

# A LECTURE ON SGD AND MARKOV CHAINS

based on a paper by Dieuleveut, Durmus, and Bach (2020)

SIMONE MARIA GIANCOLA<sup>\*†</sup>

December 17, 2024

## CONTENTS

1	Introduction	1
1.I	Motivation and Setting	1
1.II	Summary of results	2
1.III	Related work	3
2	Results	3
2.I	Assumptions	3
2.II	A baby example	3
2.III	From general to specific and back	4
3	Interesting objects	6
4	Discussion	7
	References	7
A	Assumptions	9
A.I	Remarks on assumptions	9
B	Some proofs and other statements	10
C	Baby example details	13
D	Useful definitions and theorems	13

## 1 INTRODUCTION

In this short lecture, we will review the paper in the title. After presenting the framework, we will prove the starting statement, and just briefly comment the other ones, which have very long proofs. To conclude, we will overview the Richardson-Romberg extrapolation method, which is one direct application of the results, and the ideas behind the deeper theorems.

Emphasis is on intuition and quick understanding. References are to a minimum. Computations, when performed, are explicit. Throughout, we omit the full expressions of the theorems, but specify when these are available in the original publication.

**NOTATION** We use bold for vectors, and curly latex for random variables. For example  $\boldsymbol{\theta}$  is a random vector, while  $\boldsymbol{\theta}^*$  will be our deterministic optimum. The rest of the symbols are either standard or defined when discussed first. The main takeaway is this explicit distinction between what is random and what is not.

### 1.I *Motivation and Setting*

Stochastic gradient descent (SGD) is a standard tool in machine learning (ML). However, the choice of the step-size is classical only in the deterministic case, where the algorithm degrades to gradient descent (GD). If we train over random samples, there is a *gap* between theory and practice:

---

<sup>\*</sup>Université Paris-Saclay, Orsay

<sup>†</sup>email: simonegiancola09@gmail.com

- a scaling  $O(1/k)$  is advisable, while still being non-robust to ill-conditioning;
- experiments suggest  $O(1/\sqrt{k})$  and averaging.

We wish to tackle this issue by providing intuitions on a set of results that indeed *bridges the gap*. In the spirit of isolating phenomena, let us choose constant step-size. The immediate justifications are that: (i) it is easier theoretically; (ii) there are fewer parameters to optimize (one less); (iii) initial conditions are anyway forgotten exponentially fast for the problems we consider; (iv) in practice, it gets close to the global optimum, and in ML we do not care when we are already at machine precision.

Let us introduce the notation. In this report SGD is expressed via the recursion:

$$\theta_{k+1}^{(\gamma)} = \theta_k^{(\gamma)} - \gamma \left[ \nabla f(\theta_k^{(\gamma)}) + \epsilon_{k+1}(\theta_k^{(\gamma)}) \right], \quad \theta_0 \sim \lambda, \quad (1.1)$$

↑ new estimator      ↑ true gradient      ↑ noise

where  $\lambda$  is an initial distribution and  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is the function we wish to minimize (a loss function). It will be nice and admit a global minimum  $\theta^*$ .<sup>1</sup>

#### Main Observation

The sequence  $(\theta_k^{(\gamma)})_{k \in \mathbb{N}}$  is a homogeneous Markov chain.

To remedy these oscillations, we consider a running mean:

$$\bar{\theta}_k^{(\gamma)} = \frac{1}{k+1} \sum_{j=0}^k \theta_j^{(\gamma)}. \quad (1.2)$$

We will prove under appropriate conditions on the sequence of iterates that a CLT holds, so that the running mean converges to the mean of the stationary distribution at a rate  $O(1/k)$ . The mean of the stationary distribution is obviously:

$$\bar{\theta}_\gamma = \int_{\mathbb{R}^d} \xi d\pi_\gamma(\xi). \quad (1.3)$$

Therefore, we can split the deviation of this running mean from the real (global) optimum of the function into:

$$\|\bar{\theta}_k^{(\gamma)} - \theta^*\|_2^2 \leq \underbrace{\|\bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma\|_2^2}_{\text{stochastic error}} + \underbrace{\|\bar{\theta}_\gamma - \theta^*\|_2^2}_{\text{deterministic error}}. \quad (1.4)$$

#### 1.II Summary of results

We will briefly argue that for quadratic functions the second term is null, but in general the oscillations are of order  $O(\gamma)$ . Then, we will derive an **explicit asymptotic expansion** of these oscillations, i.e. an expansion of  $\bar{\theta}_\gamma - \theta^*$  in the parameters of the algorithm. Like in classic ML, through a bias-variance decomposition, a quantitative CLT to expand the stochastic part  $\mathbb{E} \left[ \|\bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma\|_2^2 \right]$  is derived. In particular the bias depends on the initial conditions and the variance depends on the structure of the noise. Building further on that, we will discuss a **non-asymptotic** expansion between the stationary distribution  $\pi_\gamma$  and the Dirac at  $\theta^*$  in terms of the step-size. This means in practice that for a nice class of functions  $g$ , the integral with respect to the stationary measure is decomposed in Taylor-style centered at  $g(\theta^*)$ . All together, these results will describe the invariant distribution and how it is reached from any measure. In doing so, we will also justify a nice numerical trick,

<sup>1</sup> The original paper is more general and can work with tensors, but we avoid this for the sake of simplicity.

the Richardson-Romberg extrapolation, and find the excuse to discuss how to recover concentration from implicit functions depending on the problem.

#### Advantages

- (A1) Asymptotic expansions are explicit in the parameters! Full quantitative phenomenology.
- (A2) We can build confidence intervals for  $\theta^*$  (Chen et al. 2020; Su and Zhu 2018).
- (A3) More informed design of automatic restart schemes.

### 1.III Related work

We keep references here to a minimum, and refer to the original publication for valuable comments.

The most direct relative of this work is (Aguech, Moulines, and Priouret 2000), which develops a theory that matches the results for linear regression. For the perspective of SGD as a discretized gradient flow and its connections to Markov chains theory there are many important works. In particular, Fort and Pagès (1999) find tightness of the invariant distribution in a neighborhood of  $\gamma$  small and invariance with respect to the gradient flow of the limit distribution. However, they assume the process is Feller, while this work rather puts assumptions on the objective function (namely, strong convexity).

For convergence of SGD, there are many works studying the bias variance trade-off, with an emphasis on the dependence on initial conditions. Concerning the perspective of comparing the discretized invariant distribution and the continuous version and the Richardson-Romberg trick there are earlier works that combined the two.

The first part of corollary 2.10 matches the result of Defossez and Bach 2015, while proposition 2.6 is in part an extension of the work of Ljung, Pflug, and Walk (2012), in the sense that it quantifies their result that

$$\sqrt{\gamma}(\pi_\gamma - \delta_{\theta^*})_{\gamma>0} \xrightarrow[\mathcal{D}]{\gamma \rightarrow 0} \mathcal{N}(0, 1).$$

## 2 RESULTS

### 2.I Assumptions

We present here a quick overview of the assumptions, and postpone explicit writing to the appendix A.

(on  $f$ )  $\mu$ -strong convexity, in  $C^5(\mathbb{R}^d, \mathbb{R})$  with uniform bounds,<sup>2</sup> with  $L$ -co-coercive derivative (see appendix D for a refresher).

(on  $(\epsilon_k)_{k \in \mathbb{N}}$ ) The sequence of noise terms is adapted wrt to a filtration and is a Markov chain wrt to it, the  $p^{\text{th}}$  norm of the noise is controlled by  $\tau_p$  and the covariance of the noise is  $C^3(\mathbb{R}^d, \mathbb{R})$  with operator norms bounded by  $M_\epsilon(1 + \|\theta - \theta^*\|_2^{k_\epsilon})$  when  $\epsilon \equiv \epsilon(\theta)$ .

(misc) The initial iterate is measurable  $\theta_0 \in \mathcal{F}_0$ , and we have access to unbiased estimates, but the noise may depend on the current iterate, i.e. for all  $\theta \in \mathbb{R}^d$ :

$$\nabla f_{k+1}(\theta) = \nabla f(\theta) + \epsilon_{k+1}, \quad , \epsilon_{k+1} \equiv \epsilon_{k+1}(\theta). \quad (2.1)$$

Here  $\nabla f$  is effectively the random gradient of  $f$  which is random and  $\nabla f$  is the true gradient at  $\theta$ . In particular, this implies that the noise may be not i.i.d. and this makes  $\theta_k \in \mathcal{F}_k$ .

