



HAL
open science

An SDE perspective on stochastic convex optimization

Rodrigo Maulen-Soto, Jalal Fadili, Hedy Attouch

► **To cite this version:**

Rodrigo Maulen-Soto, Jalal Fadili, Hedy Attouch. An SDE perspective on stochastic convex optimization. 2022. hal-03716996

HAL Id: hal-03716996

<https://hal.science/hal-03716996v1>

Preprint submitted on 8 Jul 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

An SDE perspective on stochastic convex optimization

Rodrigo Maulen S.*

Jalal Fadili[†]

Hedy Attouch[‡]

Abstract. In this paper, we analyze the global and local behavior of gradient-like flows under stochastic errors towards the aim of solving convex optimization problems with noisy gradient input. We first study the unconstrained differentiable convex case, using a stochastic differential equation where the drift term is minus the gradient of the objective function and the diffusion term is either bounded or square-integrable. In this context, under Lipschitz continuity of the gradient, our first main result shows almost sure convergence of the objective and the trajectory process towards a minimizer of the objective function. We also provide a comprehensive complexity analysis by establishing several new pointwise and ergodic convergence rates in expectation for the convex, strongly convex and (local) Łojasiewicz case. The latter, which involves local analysis, is challenging and requires non-trivial arguments from measure theory. Then, we extend our study to the constrained case and more generally to certain nonsmooth situations. We show that several of our results have natural extensions obtained by replacing the gradient of the objective function by a cocoercive monotone operator. This makes it possible to obtain similar convergence results for optimization problems with an additively "smooth + non-smooth" convex structure. Finally, we consider another extension of our results to non-smooth optimization which is based on the Moreau envelope.

Key words. Convex optimization, Stochastic Differential Equation, Stochastic gradient descent, Łojasiewicz inequality, KL inequality, Convergence rate, Asymptotic behavior.

AMS subject classifications. 37N40, 46N10, 49M30, 65B99, 65K05, 65K10, 90B50, 90C25

1 Introduction

1.1 Problem Statement

We aim to solve convex minimization problems by means of stochastic differential equations whose drift term is driven by the gradient of the objective function. This allows for noisy (inaccurate) gradient input to be taken into account. Consider the minimization problem

$$\min_{x \in \mathbb{R}^d} f(x), \tag{P}$$

where the objective f satisfies the following standing assumptions:

$$\left\{ \begin{array}{l} f \text{ is continuously differentiable and convex with } L\text{-Lipschitz continuous gradient;} \\ \mathcal{S} \stackrel{\text{def}}{=} \operatorname{argmin}(f) \neq \emptyset. \end{array} \right. \tag{H_0}$$

We will also later deal with the constrained case, and more generally with additively structured "smooth + nonsmooth" convex optimization.

*Normandie Université, ENSICAEN, UNICAEN, CNRS, GREYC, France. E-mail: rodrigo.maulen@ensicaen.fr

[†]Normandie Université, ENSICAEN, UNICAEN, CNRS, GREYC, France. E-mail: Jalal.Fadili@ensicaen.fr

[‡]IMAG, CNRS, Université Montpellier, France. E-mail: hedy.attouch@umontpellier.fr

Let us first recall some basic facts about the deterministic case. To solve **(P)**, a fundamental dynamic to consider is the gradient flow of f , *i.e.* the gradient descent dynamic with initial condition $X_0 \in \mathbb{R}^d$:

$$\begin{cases} \dot{x} = -\nabla f(x), & t > 0 \\ x(0) = X_0. \end{cases} \quad (\text{GF})$$

It is well known since the founding papers of Brezis, Baillon, Bruck in the 1970s that, if the solution set $\text{argmin } f$ of **(P)** is non-empty, then each solution trajectory of **(GF)** converges, and its limit belongs to $\text{argmin } f$. In fact, this result is true in a more general setting, simply assuming that the objective function f is convex, lower semicontinuous (lsc) and proper (in which case we must consider the differential inclusion obtained by replacing in **(GF)** the gradient of f by the sub-differential ∂f).

