 \wedge M\} \pi_\gamma(d\theta) \leq \gamma \tilde{\tau}_2^2 / (\mu(1 - \gamma L)).$$

Taking  $M \rightarrow +\infty$  and applying the monotone convergence theorem concludes the proof.  $\square$

Using Lemma 3, we can extend Lemma 1 to functions  $\phi$  which are locally Lipschitz.

**Lemma 4.** Assume Assumption 1-Assumption 2-Assumption 3-Assumption 4(4). Let  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$  be a function satisfying there exists  $L_\phi \geq 0$  such that for any  $x, y \in \mathbb{R}^d$ ,

$$|\phi(x) - \phi(y)| \leq L_\phi \|x - y\| \{1 + \|x\| + \|y\|\}. \quad (31)$$

For any step-size  $\gamma \in (0, 1/L)$ , it holds:



(a) there exists  $C \geq 0$  such that for all  $\theta \in \mathbb{R}^d$ ,  $k \in \mathbb{N}^*$ :

$$\left| R_\gamma^k \phi(\theta) - \pi_\gamma(\phi) \right| \leq CL_\phi (1 - 2\mu\gamma(1 - \gamma L))^{k/2} \left\{ 1 + \|\theta - \theta^*\|^2 \right\} ;$$

(b) the function  $\psi_\gamma : \mathbb{R}^d \rightarrow \mathbb{R}$  defined for all  $\theta \in \mathbb{R}^d$  by (26) is well-defined satisfies  $(\text{Id} - R_\gamma)\psi_\gamma = \phi$ ,  $\pi_\gamma(\psi_\gamma) = 0$  and there exists  $L_\psi \geq 0$  such that for any  $x, y \in \mathbb{R}^d$ ,

$$|\psi(x) - \psi(y)| \leq L_\psi \|x - y\| \{1 + \|x\| + \|y\|\} . \quad (32)$$

*Proof.* In this proof,  $C \geq 0$  is a constant which can change from line to line.

(a) Let  $\gamma \in (0, 1/L)$ . Consider the two processes  $(\theta_k^{(1)})_{k \geq 0}, (\theta_k^{(2)})_{k \geq 0}$  defined by (11) with  $\lambda_1 = \delta_\theta$  and  $\lambda_2 = \pi_\gamma$ . Using (31), the Cauchy-Schwarz inequality,  $\pi_\gamma R_\gamma = \pi_\gamma$  and (13) we have for any  $k \in \mathbb{N}^*$ :

$$\begin{aligned} \left| R_\gamma^k \phi(\theta) - \pi_\gamma(\phi) \right|^2 &\leq \left| \mathbb{E} \left[ \phi(\theta_k^{(1)}) - \phi(\theta_k^{(2)}) \right] \right|^2 \\ &\leq L_\phi^2 \mathbb{E} \left[ \left\| \theta_k^{(1)} - \theta_k^{(2)} \right\|^2 \right] \mathbb{E} \left[ 1 + \left\| \theta_k^{(1)} \right\|^2 + \left\| \theta_k^{(2)} \right\|^2 \right] \\ &\leq CL_\phi^2 (1 - 2\mu\gamma(1 - \gamma L/2))^k \int \|\theta - \vartheta\|^2 d\pi_\gamma(\vartheta) \\ &\quad \times \left( 1 + (1 - 2\mu\gamma(1 - \gamma L))^k \|\theta - \theta^*\|^2 \right) , \end{aligned}$$

where we have Lemma 3 for the last inequality. Then the proof is concluded using for all  $x, y \in \mathbb{R}^d$ ,  $\|x + y\|^2 \leq 2(\|x\|^2 + \|y\|^2)$  and Lemma 3 again.

(b) Let  $\gamma \in (0, 1/L)$ . By (a),  $\{\theta \mapsto \sum_{i=1}^k (R_\gamma^i \phi(\theta) - \pi_\gamma(\phi))\}_{k \geq 0}$  converges absolutely on all compact sets of  $\mathbb{R}^d$ . Therefore  $\psi_\gamma$  given by (26) is well-defined. Let  $(\theta, \vartheta) \in \mathbb{R}^d \times \mathbb{R}^d$ . Consider now the two processes  $(\theta_{k,\gamma}^{(1)})_{k \geq 0}, (\theta_{k,\gamma}^{(2)})_{k \geq 0}$  defined by (11) with  $\lambda_1 = \delta_\theta$  and  $\lambda_2 = \delta_\vartheta$ . Then (31), the Cauchy-Schwarz inequality and (13), for any  $k \in \mathbb{N}^*$ , we get:

$$\begin{aligned} \left| R_\gamma^k \phi(\theta) - R_\gamma^k \phi(\vartheta) \right|^2 &\leq \left| \mathbb{E} \left[ \phi(\theta_k^{(1)}) - \phi(\theta_k^{(2)}) \right] \right|^2 \\ &\leq CL_\phi^2 (1 - 2\mu\gamma(1 - \gamma L))^{k/2} \|\theta - \vartheta\| \left\{ 1 + \|\theta\|^2 + \|\vartheta\|^2 \right\} . \end{aligned}$$

By definition (26),  $\psi_\gamma$  satisfies (32). Finally, it is straightforward to verify that  $\psi_\gamma$  satisfies the stated properties. □

It is worth pointing out that under Assumption 8 (the “semi-stochastic” assumption), a slightly different result holds. The following result underlines the difference between a stochastic noise and a semi-stochastic noise, especially the fact that the maximal step-size differs depending on this assumption made.