### 2.II A baby example

Let us degrade the generality of the paper to a practical scenario, which is a special case discussed in (Dieuleveut, Durmus, and Bach 2020, example 1). Consider a loss  $\mathcal{L} : \mathbb{R}^d \times \mathbb{R} \times \mathbb{R}^d \rightarrow \mathbb{R}$  that maps the triplet of feature vector, observed signal and proposed latent vector  $(x, y, \theta)$  into a scalar. If we allow the features and the labels to be random, as when they come from a dataset  $\mathcal{D}$  of observations, then we are interested in the generalization loss  $\mathcal{E}(\theta; \mathcal{L}) := \mathbb{E}_{(x,y)} [\mathcal{L}(x, y, \theta)]$ . The classic way to propose an estimator iteratively is SGD,

<sup>2</sup> note this implies  $L$ -smoothness

for which it is well known that SGD is the discretized gradient flow (i.e. a “discrete derivative”) with respect to the loss function. Since we train with random iterates, the sequence  $(\theta_k^{(\gamma)})_{k \in \mathbb{N}}$  will be random itself. We provide more details in appendix C.

### 2.III From general to specific and back

We are now ready to present the main statements. We will not prove all of them but give intuition. The common observation is that all results attempt to take a clever asymptotic/limit. To begin, we show that the object of interest is well-defined. Throughout, we will take the assumptions as granted (let all of them hold in each statement for proper parameters). For further details, we refer to the original text.

**Proposition 2.2** (Prop.2 in (Dieuleveut, Durmus, and Bach 2020)). *Let  $\gamma \in (0, 2/L)$ . The iterations of SGD seen as a Markov chain admit a unique stationary distribution  $\pi_\gamma$  that has finite second moments, i.e.  $\pi_\gamma \in \mathcal{P}_2(\mathbb{R}^d)$ . We can also **quantify** the rate of convergence in two ways: via the Wasserstein distance over probability measures (def. D.16) and in terms of “nice” functions integrated out by the kernel over time. Mathematically for all  $\theta \in \mathbb{R}^d$  and  $k \in \mathbb{N}$ :*

$$W_2^2(\overset{\text{initial distribution}}{R_\gamma^k(\theta, \cdot)}, \overset{\text{invariant measure}}{\pi_\gamma}) \leq \overset{\text{rate function; goes exp fast to zero}}{[\hbar(L, \mu, \gamma)]^k} \int_{\mathbb{R}^d} \|\theta - \xi\|_2^2 d\pi_\gamma(\xi), \quad (2.3)$$

mean square distance in invariant measure (initial conditions)

and for all  $\theta \in \mathbb{R}^d$ ,  $k \in \mathbb{N}$  and  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$  Lipschitz with constant  $L_\phi$ :

$$|R_\gamma^k \phi(\theta) - \pi_\gamma(\phi)| \leq L_\phi [\hbar(L, \mu, \gamma)]^{k/2} \sqrt{\int_{\mathbb{R}^d} \|\theta - \xi\|_2^2 d\pi_\gamma(\xi)}. \quad (2.4)$$

Naturally, the function  $\hbar \equiv \hbar(L, \mu, \gamma, k)$  is a function that is less than one exactly because we have the stability condition on the step-size.

*Proof.* See appendix B. □

Since we have a stationary distribution, it makes sense to see how the statistics behave when we start already at the stationary distribution, and when we reach it from another. This is the main lesson of the next two statements.

**Corollary 2.5** (Prop. 16 in (Dieuleveut, Durmus, and Bach 2020)). *The running mean converges to the stationary mean at rate  $O(1/k)$ .*

**Proposition 2.6** (Prop. 3 and Thm. 4 in (Dieuleveut, Durmus, and Bach 2020), for  $\theta_0 \sim \pi_\gamma$  already stationary). *Consider first our baby example. Let  $\Sigma = \mathbb{E}[xx^\top]$  be positive definite and  $\gamma \in (0, 2/L)$ . Then we have:*

- $\bar{\theta}_\gamma = \theta^*$ , so the mean of the stationary distribution is aligned with the optimum;
- the deviations are error dependent:

$$\int_{\mathbb{R}^d} [\theta - \theta^*][\theta - \theta^*]^\top d\pi_\gamma(\theta) = \gamma \mathbf{R}(\Sigma, \mathbb{P}_{\epsilon(\theta)}, \pi_\gamma), \quad (2.7)$$

where  $\mathbf{R}$  is explicit. Moreover, for this easy case we can find a nice expression in terms of the multiplicative noise, see appendix B.

In general, we find that the distribution is not guessing the optimal point right, and we will have:

$$\bar{\theta}_\gamma - \theta^* = \gamma \mathbf{g}(f, \theta^*) + O(\gamma^2), \quad \int_{\mathbb{R}^d} [\theta - \theta^*][\theta - \theta^*]^\top d\pi_\gamma(\theta) = \gamma \mathbf{P}(f, \theta^*) + O(\gamma^2). \quad (2.8)$$

Again,  $\mathbf{g}, \mathbf{P}$  are explicit.

**Remark 2.9.** The distance mean-optimum and the oscillations are of order  $\gamma, \sqrt{\gamma}$  respectively.

**Theorem 2.10** (Cor. 6 and Thm. 5 in (Dieuleveut, Durmus, and Bach 2020)), for generic  $\theta_0$ . For our baby example, we have quantitative CLT when  $\gamma$  is small enough<sup>3</sup>  $\forall \theta_0 \in \mathbb{R}^d$ :

$$\mathbb{E} [\bar{\theta}_k^{(\gamma)}] - \theta^\star = \frac{1}{k\gamma} \Sigma^{-1}(\theta_0 - \theta^\star) + O(k^k) \quad (2.11)$$

$$\mathbb{E} [\bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma][\bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma]^\top = \frac{1}{k} \mathbf{U}(\Sigma, \mathbb{P}_{\epsilon(\theta)}, \pi_\gamma) + \frac{1}{k^2\gamma^2} \mathbf{Q}(\theta_0, \theta^\star, \Sigma, \pi_\gamma, \mathbb{P}_x) + \frac{1}{k^2\gamma^2} \mathbf{S}(\theta^\star, \Sigma, \pi_\gamma) + O(k^3). \quad (2.12)$$

The bold matrices are explicit. In the general case, the formulas are expressed implicitly in terms of **Poisson solution** of functions, but exhibit the same decomposition.

**Remark 2.13.** If we start  $\theta_0 \sim \pi_\gamma$ , then the variance term is zero in the baby example!

Lastly, we make use of a nice interpretation of SGD as a discretized gradient flow. It will allow us to eventually relate the *actual chain* to the optimal value.

**Theorem 2.14** (Thm. 7 in (Dieuleveut, Durmus, and Bach 2020)). Under appropriate additional assumptions, a nice class of functions  $g$  satisfies for any  $\theta_0 \in \mathbb{R}^d$ :

$$\mathbb{E} \left[ \frac{1}{k} \sum_{i=1}^k g(\theta_i^{(\gamma)}) - g(\theta^\star) \right] = \frac{1}{k\gamma} v(\theta_0, k+1, \gamma) + \frac{\gamma}{2} b(\theta^\star, \mathbb{P}_{\epsilon}, g) - \frac{\gamma}{k} A_1(\theta_0) - \gamma^2 A_2(\theta_0, k), \quad (2.15)$$

$$A_1(\theta_0) \leq C \left( 1 + \|\theta_0 - \theta^\star\|_2^{\tilde{p}} \right), \quad A_2(\theta_0, k) \leq C \left( 1 + \frac{\|\theta_0 - \theta^\star\|_2^{\tilde{p}}}{k} \right). \quad (2.16)$$

Here  $C > 0, \tilde{p}$  are constants, the latter being explicit and  $v, b$  are implicit scalars.

**Remark 2.17.** In particular, there exists  $C_1, C_2(\theta_0) \geq 0$  such that:

$$\mathbb{E} \left[ \frac{1}{k} \sum_{i=1}^k \|\theta_i^{(\gamma)} - \theta^\star\|_2^{2q} \right] = C_1 \gamma + \frac{1}{k} C_2(\theta_0) + O(\gamma^2). \quad (2.18)$$

#### Takeaway

Under appropriate assumptions:

- (R1) initial conditions are forgotten exponentially fast and any chain converges to a unique distribution;
- (R2) there is  $d$  such that for small enough step-size:

$$\bar{\theta}_\gamma = \int_{\mathbb{R}^d} \xi d\pi_\gamma(\xi) = \theta^\star + \gamma d + r_\gamma^{(1)}, \quad \|r_\gamma^{(1)}\| \leq C\gamma^2; \quad (2.19)$$

- (R3) the bias is expanded as:

$$\mathbb{E} [\bar{\theta}_k^{(\gamma)} - \theta^\star] = \frac{A(\theta_0, \gamma)}{k} + \gamma d + r_\gamma^{(2)}, \quad \|r_\gamma^{(2)}\|_2 \leq C(\gamma^2 + e^{-k\mu\gamma}); \quad (2.20)$$

- (R4) there is a quantitative CLT for the variance terms for fixed  $\gamma$  as  $k \rightarrow \infty$ :

$$\mathbb{E} [\|\bar{\theta}_k^{(\gamma)} - \bar{\theta}_\gamma\|^2] = B_1(\gamma) \frac{1}{k} + B_2(\gamma) \frac{1}{k^2} + O\left(\frac{1}{k^3}\right). \quad (2.21)$$

Combining variance and bias, we have characterized the trade-off in SGD.