In many cases, the gradient input is subject to noise, for example, if the gradient cannot be evaluated directly, or due to some other exogenous factor. In such scenario, one can model these errors using a stochastic integral with respect to the measure defined by a continuous Itô martingale. This entails the following stochastic differential equation as a stochastic counterpart of **(GF)**:

$$\begin{cases} dX(t) = -\nabla f(X(t))dt + \sigma(t, X(t))dW(t), & t > 0 \\ X(0) = X_0, \end{cases} \quad (\text{SDE})$$

defined over a filtered probability space $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$, where the diffusion (volatility) term $\sigma : \mathbb{R}_+ \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$ is matrix-valued measurable function, and W is the m -dimensional Brownian motion.

Our goal is to study the dynamic of **(SDE)** and its long time behavior in order to solve **(P)**. To identify the assumptions necessary to hope for such a behavior to occur, remember that when the diffusion term σ is a positive real constant, it is well-known that $X(t)$ in this case is a continuous-time diffusion process known as Langevin diffusion, and has a unique invariant probability measure π_σ with density $\propto e^{-2f(x)/\sigma^2}$ [1]. It is easy to see that the measure π_σ gets concentrated around $\text{argmin } f$ as σ tends to 0^+ with $\lim_{\sigma \rightarrow 0^+} \pi_\sigma(\text{argmin } f) = 1$; see *e.g.* [2].

Motivated by this observation, our paper will then mostly focus on the case where $\sigma(\cdot, x)$ vanishes sufficiently fast as $t \rightarrow +\infty$ uniformly in x , and some guarantees will also be provided for uniformly bounded σ . Therefore, throughout, the entries σ_{ik} are assumed to satisfy:

$$\begin{cases} \sup_{t \geq 0, x \in \mathbb{R}^d} |\sigma_{ik}(t, x)| < \infty, \\ |\sigma_{ik}(t, x') - \sigma_{ik}(t, x)| \leq l_0 \|x' - x\|, \end{cases} \quad (\text{H})$$

for some $l_0 > 0$ and for all $t \geq 0, x, x' \in \mathbb{R}^d$. The Lipschitz continuity assumption is mild and required to ensure well-posedness of **(SDE)**.

1.2 Contributions

We study the properties of the process $X(t)$ and $f(X(t))$ for the stochastic differential equation **(SDE)** from an optimization perspective, under the assumptions **(H₀)** and **(H)**. When the diffusion term is uniformly bounded, we show convergence of $\mathbb{E}[f(X(t)) - \min f]$ to a noise-dominated region both for the convex and strongly convex case. When the diffusion term is square-integrable, we show in Theorem 3.1 that $X(t)$ converges almost surely to a solution of **(P)**, which is a new result to the best of our knowledge. Moreover, in Theorem 3.2 and Proposition 3.3, we provide new ergodic and pointwise convergence rates of the objective in expectation, again for both the convex and strongly convex case.

Then we turn to a local analysis relying on the Łojasiewicz inequality and its strong ties with error bounds. Since this property is most often satisfied only locally, we deepen the discussion on the long time localization of the process. This is fundamental, because in the recent literature on local convergence properties of stochastic gradient descent, strong assumptions are imposed, such as $X(t)$ or $f(X(t))$ is locally bounded almost surely. Such assumptions are unfortunately unrealistic due to the presence of the Brownian Motion. We manage to circumvent this problem by using arguments from measure theory, in particular Egorov's theorem. In turn, under the Łojasiewicz inequality assumption with exponent $q \geq 1/2$, this allows us to show local convergence rates of the objective and the trajectory itself in expectation over a set of events whose probability is arbitrarily close to 1 (see Theorem 4.5).

Table 1 summarizes the local and global convergence rates obtained for $\mathbb{E}[f(X(t)) - \min f]$. In this table, $\delta > 0$ is a parameter which is intended to be taken arbitrarily close to 0 but different from it, $\sigma_* > 0$ and $\sigma_\infty(\cdot)$ are defined as

$$\|\sigma(t, x)\|_F^2 \leq \sigma_*^2, \quad \forall t \geq 0, \forall x \in \mathbb{R}^d, \quad \text{and} \quad \sigma_\infty(t) \stackrel{\text{def}}{=} \sup_{x \in \mathbb{R}^d} \|\sigma(t, x)\|_F, \quad (1.1)$$

and $\sigma_\infty(\cdot)$ is a decreasing function. $\mathcal{L}^q(\mathcal{S})$ is the class of functions satisfying the Łojasiewicz inequality with exponent $q \in [0, 1]$ at each point of \mathcal{S} (see Definition 4.1)¹.