**Lemma 5.** Assume Assumption 1-Assumption 2-Assumption 3-Assumption 8. Then for any initial condition  $\theta_0^{(\gamma)} \in \mathbb{R}^d$ , we have for any  $\gamma \in (0, 2/(\mu + L)]$ ,

$$\mathbb{E} \left[ \left\| \theta_{k+1}^{(\gamma)} - \theta^* \right\|^2 \middle| \mathcal{F}_k \right] \leq (1 - 2\gamma\mu L/(\mu + L)) \left\| \theta_k^{(\gamma)} - \theta^* \right\|^2 + \gamma^2 \tau^2 ,$$

where  $(\theta_k^{(\gamma)})_{k \geq 0}$  is given by (1).

*Proof.* First, note that since  $f$  satisfies Assumption 1 and Assumption 2, by [43, Chapter 2, (2.1.24)], for all  $x, y \in \mathbb{R}^d$ ,

$$\left\langle f'(x) - f'(y), x - y \right\rangle \geq \frac{L\mu}{L + \mu} \|x - y\|^2 + \frac{1}{L + \mu} \|f'(x) - f'(y)\|^2. \quad (33)$$

Besides, under Assumption 8, we have:

$$\begin{aligned} \mathbb{E} \left[ \left\| f'_{k+1}(\theta_k^{(\gamma)}) \right\|^2 \middle| \mathcal{F}_k \right] &= \left\| f'(\theta_k^{(\gamma)}) \right\|^2 + \mathbb{E} \left[ \left\| f'_{k+1}(\theta_k^{(\gamma)}) - f'(\theta_k^{(\gamma)}) \right\|^2 \right] \\ &\leq \left\| f'(\theta_k^{(\gamma)}) \right\|^2 + \tau^2. \end{aligned}$$

So that finally, using (29), Assumption 3, (33), Assumption 2 and rearranging terms we get

$$\begin{aligned} \mathbb{E} \left[ \left\| \theta_{k+1}^{(\gamma)} - \theta^* \right\|^2 \middle| \mathcal{F}_k \right] &\leq (1 - 2\gamma\mu L/(\mu + L)) \left\| \theta_k^{(\gamma)} - \theta^* \right\|^2 + \gamma^2 \tau^2 \\ &\quad - 2 \frac{\gamma}{L + \mu} \left\| f'(\theta_k^{(\gamma)}) \right\|^2 + \gamma^2 \left\| f'(\theta_k^{(\gamma)}) \right\|^2. \end{aligned}$$

Using that  $\gamma \leq 2/(m + L)$  concludes the proof.  $\square$

We give uniform bound on the moments of the chain  $(\theta_k^{(\gamma)})_{k \geq 0}$  for  $\gamma > 0$ . For  $p \geq 1$ , recall that under Assumption 4(2p), the noise at optimal point has a moment of order  $2p$  and we denote

$$\tau_{2p} = \mathbb{E}^{1/2p} \left[ \left\| \varepsilon_1(\theta^*) \right\|^{2p} \right]. \quad (34)$$

We give a bound on the  $p$ -order moment of the chain, under the assumption that the noise has a moment of order  $2p$ .

For moment of order larger than 2, we have the following result.

**Lemma 6.** *Assume Assumption 1-Assumption 2-Assumption 3-Assumption 4(2p), for  $p \geq 1$ . There exist numerical constants  $C_p, D_p \geq 2$  that only depend on  $p$ , such that, if  $\gamma \in (0, 1/(LC_p))$ , for all  $k \in \mathbb{N}^*$  and  $\theta_0 \in \mathbb{R}^d$*

$$\mathbb{E}^{1/p} \left[ \left\| \theta_k^{(\gamma)} - \theta^* \right\|^{2p} \right] \leq (1 - 2\gamma\mu(1 - C_p\gamma L/2))^k \mathbb{E}^{1/p} \left[ \left\| \theta_0 - \theta^* \right\|^{2p} \right] + \frac{D_p \gamma \tau_{2p}^2}{\mu},$$

where  $(\theta_k^{(\gamma)})_{k \in \mathbb{N}}$  is defined by (1) with initial condition  $\theta_0^{(\gamma)} = \theta_0$ . Moreover, the following bound holds

$$\int_{\mathbb{R}^d} \left\| \theta - \theta^* \right\|^{2p} \pi_\gamma(d\theta) \leq (D_p \gamma \tau_{2p}^2 / \mu)^p. \quad (35)$$

**Remark.** • Notably, Lemma 6 implies that  $\int_{\mathbb{R}^d} \left\| \theta - \theta^* \right\|^4 \pi_\gamma(d\theta) = O(\gamma^2)$ , and thus  $\int_{\mathbb{R}^d} \left\| \theta - \theta^* \right\|^3 \pi_\gamma(d\theta) = O(\gamma^{3/2})$ . We also note that  $\int_{\mathbb{R}^d} \left\| \theta - \theta^* \right\|^2 \pi_\gamma(d\theta) = O(\gamma)$ , also implies by Jensen's inequality that  $\left\| \bar{\theta}_\gamma - \theta^* \right\|^2 = O(\gamma)$ .

- Note that there is no contradiction between (35) and Theorem 3, as for any  $p \geq 2$ , one has for  $g(\theta) = \left\| \theta - \theta^* \right\|^2$  and  $h_g$  the solution to the Poisson equation, that  $h_g''(\theta^*) = 0$ , so that the first term in the development (of order  $\gamma$ ) is indeed 0.