**Remark 2.22.** Like the Euler-Maruyama scheme the weak error expansion of SGD in step-size  $\gamma$  between  $\pi_\gamma$  and  $\delta_{\theta^\star}$  is of order  $O(\gamma)$  (Talay and Tubaro 1990).

<sup>3</sup> Needs to satisfy a bound wrt to inflated empirically covariances and true covariances of  $\mathbb{P}_x$ , see (Dieuleveut, Durmus, and Bach 2020, eqn. 9)

We provide an overview of two interesting and less traditional aspects of the work. One is numerically oriented and quick to present, the other is rather deep.

**RICHARDSON-ROMBERG** Result (R2) above suggests using Richardson-Romberg extrapolation to decrease the bias of iterations. The idea is quick to understand. Recall that we have the weak error expansion between the stationary integral and the Dirac indicator at the true minimum, that is:

$$\int_{\mathbb{R}^d} g(\xi) d\pi_\gamma(\xi) = g(\theta^*) + \gamma C_1^g + r_\gamma^g, \quad \|r_\gamma^g\|_2 \leq C_2^g \gamma^2, \quad \text{for all } g \text{ "nice"}. \quad (3.1)$$

Therefore, if we run two chains at  $\gamma, 2\gamma$  step-size, we will have that  $\bar{\theta}_k^{(\gamma)} \rightarrow \bar{\theta}_\gamma, \bar{\theta}_k^{(2\gamma)} \rightarrow \bar{\theta}_{2\gamma}$  with distance to the optimal given by (choose  $g = \text{Id}$  the identity map):

$$\bar{\theta}_\gamma = \theta^* + \gamma d_1^{\text{Id}} + r_\gamma^{\text{Id}}, \quad \bar{\theta}_{2\gamma} = \theta^* + 2\gamma d_1^{\text{Id}} + r_{2\gamma}^{\text{Id}}, \quad \max \left\{ \|2r_\gamma^{\text{Id}}\|, \|r_{2\gamma}^{\text{Id}}\| \right\} \leq 2C\gamma^2. \quad (3.2)$$

Taking the combined chain running means:  $\left( 2\bar{\theta}_k^{(\gamma)} - \bar{\theta}_k^{(2\gamma)} \right)_{k \in \mathbb{N}}$  will then cancel out the linear  $O(\gamma)$  term and improve the distance to  $\theta^*$  up to factors of  $O(\gamma^2)$ .

**Remark 3.3.** *The cost is an increase in the variance which stays of the same order.*

**WHAT IS A POISSON SOLUTION?** Lastly, we try to give some intuition on generators and Poisson solutions. If we rescale time, we can see equation eq. (1.1) as a noisy gradient flow, which would be just  $\dot{\theta}_t = -\nabla f(\theta_t)$ . We could look at a set of functions that may represent a meaningful derivative of the flow of points. Denoting such flow as  $\varphi_t(\theta)$ , we define the infinitesimal generator  $\mathcal{A}$  as:

$$\mathcal{A}h(\theta) := \lim_{t \downarrow 0} \frac{1}{t} (h(\varphi_t(\theta)) - h(\theta)), \quad (3.4)$$

where we apply functions  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  that are such that the limit exists. For simplicity, we will say that  $\mathcal{A}$  has domain  $D(\mathcal{A})$  with these nice functions. From the definition of Markov property, we have that  $\mathbb{E}[h(\theta_{k+1}) | \mathcal{F}_k] = R_\gamma h(\theta_k)$ , and by the construction proposed we will also have the approximate relation:

$$\mathbb{E}[h(\theta_{t+s}) | \mathcal{F}_t] \approx h(\theta_t) + t\mathcal{A}h(\theta_t). \quad (3.5)$$

With this in hand, one can show that  $d/dt h(\varphi_t(\theta)) = \mathcal{A}h(\varphi_t(\theta))$  and that we can make the interpretation of a derivative in time to conclude that  $\varphi_t(\theta) = u(t, \theta)$  is a solution to the PDE:

$$\frac{\partial}{\partial t} u(t, \theta) = \mathcal{A}u(t, \theta). \quad (3.6)$$

In a somewhat more interesting perspective, from these observations it can be shown that:

$$\mathbf{m}_t^h := h(\theta_t) - h(\theta_0) - \int_0^t \mathcal{A}h(\theta_s) ds, \quad (3.7)$$

is a martingale, and thus obeys a central limit theorem. Let us see this through a quicker example that removes many details. Under our assumptions certain Poisson equations are well-behaved and there is a more direct link to at least an asymptotic form of the CLT. Given a kernel  $K_\gamma$  with a unique invariant distribution and a function  $h \in L^1(\pi)$ , we reorder the random sum:

$$s_n(h) := \sum_{k=0}^{n-1} h(\theta_k) = h(\theta_0) - R_\gamma h(\theta_0) + R_\gamma h(\theta_0) + h(\theta_{n-1}) - R_\gamma^2 h(\theta_0) - R_\gamma h(\theta_{n-1}) + \dots, \quad (3.8)$$

stressing that we might decompose each iteration steps into martingales. In this spirit, we say  $\hat{h}$  is a solution to the Poisson equation if (Douc et al. 2018, chaps. 22-23):

$$\hat{h} - R_\gamma \hat{h} = h - \pi(h), \quad \pi\text{-a.e.}, h \in L^1(R_\gamma). \quad (3.9)$$

The intuitive interpretation is as follows. The transition of the law of a Markov chain can be seen through the lenses of the (functional) heat equation  $\partial_t \hat{h} = (I - K_\gamma) \hat{h} = \Delta \hat{h}$ . If  $\Delta \hat{h} = 0$ , the function is said to be harmonic, while if  $\Delta \hat{h} = 0$  for all  $\hat{h}$ , the chain is trivially stationary. In this case, we give a “derivative”  $h$  and a kernel  $K_\gamma$  and want to find which function we were actually working on. A trivial substitution of the identity shows us that:

$$s_n(h) = \hat{h}(\theta_0) - \hat{h}(\theta_n) + \underbrace{\sum_{k=1}^n \hat{h}(\theta_k) - R_\gamma \hat{h}(\theta_{k-1})}_{m_n^h}, \quad (3.10)$$

where  $m_n^h$  is a martingale.<sup>4</sup> The normalization by  $1/\sqrt{n}$  and the CLT for martingales comes into rescue (Douc et al. 2018, chaps. 21-22). In appendix D we report a continuous time generalization. Now the problem is quantifying this convergence.

#### 4 DISCUSSION

While we have only scratched the surface, it is possible to see that under our rather restrictive assumption the theory is basically closed, as we have a full characterization of the stochastic nature of the phenomenon that is non-asymptotic. As common in the theory of optimization, going beyond the combination of  $L$ -smooth and  $\mu$ -strong convexity while retaining an explicit convergence rate is hard. One potential expansion in this direction could be modelling the quadratic nature of the function locally, at the cost of largely complicating the expressions. A potentially related theoretical principle is that of the resolvent method, which expands the concept of Poisson equation to a larger class of functions (see Douc et al. (2018, chaps. 21-22)). The real obstacle being the objective function, one could also slightly decrease the generality of the results and still inspect SGD under different function classes (as was done over the years), but still apply the generic tools present. In this regard, it is understandable that the result of uniqueness of the invariant distribution is striking and useful, therefore it would be crucial to see if it can be kept, or it has to be given up.

Lastly, there is a renewed interest on one-pass SGD for high-dimensional problems and its interplay with hardness of inference. There, the attempt is to classify functions based on the computational time needed to minimize them. While the setting is very different, Collins-Woodfin et al. (2023) use a resolvent method to transform a very implicit ODE into a solvable complex ODE, and express their solutions as a set of implicit equations.