Property of f	Gradient Flow	SDE ($\sup_{t>0} \sigma_\infty(t) \leq \sigma_*$)	SDE ($\sigma_\infty \in \mathcal{L}^2(\mathbb{R}_+)$)
Convex	t^{-1}	$t^{-1} + \sigma_*^2$	t^{-1}
μ -Strongly Convex	$e^{-2\mu t}$	$e^{-2\mu t} + \sigma_*^2$	$\max\{e^{-2\mu t}, \sigma_\infty^2(t)\}$
Convex $\cap \mathcal{L}^{1/2}(\mathcal{S})$ (coef. μ)	$e^{-\mu^2 t}$	✗	$\max\{e^{-\mu^2 t}, \sigma_\infty^2(t)\} + \sqrt{\delta}$
Convex $\cap \mathcal{L}^q(\mathcal{S})$, $q \in (\frac{1}{2}, 1)$	$t^{-\frac{1}{2q-1}}$	✗	$t^{-\frac{1}{2q-1}} + \sqrt{\delta}$ ²

Table 1: Summary of local and global convergence rates obtained for $\mathbb{E}[f(X(t)) - \min f]$.

Although it is natural to think that we can take the limit when δ goes to 0^+ , the time from which these convergence rates are valid depends on δ and increases (potentially to $+\infty$) as δ approaches 0^+ . Assuming only the boundedness of the diffusion and the Łojasiewicz inequality, we could not find better results (cells marked with **✗**) than those presented in the convex case. Since the Łojasiewicz inequality is local, a natural approach would be to localize the process in the long term with high probability. However, it is not clear how to achieve this.

In Section 5, we turn to extending some of the preceding results to the structured convex minimization problem

$$\min_{x \in \mathbb{R}^d} f(x) + g(x), \quad (\text{P}_c) \tag{P_c}$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies (H_0) , $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is proper, lsc and convex and $\text{argmin}(f + g) \neq \emptyset$. This obviously covers the case of constrained minimization of f over a non-empty closed convex set. We take two different routes leading to different SDEs.

¹Semialgebraic and more generally analytic functions is a typical class verifying the Łojasiewicz inequality at each point [3, 4].

²This is not yet proven, our conjecture is that it is true when $\sigma_\infty = \mathcal{O}((t+1)^{-\frac{q}{2q-1}})$ (see the detailed discussion in Conjecture 4.11).

The first approach consists in reformulating (P_c) as finding for zeros of the operator $M_\mu : \mathbb{R}^d \rightarrow \mathbb{R}^d$

$$M_\mu(x) = \frac{1}{\mu} (x - \text{prox}_{\mu g}(x - \mu \nabla f(x))),$$

where $\mu > 0$ and $\text{prox}_{\mu g}$ is the proximal mapping of μg . It is well-known that the operator M_μ is cocoercive [5], hence monotone and Lipschitz continuous, and $M_\mu = \nabla f$ when g vanishes. The idea is then to replace the operator ∇f in (SDE) by M_μ leading to an SDE which will have many of the convergence properties obtained in the smooth convex case. This approach is in accordance with the deterministic theory for monotone cocoercive operators (see [6, 7, 5]).

The second approach regularizes the nonsmooth component g of the objective function using its Moreau envelope

$$g_\theta(x) = \min_{z \in \mathbb{R}^d} g(z) + \frac{1}{2\theta} \|x - z\|^2.$$

This leads to studying the dynamic (SDE) with the function $f + g_\theta$, which has a continuous Lipschitz gradient. This approximation method leads to a non-autonomous SDE. Note, however, that the noise in this case can be considered on the evaluation of $\nabla f(x)$, while it is on $M_\mu(x)$ in the first approach.

1.3 Relation to prior work

The gradient system (GF), which is valid on a general real Hilbert space \mathcal{H} , is a dissipative dynamical system, whose study dates back to Cauchy [8]. It plays a fundamental role in optimization: it transforms the problem of minimizing f into the study of the asymptotic behavior of the trajectories of (GF). This example was the precursor to the rich connection between continuous dissipative dynamical systems and optimization. Its Euler forward discretization (with stepsize $\gamma_k > 0$) is the celebrated gradient descent scheme

$$x_{k+1} = x_k - \gamma_k \nabla f(x_k). \tag{GD}$$

Under (H_0) , and for $(\gamma_k)_{k \in \mathbb{N}} \subset]0, 2/L[$, then we have both the convergence of the values $f(x_k) - \min f = \mathcal{O}(1/k)$ (in fact even $o(1/k)$), and the weak convergence of iterates $(x_k)_{k \in \mathbb{N}}$ to a point in $\text{argmin} f$. Moreover, if the Łojasiewicz inequality (4.1) (see [9]) is satisfied, then we can ensure the strong convergence of $(x_k)_{k \in \mathbb{N}}$ to a point in $\text{argmin} f$ and faster convergence rates than those ensured by the simple convexity hypothesis (see [10, 11]).

Now, let us focus on the finite-dimensional case ($\mathcal{H} = \mathbb{R}^d$). Although the Gradient Descent is a classical algorithm to solve the convex minimization problem, with the need to handle large-scale problems (such as in various areas of data science and machine learning), there has become necessary to find ways to get around the high computational cost per iteration that these problems entail. The Robbins-Monro stochastic approximation algorithm [12] is at the heart of Stochastic Gradient Descent (SGD), which, roughly speaking, consists in cheaply and randomly approximating the gradient at the price of obtaining a random noise in the solutions. Given an initial point $x_0 \in \mathbb{R}^d$, (SGD) updates the iterates according to

$$x_{k+1} = x_k - \gamma \nabla f(x_k) - \gamma \xi_k, \tag{SGD}$$

where ξ_k denotes the (random) noise term at the k -th iteration.

Recent work (see [13, 14, 15, 16, 17, 18]) has linked algorithm (SGD) with dynamic (SDE), showing the context under which (SDE) can be seen as an approximation (under a specific error) of (SGD) and vice-versa. For example, (SDE) can be interpreted as the pathwise solution to the Fokker-Planck equation (see [19]).

The Euler forward discretization (with stepsize $\gamma > 0$) of (SDE) when $d = m$ and $\sigma = \sqrt{2}I_d$ is the following algorithm:

$$X_{k+1} = X_k - \gamma \nabla f(X_k) + \sqrt{2\gamma} \xi_k, \quad (\text{LMC})$$

where $\xi_k \sim \mathcal{N}(0, I_d)$ (multivariate standard normal distribution). This algorithm, which is known as Langevin Monte Carlo (see [20]), is a standard sampling scheme, whose purpose is to generate samples from an approximation of a target distribution, in our case, proportional to $e^{-f(x)}$. Under appropriate assumptions on f , when γ is small and k is large such that $k\gamma$ is large, the distribution of X_k converges in different topologies or is close in various metrics to the target distribution with density $\propto e^{-f(x)}$. Asymptotic and non-asymptotic (with convergence rates) results of this kind have been studied in a number of papers under various conditions; see [21, 22, 23, 24, 25, 26] and references therein. By rescaling the problem, relation between sampling (*i.e.* (LMC)) and optimization (*i.e.* (SGD)) has been also investigated for the strongly convex case in *e.g.* [21].

Concerning (SDE), one can easily infer from [27, Proposition 7.4] that assuming $\sup_{x \in \mathbb{R}^d} \|\sigma(t, x)\|_F = o(1/\sqrt{\log(t)})$, and conditioning on the event that $X(t)$ is bounded, we have almost surely that the set of limits of convergent sequences $X(t_k)$, $t_k \rightarrow +\infty$ is contained in $\operatorname{argmin} f$. Using results on asymptotic pseudo-trajectories from [27], the work of [28, 29, 30] analyzed the behavior of the Stochastic Mirror Descent dynamics:

$$\begin{aligned} dY(t) &= -\nabla f(X(t))dt + \sigma(t, X)dW(t), \\ X(t) &= Q(\eta Y(t)), \end{aligned} \quad (\text{SMD})$$

where $\mathcal{X} \subset \mathbb{R}^d$ is a closed convex feasible region, f is convex with Lipschitz continuous gradient on \mathcal{X} , $Q : \mathbb{R}^d \rightarrow \mathcal{X}$ is the mirror map induced by some strongly convex entropy, and $\eta > 0$ is a sensitivity parameter. In [28, Theorem 4.1], it is shown that if \mathcal{X} is also assumed bounded, that $\sup_{x \in \mathbb{R}^d} \|\sigma(t, x)\|_F = o(1/\sqrt{\log(t)})$, and Q satisfies some continuity assumptions³, then the solution $X(t)$ (SMD) converges to a point in $\operatorname{argmin} f$ almost surely. Similar assumptions can be found in [30] to obtain almost sure convergence on the objective. Let us observe that all these results do not apply to our setting. Indeed, if $\mathcal{X} = \mathbb{R}^d$ (unconstrained problem), $Q(x) = x$ and $\eta = 1$, we recover (SDE). Our work does not assume any boundedness whatsoever to establish our results. This comes however at somewhat stronger assumptions on $\sigma(\cdot, \cdot)$.