#### REFERENCES

- Aguech, Rafik, Eric Moulines, and Pierre Priouret (Jan. 2000). “On a Perturbation Approach for the Analysis of Stochastic Tracking Algorithms”. In: *SIAM Journal on Control and Optimization* 39.3, pp. 872–899. ISSN: 0363-0129. DOI: 10.1137/S0363012998333852. (Visited on 12/11/2024).
- Chen, Xi et al. (Feb. 2020). “Statistical Inference for Model Parameters in Stochastic Gradient Descent”. In: *The Annals of Statistics* 48.1, pp. 251–273. ISSN: 0090-5364, 2168-8966. DOI: 10.1214/18-AOS1801. (Visited on 12/11/2024).
- Collins-Woodfin, Elizabeth et al. (2023). *Hitting the High-Dimensional Notes: An ODE for SGD learning dynamics on GLMs and multi-index models*. arXiv: 2308.08977 [math.OA]. URL: <https://arxiv.org/abs/2308.08977>.
- Defossez, Alexandre and Francis Bach (Feb. 2015). “Averaged Least-Mean-Squares: Bias-Variance Trade-offs and Optimal Sampling Distributions”. In: *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*. PMLR, pp. 205–213. (Visited on 12/11/2024).
- Dieuleveut, Aymeric, Alain Durmus, and Francis Bach (June 2020). “Bridging the Gap between Constant Step Size Stochastic Gradient Descent and Markov Chains”. In: *The Annals of Statistics* 48.3, pp. 1348–1382. ISSN: 0090-5364, 2168-8966. DOI: 10.1214/19-AOS1850. (Visited on 12/11/2024).
- Dieuleveut, Aymeric, Nicolas Flammarion, and Francis Bach (2017). “Harder, Better, Faster, Stronger Convergence Rates for Least-Squares Regression”. In: *Journal of Machine Learning Research* 18.101, pp. 1–51. ISSN: 1533-7928. (Visited on 12/11/2024).
- Douc, Randal et al. (2018). *Markov Chains*. Springer Series in Operations Research and Financial Engineering. Cham: Springer International Publishing. ISBN: 978-3-319-97703-4 978-3-319-97704-1. DOI: 10.1007/978-3-319-97704-1. (Visited on 12/11/2024).

<sup>4</sup> Just verify that  $\mathbb{E}[m_n^h | \mathcal{F}_{n-1}] - m_{n-1}^h = \mathbb{E}[\hat{h}(\theta_n) | \mathcal{F}_{n-1}] - R_\gamma \hat{h}(\theta_{n-1}) = 0$ .

- Fort, Jean-Claude and Gilles Pagès (Jan. 1999). “Asymptotic Behavior of a Markovian Stochastic Algorithm with Constant Step”. In: *SIAM Journal on Control and Optimization* 37.5, pp. 1456–1482. ISSN: 0363-0129. DOI: 10.1137/S0363012997328610. (Visited on 12/11/2024).
- Le Gall, Jean-François (2016). “General Theory of Markov Processes”. In: *Brownian Motion, Martingales, and Stochastic Calculus*. Ed. by Jean-François Le Gall. Cham: Springer International Publishing, pp. 151–184. ISBN: 978-3-319-31089-3. DOI: 10.1007/978-3-319-31089-3\_6. (Visited on 12/16/2024).
- Ljung, L., G. Pflug, and H. Walk (Dec. 2012). *Stochastic Approximation and Optimization of Random Systems*. Birkhäuser. ISBN: 978-3-0348-8609-3.
- Su, Weijie J. and Yuancheng Zhu (Aug. 2018). *Uncertainty Quantification for Online Learning and Stochastic Approximation via Hierarchical Incremental Gradient Descent*. DOI: 10.48550/arXiv.1802.04876. arXiv: 1802.04876 [stat]. (Visited on 12/11/2024).
- Talay, Denis and Luciano Tubaro (Jan. 1990). “Expansion of the Global Error for Numerical Schemes Solving Stochastic Differential Equations”. In: *Stochastic Analysis and Applications* 8.4, pp. 483–509. ISSN: 0736-2994. DOI: 10.1080/07362999008809220. (Visited on 12/11/2024).
- Villani, Cédric (2009). *Optimal Transport*. Ed. by M. Berger et al. Vol. 338. Grundlehren Der Mathematischen Wissenschaften. Berlin, Heidelberg: Springer. ISBN: 978-3-540-71049-3 978-3-540-71050-9. DOI: 10.1007/978-3-540-71050-9. (Visited on 12/11/2024).



## A ASSUMPTIONS

In this section, we rewrite the assumptions for completeness.

**Assumption A.1.** The function  $f$  is  $\mu$ -strongly convex, in the sense of definition D.1.

**Assumption A.2.** The functions  $f$  is five times differentiable with continuous and uniformly bounded derivatives. In particular, it is  $L$ -smooth.

**Assumption A.3.** There is a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  and a filtration  $(\mathcal{F}_k)_{k \in \mathbb{N}}$  for which:

- $\epsilon_k(\theta)$  is adapted for all  $\theta \in \mathbb{R}^d$ ;
- $\mathbb{E} [\epsilon_{k+1}(\theta) \mid \mathcal{F}_k] = 0$  for all  $\theta \in \mathbb{R}^d$ ;
- $(\epsilon_k)_{k \in \mathbb{N}}$  are i.i.d. random  $\mathbb{R}^d \mapsto \mathbb{R}^d$  functions (i.e. random fields);<sup>5</sup>
- $\theta_0 \in \mathcal{F}_0$ .

**Assumption A.4.** This condition depends on some  $p$  to be chosen. The function  $f_k$  is almost surely  $C$ -co-coercive (def. D.13) with  $C = L$  the smoothness constant. For the error, the  $p$  norm at the optimum is controlled as:

$$\left( \mathbb{E} [\|\epsilon_k(\theta^*)\|_2^p] \right)^{\frac{1}{p}} \leq \tau_p, \quad \text{for some } \tau_p.$$

**Assumption A.5.** The function  $C(\theta) := \mathbb{E} [\epsilon(\theta)\epsilon(\theta)^\top]$  is in  $C^3(\mathbb{R}^d, \mathbb{R}^{d \times d})$  and such that:

$$\max_{i \in \{1,2,3\}} \|C^{(i)}(\theta)\| \leq M_\epsilon \left( 1 + \|\theta - \theta^*\|_2^{k_\epsilon} \right), \quad \text{for some } k_\epsilon, M_\epsilon, \text{ for all } \theta \in \mathbb{R}^d,$$

where  $C^{(i)}$  are the (tensor) derivatives of the matrix  $C$ .

**Assumption A.6.** This condition depends on the tuple  $(\ell, p)$ . Let  $g$  be polynomially locally Lipschitz, in the sense of definition D.11. There are positive  $a_g, b_g$  such that  $g \in C^\ell(\mathbb{R}^d)$  and:

$$\|g^{(i)}(\theta)\| \leq a_g \left( b_g + \|\theta - \theta^*\|_2^p \right), \quad \forall \theta \in \mathbb{R}^d, i \in \{1, \dots, \ell\}, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d,$$

where  $g^{(i)}$  are the  $i^{\text{th}}$  derivatives of  $g$ , which become tensors for  $i > 2$ .

**Assumption A.7 (A.4 bis).** Alternatively, suppose:

- for some  $\tilde{\tau}_p \geq 0$  it holds that  $\left( \mathbb{E} [\|\epsilon(\theta)\|_2^p] \right)^{1/p} \leq \tilde{\tau}_p$ ;
- the smoothness constant  $L$  is such that:

$$\mathbb{E} [\|\nabla f_1(\mathbf{x}) - \nabla f_1(\mathbf{y})\|_2^q] \leq L^{q-1} \|\mathbf{x} - \mathbf{y}\|_2^{q-2} \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, \forall q \in \{2, \dots, p\}.$$

where we write  $f_1$  for the random observed gradient at the 1<sup>st</sup> iteration (they are i.i.d. so it does not matter).

**Assumption A.8 (A.4 tris).** Alternatively, suppose there exists a **global**  $\tau \geq 0$  for which  $\sup_{\theta \in \mathbb{R}^d} \left( \mathbb{E} [\|\epsilon(\theta)\|_2^4] \right)^{1/4} \leq \tau$ .

### A.I Remarks on assumptions

The assumption of random fields is weaker than assuming  $(\epsilon_k(\theta))_{k \in \mathbb{N}}$  is i.i.d. for all  $\theta$ , since it is not at fixed  $\theta$  but globally.

Assumption A.6 is only needed for theorem 2.14.

The reason to introduce Assumptions A.7-A.8 is to make clear how to avoid the bounded assumption of A.4. In particular, assumption A.7 is the weakest, and assumption A.8 is the strongest. Moreover, if we make the assumption of i.i.d. errors as above, even just for the sequence  $(\epsilon_k(\theta_{k-1}))_{k \in \mathbb{N}}$  is sufficient to imply assumption A.8 alone. In most works, the noise is taken to be “completely independent” as in this example, which is technically termed **semi-stochastic**. In such regard, this work is far more general than classical ML analysis.

Let us now make clear which is needed for which.

<sup>5</sup> This condition is global: we see  $\epsilon$  as a random mapping  $\mathbb{R}^d \ni \theta \mapsto \epsilon(\theta) \in \mathbb{R}^d$ , so we do not state  $\forall \theta \in \mathbb{R}^d$  merely because it is not fixed!

- For proposition 2.2 we need A.1-A.2-A.3-A.4 with  $p = 2$ ;
- for proposition 2.6 we need A.1-A.2-A.3-A.4, with  $p = 4$  in the special case and  $p = \max\{6, 2(k_\epsilon + 1)\}$  in general;
- for theorem 2.10 we need A.1-A.2-A.3-A.4 with  $p = 4$ ;
- for theorem 2.14 we need A.1-A.2-A.4-A.5-A.6 with:

– an additional condition on the noise, such that for some  $q \in \mathbb{N}, C \geq 0$  for all  $\theta \in \mathbb{R}^d$  one has:

$$\mathbb{E} \left[ \|\epsilon_1(\theta)\|_2^{p+k_\epsilon} \right] \leq C(1 + \|\theta - \theta^*\|_2^q); \quad (\text{A.9})$$

- A.6 holding with  $\ell = 5, p$ , i.e. the five times continuous differentiability and the  $p$  from the other assumptions match;
- A.4 holding for a slightly tweaked  $\tilde{p} = \max\{p + 3 + q, k_\epsilon\}$ ;
- the step size being changed to  $\gamma \in (0, 1/\zeta L)$  for  $\zeta \equiv \zeta(\tilde{p}) > 0$ .

It is clear that the last result is more technical than the others.

## B SOME PROOFS AND OTHER STATEMENTS

Let us recall the full statement to have it ready.

**Proposition 2.2** (Prop.2 in (Dieuleveut, Durmus, and Bach 2020)). *Let  $\gamma \in (0, 2/L)$ . The iterations of SGD seen as a Markov chain admit a unique stationary distribution  $\pi_\gamma$  that has finite second moments, i.e.  $\pi_\gamma \in \mathcal{P}_2(\mathbb{R}^d)$ . We can also **quantify** the rate of convergence in two ways: via the Wasserstein distance over probability measures (def. D.16) and in terms of “nice” functions integrated out by the kernel over time. Mathematically for all  $\theta \in \mathbb{R}^d$  and  $k \in \mathbb{N}$ :*

$$W_2^2(R_\gamma^k(\theta, \cdot), \pi_\gamma) \leq [\hbar(L, \mu, \gamma)]^k \int_{\mathbb{R}^d} \|\theta - \xi\|_2^2 d\pi_\gamma(\xi), \quad (2.2)$$

initial distribution (points to  $R_\gamma^k(\theta, \cdot)$ )  
rate function; goes exp fast to zero (points to  $[\hbar(L, \mu, \gamma)]^k$ )  
invariant measure (points to  $\pi_\gamma$ )  
mean square distance in invariant measure (initial conditions) (points to the integral)

and for all  $\theta \in \mathbb{R}^d, k \in \mathbb{N}$  and  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$  Lipschitz with constant  $L_\phi$ :

$$|R_\gamma^k \phi(\theta) - \pi_\gamma(\phi)| \leq L_\phi [\hbar(L, \mu, \gamma)]^{k/2} \sqrt{\int_{\mathbb{R}^d} \|\theta - \xi\|_2^2 d\pi_\gamma(\xi)}. \quad (2.3)$$

Naturally, the function  $\hbar \equiv \hbar(L, \mu, \gamma, k)$  is a function that is less than one exactly because we have the stability condition on the step-size.

*Proof.* The idea of the proof is very simple: we know that the Wasserstein distance is very nice: namely it makes  $\mathcal{P}_2(\mathbb{R}^d)$  complete, therefore reducing the question of convergence to proving that the sequence is Cauchy. In our case, we have a sequence of probability measures generated by multiple applications of the same kernel  $R_\gamma$ , therefore reducing the question to possibly bounding a single kernel.<sup>6</sup> Then, we have plenty of assumptions on the function and the type of randomness we allow, which will hopefully bring us from a rather abstract distance over probability measures to a friendlier Euclidean distance, that is also crucially at one step before. As we will see, finding a direct bound of the Wasserstein distance in terms of Euclidean distance is by a routine argument, but the assumptions on our problem are really what makes it possible to come back by one-step while keeping the inequality. Once we forget one step, by the fact that again we apply the same kernel, we will be able to iterate infinitely these steps. Hopefully, the multiplied discount at each iteration will be fast enough to establish convergence. Once this fact is established, existence and uniqueness of the invariant probability measure follow trivially by special choices of our more general master equation of

<sup>6</sup> Notice that here it is crucial to have the same step-size to make the reasoning this simple.

“fast convergence”.

We enforce the stability condition that  $\gamma \in (0, 2/L)$ , and consider two starting distributions  $\lambda_1, \lambda_2 \in \mathcal{P}_2(\mathbb{R}^d)$ . A standard result in optimal transport (see Villani (2009, thm. 4.1)) states that the Wasserstein distance can be represented as a distance of random variables  $\theta_0^{(1)}, \theta_0^{(2)}$  independent of the noise field  $(\epsilon_k)_{k \in \mathbb{N}}$ . In other words, we have  $W_2^2(\lambda_1, \lambda_2) = \mathbb{E} \left[ \left\| \theta_0^{(1)} - \theta_0^{(2)} \right\|_2^2 \right]$ . Let us run these two chains with the same step size  $\gamma$ , omitted in the notation for simplicity. Let us run the two chains with the same noise field, so that their gradients are respectively:

$$\nabla f_k(\theta_k^{(1)}) = \nabla f(\theta_k^{(1)}) + \epsilon_{k+1}(\theta_k^{(1)}); \quad (\text{B.1})$$

$$\nabla f_k(\theta_k^{(2)}) = \nabla f(\theta_k^{(2)}) + \epsilon_{k+1}(\theta_k^{(2)}). \quad (\text{B.2})$$

Let us use the independence of chain and noise to show that noise is de-correlated across chains. As a canonical example we have:

$$\mathbb{E} \left[ \left\langle \theta_0^{(1)}, \epsilon(\theta_0^{(2)}) \right\rangle \right] = \mathbb{E} \left[ \mathbb{E} \left[ \left\langle \theta_0^{(1)}, \epsilon(\theta_0^{(2)}) \right\rangle \mid \mathcal{F}_0 \right] \right] \stackrel{\text{A.3}}{=} \mathbb{E} \left[ \left\langle \theta_0^{(1)}, \mathbb{E} \left[ \epsilon(\theta_0^{(2)}) \mid \mathcal{F}_0 \right] \right\rangle \right] \stackrel{\text{A.3}}{=} 0. \quad (\text{B.3})$$

By symmetry, the result will hold if we swap the roles. A moment of thought also shows that the result holds if we consider the chain and the noise of the same chain, or any measurable function obviously. By construction, the combined chain distribution belongs to the set of couplings  $\Pi(\lambda R_\gamma^k, \lambda_2 R_\gamma^k)$  and since the Wasserstein distance is an infimum, we can write a one-step inequality:

$$W_2^2(\lambda_1 R_\gamma, \lambda_2 R_\gamma) \leq \mathbb{E} \left[ \left\| \theta_1^{(1)} - \theta_1^{(2)} \right\|_2^2 \right] \quad (\text{B.4})$$

$$= \mathbb{E} \left[ \left\| \theta_0^{(1)} - \gamma \nabla f_1(\theta_0^{(1)}) - \theta_0^{(2)} + \gamma \nabla f_1(\theta_0^{(1)}) \right\|_2^2 \right], \quad (\text{B.5})$$

$$= \mathbb{E} \left[ \left\| \theta_0^{(1)} - \gamma \nabla f(\theta_0^{(1)}) - \gamma \epsilon_1(\theta_0^{(1)}) - \theta_0^{(2)} + \gamma \nabla f(\theta_0^{(1)}) + \gamma \epsilon_1(\theta_0^{(1)}) \right\|_2^2 \right], \quad (\text{B.6})$$

where unrolled the step of the recursion. It is now natural to seek to exploit the orthogonality we just found. Let us then square the difference of starting points minus the difference of stochastic gradients. We find:

$$\mathbb{E} \left[ \left\| \theta_1^{(1)} - \theta_1^{(2)} \right\|_2^2 \right] = \mathbb{E} \left[ \left\| \theta_0^{(1)} - \theta_0^{(2)} \right\|_2^2 \right] + \gamma^2 \mathbb{E} \left[ \left\| \nabla f_1(\theta_0^{(1)}) - \nabla f_2(\theta_0^{(2)}) \right\|_2^2 \right] \quad (\text{B.7})$$

$$- 2\gamma \mathbb{E} \left[ \left\langle \underbrace{\nabla f_1(\theta_0^{(1)}) - \nabla f_1(\theta_0^{(2)})}_{= \nabla f(\theta_0^{(1)}) + \epsilon_1(\theta_0^{(1)}) - \nabla f(\theta_0^{(2)}) - \epsilon_1(\theta_0^{(2)})}, \theta_0^{(1)} - \theta_0^{(2)} \right\rangle \right]. \quad (\text{B.8})$$