While finalizing this work, we became aware of the recent work of [31], which analyzes the behavior of (SDE) for $f \in C^2(\mathbb{R}^d)$ not necessarily convex and which satisfies $\sup_{x \in \mathbb{R}^d} \|\sigma(\cdot, x)\|_F \in L^2(\mathbb{R}_+)$. Conditioning on the event that $\limsup_{t \rightarrow \infty} \|X(t)\| < \infty$, they showed that $\nabla f(X(t)) \rightarrow 0$ almost surely, almost sure convergence of $f(X(t))$, and if the objective f is semialgebraic (and more generally tame), they also showed almost sure convergence of $X(t)$ towards a critical point of f . They also made attempt to get local convergence rates under the Łojasiewicz inequality that are less transparent than ours. Our analysis on the other hand leverages convexity of f to establish stronger results.

1.4 Organization of the paper

Section 2 introduces notations and reviews some necessary material from convex and stochastic analysis. Section 3 states our main convergence results in the case of a convex differentiable objective function whose gradient is Lipschitz continuous. We first show the almost sure convergence of the process towards the set of minimizers, then we establish convergence rates for the values. Section 4 introduces further geometric properties of the objective functions, namely Łojasiewicz property and related error bound, which allows to

³Compactness of \mathcal{X} and the condition on $\sigma(\cdot, \cdot)$ are clearly reminiscent of [27, Proposition 7.4], though the latter is not discussed in [28].

obtain improved (local) convergence rates. This covers in particular the (locally) strongly convex case. In section 5, we extend some results to the nonsmooth case by considering the additively structured "smooth + nonsmooth" convex minimization. We develop new stochastic differential equations that naturally lend themselves to splitting techniques. Technical lemmas and theorems that are needed throughout the paper are collected in the appendix.

2 Notation and Preliminaries

We will use the following shorthand notations: given $d, n \in \mathbb{N}$, $[n] \stackrel{\text{def}}{=} \{1, \dots, n\}$, $\mathbb{R}^{d \times n}$ is the set of real matrices of size $d \times n$, and I_d is the identity matrix of dimension d . For $M \in \mathbb{R}^{d \times n}$, $M^\top \in \mathbb{R}^{n \times d}$ is its transpose matrix and $\|M\|_F$ is its Frobenius norm. For $M, M' \in \mathbb{R}^{d \times d}$, $M \preceq M'$ if and only if $u^\top(M' - M)u \geq 0$ for every $u \in \mathbb{R}^d$. For a set \mathcal{D} , we denote its power set as $\mathcal{P}(\mathcal{D}) \stackrel{\text{def}}{=} \{\mathcal{C} : \mathcal{C} \subseteq \mathcal{D}\}$. The sublevel of f at height $r \in \mathbb{R}$ is denoted $[f \leq r] \stackrel{\text{def}}{=} \{x \in \mathbb{R}^d : f(x) \leq r\}$.

2.1 On convex analysis

Let us recall some important definitions and results from convex analysis in the finite-dimensional case; for a comprehensive coverage, we refer the reader to [32].

We denote by $\Gamma_0(\mathbb{R}^d)$ the class of proper lsc and convex functions on \mathbb{R}^d taking values in $\mathbb{R} \cup \{+\infty\}$. For $\mu > 0$, $\Gamma_\mu(\mathbb{R}^d) \subset \Gamma_0(\mathbb{R}^d)$ is the class of μ -strongly convex functions. We denote by $C^s(\mathbb{R}^d)$ the class of s -times continuously differentiable functions on \mathbb{R}^d . For $L \geq 0$, $C_L^{1,1}(\mathbb{R}^d) \subset C^1(\mathbb{R}^d)$ is the set of functions on \mathbb{R}^d whose gradient is L -Lipschitz continuous.

The following *Descent Lemma* which is satisfied by this class of functions plays a central role in optimization.