The second term is not really friendly, but we can use our result quite directly by cancelling the noise dependent terms! We therefore find a more amenable expression:

$$\mathbb{E} \left[ \left\| \theta_1^{(1)} - \theta_1^{(2)} \right\|_2^2 \right] = \mathbb{E} \left[ \left\| \theta_0^{(1)} - \theta_0^{(2)} \right\|_2^2 \right] + \gamma^2 \mathbb{E} \left[ \left\| \nabla f_1(\theta_0^{(1)}) - \nabla f_2(\theta_0^{(2)}) \right\|_2^2 \right] \quad (\text{B.9})$$

$$- 2\gamma \mathbb{E} \left[ \left\langle \nabla f(\theta_0^{(1)}) - \nabla f(\theta_0^{(2)}), \theta_0^{(1)} - \theta_0^{(2)} \right\rangle \right]. \quad (\text{B.10})$$

Let us stress that inside the norm we will have the *noisy* gradient while in the inner product we have the *unbiased* gradient evaluated at the random iterates. At this moment we need to exploit the assumptions on the function class we chose to get rid of the random gradients. By  $L$ -co-coercivity of  $f_1$  (def. D.13) assumed in A.4 we find that:

$$\mathbb{E} \left[ \left\| \nabla f_1(\theta_0^{(1)}) - \nabla f_2(\theta_0^{(2)}) \right\|_2^2 \right] \leq L \mathbb{E} \left[ \left\langle \nabla f_1(\theta_0^{(1)}) - \nabla f_2(\theta_0^{(2)}), \theta_0^{(1)} - \theta_0^{(2)} \right\rangle \right], \quad (\text{B.11})$$

and another application of independence allows us to say that the noise cancels with the vectors in the expectation of the inner product. Thus,

$$\mathbb{E} \left[ \left\| \theta_1^{(1)} - \theta_1^{(2)} \right\|_2^2 \right] \leq \mathbb{E} \left[ \left\| \theta_0^{(1)} - \theta_0^{(2)} \right\|_2^2 \right] - 2\gamma \left( 1 - \gamma \frac{L}{2} \right) \mathbb{E} \left[ \left\langle \nabla f(\theta_0^{(1)}) - \nabla f(\theta_0^{(2)}), \theta_0^{(1)} - \theta_0^{(2)} \right\rangle \right]. \quad (\text{B.12})$$

Now that we have recovered an expression that depends only on  $f$ , we can use the strong convexity from assumption A.1 to conclude:

$$\mathbb{E} \left[ \left\| \boldsymbol{\theta}_1^{(1)} - \boldsymbol{\theta}_1^{(2)} \right\|_2^2 \right] \leq \left[ 1 - 2\mu\gamma \left( 1 - \gamma \frac{L}{2} \right) \right] \mathbb{E} \left[ \left\| \boldsymbol{\theta}_0^{(1)} - \boldsymbol{\theta}_0^{(2)} \right\|_2^2 \right]. \quad (\text{B.13})$$

What we find is that a one-step application of the kernel on both distributions has a bound with respect to the starting distance. Inducting over this, at the  $k^{\text{th}}$  step we will have that:

$$W_2^2(\lambda_1 R_\gamma^k, \lambda_2 R_\gamma^k) \leq \mathbb{E} \left[ \left\| \boldsymbol{\theta}_k^{(1)} - \boldsymbol{\theta}_k^{(2)} \right\|_2^2 \right] \leq \left[ \underbrace{1 - 2\mu\gamma \left( 1 - \gamma \frac{L}{2} \right)}_{:=\hbar} \right]^k W_2^2(\lambda_1, \lambda_2). \quad (\text{B.14})$$

In words, starting from any distribution, we will get exponentially close to the initial distance as steps go on, if and only if  $\gamma \in (0, 2/L)$ , which we indeed assumed. Choosing specifically  $\lambda_2 = \lambda_1 R_\gamma$  to be one step later,<sup>7</sup> we obtain that:

$$\sum_{k=1}^N W_2^2(\lambda_1 R_\gamma^k, \lambda_2 R_\gamma^k) = \sum_{k=1}^N W_2^2(\lambda_1 R_\gamma^k, \lambda_1 R_\gamma^{k+1}) \leq \sum_{k=1}^N \hbar^k W_2^2(\lambda) \quad (\text{B.15})$$

Since we found that the elements are decaying geometrically, the sum on the left-hand side is finite, therefore making the sequence  $(\lambda_1 R_\gamma^k)_{k \in \mathbb{N}}$  Cauchy. Obviously, our distance is nice and makes  $\mathcal{P}_2(\mathbb{R}^d)$  complete (see Villani (2009, thm. 6.16)), so Cauchyness is characterized by having a limit in  $\mathcal{P}_2(\mathbb{R}^d)$ , which we call  $\pi_\gamma^{\lambda_1} \in \mathcal{P}_2(\mathbb{R}^d)$ .

We have more, since from any two starting distributions their distance decays quickly, so we can automatically obtain uniqueness of the limit by contradiction. Indeed, assume there exists another limit  $\pi_\gamma^{\lambda_2}$  for  $\lambda_2 \neq \lambda_1$ , then:

$$W_2(\pi_\gamma^{\lambda_1}, \pi_\gamma^{\lambda_2}) \leq W_2(\pi_\gamma^{\lambda_1}, \lambda_1 R_\gamma^k) + W_2(\lambda_1 R_\gamma^k, \pi_\gamma^{\lambda_2}) \quad \text{tr. ineq.;} \quad (\text{B.16})$$

$$\leq W_2(\pi_\gamma^{\lambda_1}, \lambda_1 R_\gamma^k) + W_2(\lambda_1 R_\gamma^k, \lambda_2 R_\gamma^k) + W_2(\lambda_2 R_\gamma^k, \pi_\gamma^{\lambda_2}) \quad \text{tr. ineq.;} \quad (\text{B.17})$$

$$\xrightarrow{k \rightarrow \infty} 0, \quad (\text{B.18})$$

since all the three terms go to zero: the outer by assumption, the middle one by the general decaying property we found. By the fact that the Wasserstein is a distance, *a fortiori* we have  $\pi_\gamma^{\lambda_1} = \pi_\gamma^{\lambda_2}$  contradicting the hypothesis.

Let us now pass to the quantification of convergence. We aim to apply our master equation eq. (B.14) again, for a particular choice. Taking  $\lambda_1 = \delta_\theta$ ,  $\lambda_2 = \pi_\gamma$ , we will have that  $\lambda_2 R_\gamma^k = \pi_\gamma R_\gamma^k = \pi_\gamma \in \mathcal{P}_2(\mathbb{R}^d)$  by stationarity, and thus will look at a distance from an arbitrary starting point and the final distribution. Specializing the bound, the claim is proved by definition of Wasserstein distance (def. D.16). Also, the bound is non-trivial since:  $\int_{\mathbb{R}^d} \|\boldsymbol{\theta} - \boldsymbol{\xi}\|_2^2 d\pi_\gamma(\boldsymbol{\xi}) \leq 2\|\boldsymbol{\theta}\|^2 + 2 \int_{\mathbb{R}^d} \|\boldsymbol{\xi}\|_2^2 d\pi_\gamma(\boldsymbol{\xi}) < \infty$ .

For the Lipschitz criterion, a little more does the game. With the same particular measures, we evaluate expectations of Lipschitz functions  $\phi$  with  $L_\phi$  constant:

$$\left| R_\gamma^k(\phi(\boldsymbol{\theta})) - \pi_\gamma \phi \right| = \left| \mathbb{E} \left[ \phi(\boldsymbol{\theta}_{k,\gamma}^{(1)}) - \phi(\boldsymbol{\theta}_{k,\gamma}^{(2)}) \right] \right| \quad (\text{B.19})$$

$$\leq L_\phi \mathbb{E} \left[ \left\| \boldsymbol{\theta}_{k,\gamma}^{(1)} - \boldsymbol{\theta}_{k,\gamma}^{(2)} \right\|_2 \right] \quad (\text{B.20})$$

$$\leq L_\phi \sqrt{\mathbb{E} \left[ \left\| \boldsymbol{\theta}_{k,\gamma}^{(1)} - \boldsymbol{\theta}_{k,\gamma}^{(2)} \right\|_2^2 \right]} \quad (\text{B.21})$$

$$\leq L_\phi \sqrt{\hbar^k} \sqrt{\int_{\mathbb{R}^d} \|\boldsymbol{\theta} - \boldsymbol{\xi}\| d\pi_\gamma(\boldsymbol{\xi})} \quad \text{master equation B.14.} \quad (\text{B.22})$$

□

<sup>7</sup> We can do this since by A.2-A.3-A.4 it holds that  $\lambda_1 R_\gamma \in \mathcal{P}_2(\mathbb{R}^d)$ , by an easy triangle inequality and the uniform bound on the gradients basically.

**Proposition B.23** (Prop. 17 in (Dieuleveut, Durmus, and Bach 2020)). *Recall the statement for the baby example of proposition 2.6. There, we can further say that for all  $\gamma \in (0, 1/r^2)$ , where  $r^2$  is defined in appendix C that:*

$$\mathbf{R}^{-1} \int_{\mathbb{R}^d} [\boldsymbol{\theta} - \boldsymbol{\theta}^*][\boldsymbol{\theta} - \boldsymbol{\theta}^*]^\top d\pi_\gamma(\boldsymbol{\theta}) = \gamma \mathbb{E} [\boldsymbol{\xi}_1 \boldsymbol{\xi}_1^\top], \quad (\text{B.24})$$

where  $\boldsymbol{\xi}$  is the multiplicative error of equation eq. (C.6).

## C BABY EXAMPLE DETAILS

In particular, if we take the square loss and a linear model  $y = \langle \mathbf{x}, \boldsymbol{\theta}^* \rangle + \eta$ , we find  $\mathcal{L}(\mathbf{x}, y, \boldsymbol{\theta}) = (\langle \mathbf{x}, \boldsymbol{\theta} \rangle - y)^2$  which has an explicit generalization error:<sup>8</sup>

$$\mathcal{E}(\boldsymbol{\theta}; \mathbf{L}) = \mathbb{E}_{(\mathbf{x}, y)} [\langle \mathbf{x}, \boldsymbol{\theta} \rangle - y]^2 = \mathbb{E}_{(\mathbf{x}, y)} [\langle \mathbf{x}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle - \eta]^2 = \|\boldsymbol{\Sigma}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|_2^2 + \sigma^2, \quad \boldsymbol{\Sigma} := \mathbb{E} [\mathbf{x} \mathbf{x}^\top]. \quad (\text{C.1})$$

Instead, the iterations of SGD will see a gradient that is unbiased:

$$\mathbb{E}_{(x_k, y_k)} [\nabla \mathcal{L}(\mathbf{x}_k, y_k, \boldsymbol{\theta}_k^{(\gamma)})] = \mathbb{E}_{(x_k, \eta_k)} [2x_k (\langle x_k, \boldsymbol{\theta}_k^{(\gamma)} - \boldsymbol{\theta}^* \rangle + \eta_k)] = 2\boldsymbol{\Sigma}(\boldsymbol{\theta}_k^{(\gamma)} - \boldsymbol{\theta}^*) = \nabla \mathcal{E}(\boldsymbol{\theta}_k^{(\gamma)}; \mathcal{L}). \quad (\text{C.2})$$

Moreover, it is interesting to express the gradient in terms of the true gradient by making the right terms appear:

$$\frac{1}{2} \nabla \mathcal{L}(\mathbf{x}_k, y_k, \boldsymbol{\theta}_k^{(\gamma)}) = \mathbf{x}_k \mathbf{x}_k^\top (\boldsymbol{\theta}_k^{(\gamma)} - \boldsymbol{\theta}^*) - \mathbf{x}_k \eta_k \quad (\text{C.3})$$

$$= \boldsymbol{\Sigma}(\boldsymbol{\theta}_k^{(\gamma)} - \boldsymbol{\theta}^*) + [\mathbf{x}_k \mathbf{x}_k^\top - \boldsymbol{\Sigma}](\boldsymbol{\theta}_k^{(\gamma)} - \boldsymbol{\theta}^*) - \mathbf{x}_k \eta_k \quad (\text{C.4})$$

$$= \frac{1}{2} \nabla \mathcal{E}(\boldsymbol{\theta}_k^{(\gamma)}; \mathbf{L}) + [\mathbf{x}_k \mathbf{x}_k^\top - \boldsymbol{\Sigma}](\boldsymbol{\theta}_k^{(\gamma)} - \boldsymbol{\theta}^*) - \mathbf{x}_k \eta_k. \quad (\text{C.5})$$

We can further make sense of the error term by a decomposition into an *additive* and a *multiplicative* term as  $\epsilon_k \equiv \epsilon_k(\boldsymbol{\theta}_k^{(\gamma)}) = \rho_k(\boldsymbol{\theta}_k^{(\gamma)}) + \xi_k$ , where

$$\rho_k \equiv \rho_k(\boldsymbol{\theta}) := [\mathbf{x}_k \mathbf{x}_k^\top - \boldsymbol{\Sigma}](\boldsymbol{\theta}_k^{(\gamma)} - \boldsymbol{\theta}^*), \quad \xi_k := -\mathbf{x}_k \eta_k = (\mathbf{x}_k^\top \boldsymbol{\theta}^* - y_k) \mathbf{x}_k. \quad (\text{C.6})$$

In particular, the multiplicative noise is independent of the current iterate. Both appear in the proof of proposition 2.10.

**ABOUT THE ASSUMPTIONS** It is not automatic that the baby example of section 2.II satisfies the assumptions. In particular:

- for  $p \geq 2$ , assumption A.4 will hold;
- assumption A.5 requires almost sure boundedness of observations;
- for A.4 to hold, we might require that there exists  $r \geq 0$  such that  $\mathbb{E} [\|\mathbf{x}_k\|_2^2 \mathbf{x}_k \mathbf{x}_k^\top] \leq r^2 \boldsymbol{\Sigma}$ , which holds if data is a.s. bounded or has bounded kurtosis, as mentioned in (Dieuleveut, Flammarion, and Bach 2017).

## D USEFUL DEFINITIONS AND THEOREMS

We report here some classical definitions for completeness.

**Definition D.1** ( $\mu$ -strong convexity). *A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex for  $\mu > 0$  when:*

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{x}') \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{x}') - (1 - \alpha) \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}, \forall \alpha \in [0, 1]. \quad (\text{D.2})$$

**Definition D.3** ( $L$ -smoothness). *A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$ -smooth if its gradient is  $L$ -Lipschitz.*

<sup>8</sup> Let  $\eta$  be square integrable and have variance  $\sigma^2$ .

**Lemma D.4** (Descent lemma). *If  $f$  is  $L$ -smooth then for all  $x \in \text{int}(\text{dom}(f))$  and  $y \in \text{dom}(f)$  it holds:*

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|x - y\|_2^2. \quad (\text{D.5})$$

*Proof.* Let  $t \in [0, 1]$ . Define the interpolating scalar function:

$$\varphi(t) := f(x + th) - f(x) - \langle \nabla f(x), th \rangle, \quad h := y - x. \quad (\text{D.6})$$

In particular,  $\varphi(0) = 0$  and we wish to prove  $\varphi(1) \leq \frac{L}{2} \|x - y\|_2^2$ . Clearly,  $\varphi$  is differentiable since  $f$  is. We can bound the derivative as:

$$\varphi'(t) = \langle h, \nabla f(x + th) \rangle - \langle h, \nabla f(x) \rangle = \langle \nabla f(x + th) - \nabla f(x), h \rangle \leq tL \|h\|_2^2. \quad (\text{D.7})$$

We can then apply the fundamental theorem of calculus:

$$f(y) - \varphi(1) - \varphi(0) = \int_0^1 \varphi'(t) dt \leq L \|h\|_2^2 \int_0^1 t dt = \frac{L}{2} \|h\|_2^2, \quad (\text{D.8})$$

which proves the claim.  $\square$

Combining the first definition and the descent lemma, we find that assuming  $\mu$ -strong convexity and  $L$ -smoothness bounds the function in a variable way as:

$$f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|x - y\|_2^2 \leq f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|x - y\|_2^2, \quad \forall x, y. \quad (\text{D.9})$$

In words, at each fixed point  $y$ , we need to bound the function value by any “parabola” parameterized by the auxiliary  $x$  vectors. While seemingly irrelevant, this is quite stringent. If we take the special case where  $\frac{\mu}{L} \equiv 1$ , then we only have quadratic functions! If we allow for  $\mu < L$  strictly, we get a non-trivial class, that is however still very stringent. On a side note, if we reorder terms we can see that our definition is essentially a bound on the convexity gap, which must be quadratic:

$$\frac{\mu}{2} \|x - y\|_2^2 \leq \underbrace{f(y) - f(x) - \langle \nabla f(x), y - x \rangle}_{:= \text{D}_{\text{Bre}}(y|x;f)} \leq \frac{L}{2} \|x - y\|_2^2, \quad (\text{D.10})$$

where the middle term is also termed Bregman divergence (hence the notation in underbrace).

**Definition D.11** (Locally-polynomially Lipschitz). *A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$  is locally-polynomially Lipschitz (locally Lipschitz for short) if there exists  $\alpha \geq 0$  such that:*

$$\|f(x) - f(x')\|_2 \leq \left(1 + \|x\|_2^\alpha + \|x'\|_2^\alpha\right) \|x - x'\|_2, \quad \forall x, x' \in \mathbb{R}^d. \quad (\text{D.12})$$

**Definition D.13** (C-co-coercivity). *A function  $g : \mathbb{R}^d \rightarrow \mathbb{R}^p$  is C-co-coercive if for all  $x, x' \in \mathbb{R}^d$  we have:*

$$C \langle g(x) - g(x'), x - x' \rangle \geq \|g(x) - g(x')\|_2^2. \quad (\text{D.14})$$

In the main text we assume that the gradient is both  $L$ -co-coercive and Lipschitz, by a simple application of the Cauchy-Schwartz inequality we will have that:

$$\|\nabla f(x) - \nabla f(x')\|_2^2 \leq L \langle \nabla f(x) - \nabla f(x'), x - x' \rangle \leq L^2 \|x - x'\|_2^2, \quad \forall x, x' \in \mathbb{R}^d. \quad (\text{D.15})$$

**Definition D.16** (Order 2-Wasserstein distance). *Let  $\lambda, \nu \in \mathcal{P}_2(\mathbb{R}^d)$  be two measures with finite second moment, and  $\Pi(\lambda, \nu)$  be the set of couplings on  $\mathbb{R}^d \times \mathbb{R}^d$  that have marginals  $(\lambda, \nu)$ . We define their 2-Wasserstein distance as:*

$$W_2(\lambda, \nu) := \inf_{\xi \in \Pi(\lambda, \nu)} \left( \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|_2^2 \xi(dx, dy) \right)^{\frac{1}{2}}. \quad (\text{D.17})$$

**Definition D.18** (Markov kernel). *A Markov kernel  $R$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  is a map such that:*

(K1) *for all  $A \in \mathcal{B}(\mathbb{R}^d)$  the mapping  $\theta \mapsto R(\theta, A)$  is Borel measurable;*

(K2) for all  $\theta \in \mathbb{R}^d$  the mapping  $A \mapsto R(\theta, A)$  is a probability measure on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ .

In practice, we will use it to abstractly represent how the law of eq. (1.1) changes across iterations. That is, we let  $R_\gamma$  depend on  $\gamma$  the step-size, and starting from an arbitrary  $\theta_0 \in \mathbb{R}^d$  (that can be sampled) consider the recursion:

$$R_\gamma^1 := R_\gamma, \quad R_\gamma^{k+1}(\theta_0, A) = \int_{\mathbb{R}^d} R_\gamma^k(\theta_0, d\theta) R_\gamma(\theta, A), \quad \forall \theta \in \mathbb{R}^d, A \in \mathcal{B}(\mathbb{R}^d), k \in \mathbb{N}, \quad (\text{D.19})$$

where the base iterate acts such that almost surely  $R_\gamma(\theta_k, A) = \mathbb{P}[\theta_{k+1} \in A \mid \theta_k]$  for all  $k \in \mathbb{N}, \forall A \in \mathcal{B}(\mathbb{R}^d)$ . In (informal!) words, the probability that we are in set  $A$  if we were at  $\theta_k$  before. Such notational construction allows us to define measures and measurable functions as follows. Since by (K1) at fixed  $A$  we have a measurable function, we can integrate out for a given  $\lambda \in \mathcal{P}(\mathbb{R}^d)$  as:

$$\lambda R_\gamma^k(\cdot) := \int_{\mathbb{R}^d} \lambda(d\theta) R_\gamma^k(\theta, \cdot) : \mathcal{B}(\mathbb{R}^d) \rightarrow \mathcal{P}(\mathbb{R}^d) \quad (\text{D.20})$$

$$A \mapsto \lambda R_\gamma^k(A) = \mathbb{P}[\theta_k^{(\gamma)} \in A \mid \theta_0 \sim \lambda]. \quad (\text{D.21})$$

Namely, if we chain from the left we obtain the probability of starting from  $\lambda$  and ending in  $A$  at the  $k^{\text{th}}$  step. If instead we chain from the right, we will obtain a measurable function representing the expectation at the  $k^{\text{th}}$  step with respect to the iterated distribution. Mathematically, for  $\phi \in F_+(\mathbb{R}^d, \mathbb{R})$  the space of positive measurable functions we have that:

$$R_\gamma^k \phi(\cdot) := \int_{\mathbb{R}^d} \phi(\theta) R_\gamma^k(\cdot, d\theta) : \mathbb{R}^d \rightarrow F_+(\mathbb{R}^d, \mathbb{R}) \quad (\text{D.22})$$

$$\theta_0 \mapsto R_\gamma^k \phi(\theta_0) = \mathbb{E} \left[ \phi \left( \theta_k^{(\gamma)} \right) \mid \theta_0 \right]. \quad (\text{D.23})$$

We therefore obtain probabilities integrating from the left and expectations integrating from the right. Note that by the notations it follows also that for all  $A \in \mathcal{B}(\mathbb{R}^d)$ :

$$\lambda \left( R_\gamma^k h \right) (\cdot) = \int_{\mathbb{R}^d} \lambda(d\theta) (R_\gamma^k h)(\theta, \cdot) \quad (\text{D.24})$$

$$= \int_{\mathbb{R}^d} \lambda(d\theta) \int_{\mathbb{R}^d} h(d\xi) R_\gamma^k(\theta, d\xi) \quad (\text{D.25})$$

$$= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \lambda(d\theta) R_\gamma^k(\theta, d\xi) h(\xi) \quad (\text{D.26})$$

$$= \int_{\mathbb{R}^d} h(\xi) \int_{\mathbb{R}^d} \lambda(d\theta) R_\gamma^k(\theta, d\xi) \quad (\text{D.27})$$

$$= \int_{\mathbb{R}^d} h(\theta) (\lambda R_\gamma^k)(\cdot, d\theta) \quad (\text{D.28})$$

$$= (\lambda R_\gamma^k)(h)(\cdot). \quad (\text{D.29})$$

**Proposition D.30** (Theorem 6.14 in (Le Gall 2016)). *Let  $h, g$  be continuous functions on a euclidean domain  $E \subset \mathbb{R}^d$  that tend to zero at infinity.<sup>9</sup> The following are equivalent:*

1.  $h \in D(\mathcal{A})$  and  $\mathcal{A}h = g$ ;
2. for all  $x \in E$  the process:

$$h(\theta_t^{\theta_0}) - \int_0^t g(\theta_s^{\theta_0}) ds, \quad (\text{D.31})$$

is a martingale with respect to the (canonical) filtration  $(\mathcal{F}_t)_{t \in [0, \infty]}$ . Here, by  $\theta_t^{\theta_0}$  we mean the Markov chain  $(\theta_t)_{t \in \mathbb{R}_+}$  that started almost surely at  $\theta_0$ .

<sup>9</sup> This is a specific definition: for all  $\epsilon > 0$  there exists a compact set  $K$  such that  $|f(\theta)| < \epsilon$  for all  $\theta \in E \setminus K$ .