Lemma 2.1. *Let $f \in C_L^{1,1}(\mathbb{R}^d)$, then*

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2, \quad \forall x, y \in \mathbb{R}^d.$$

Corollary 2.2. *Let $f \in C_L^{1,1}(\mathbb{R}^d)$ such that $\text{argmin } f \neq \emptyset$, then*

$$\|\nabla f(x)\|^2 \leq 2L(f(x) - \min f), \quad \forall x \in \mathbb{R}^d.$$

Proof. Use Lemma 2.1 for an arbitrary $x \in \mathbb{R}^d$ and $y = x - \frac{1}{L}\nabla f(x)$. Then bound

$$f\left(x - \frac{1}{L}\nabla f(x)\right) \geq \min f.$$

□

The *subdifferential* of a function $f \in \Gamma_0(\mathbb{R}^d)$ is the set-valued operator $\partial f : \mathbb{R}^d \rightarrow \mathcal{P}(\mathbb{R}^d)$ such that, for every x in \mathbb{R}^d ,

$$\partial f(x) = \{u \in \mathbb{R}^d : f(y) \geq f(x) + \langle u, y - x \rangle \quad \forall y \in \mathbb{R}^d\}.$$

When f is continuous, $\partial f(x)$ is non-empty convex and compact set for every $x \in \mathbb{R}^d$. If f is differentiable, then $\partial f(x) = \{\nabla f(x)\}$. For every $x \in \mathbb{R}^d$ such that $\partial f(x) \neq \emptyset$, the minimum norm selection of $\partial f(x)$ is the unique element $\partial^0 f(x) \stackrel{\text{def}}{=} \text{argmin}_{u \in \partial f(x)} \|u\|$.

2.2 On stochastic processes

Let us recall some elements of stochastic analysis; for a more complete account, we refer to [33, 34, 35]. Throughout the paper, $(\Omega, \mathcal{F}, \mathbb{P})$ is a probability space and $\{\mathcal{F}_t | t \geq 0\}$ is a filtration of the σ -algebra \mathcal{F} . Given $\mathcal{C} \in \mathcal{P}(\Omega)$, we will denote $\sigma(\mathcal{C})$ the σ -algebra generated by \mathcal{C} . We denote $\mathcal{F}_\infty \stackrel{\text{def}}{=} \sigma\left(\bigcup_{t \geq 0} \mathcal{F}_t\right) \in \mathcal{F}$.

The expectation of a random variable $\xi : \Omega \rightarrow \mathbb{R}^d$ is denoted by

$$\mathbb{E}(\xi) \stackrel{\text{def}}{=} \int_{\Omega} \xi(\omega) d\mathbb{P}(\omega).$$

An event $E \in \mathcal{F}$ happens almost surely if $\mathbb{P}(E) = 1$, and it will be denoted as " E , \mathbb{P} -a.s." or simply " E , a.s.". The characteristic function of an event $E \in \mathcal{F}$ is denoted by

$$\mathbb{1}_E(\omega) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } \omega \in E, \\ 0 & \text{otherwise.} \end{cases}$$

An \mathbb{R}^d -valued stochastic process is a function $X : \Omega \times \mathbb{R}_+ \rightarrow \mathbb{R}^d$. It is said to be continuous if $X(\omega, \cdot) \in C(\mathbb{R}_+; \mathbb{R}^d)$ for almost all $\omega \in \Omega$. We will denote $X(t) \stackrel{\text{def}}{=} X(\cdot, t)$. We are going to study (SDE), and in order to ensure the uniqueness of a solution, we introduce a relation over stochastic processes. Two stochastic processes $X, Y : \Omega \times [0, T] \rightarrow \mathbb{R}^d$ are said to be equivalent if $X(t) = Y(t), \forall t \in [0, T], \mathbb{P}$ -a.s. This leads us to define the equivalence relation \mathcal{R} , which associates the equivalent stochastic processes in the same class.

Furthermore, we will need some properties about the measurability of these processes. A stochastic process $X : \Omega \times \mathbb{R}_+ \rightarrow \mathbb{R}^d$ is progressively measurable if for every $t \geq 0$, the map $\Omega \times [0, t] \rightarrow \mathbb{R}^d$ defined by $(\omega, s) \rightarrow X(\omega, s)$ is $\mathcal{F}_t \otimes \mathcal{B}([0, t])$ -measurable, where \otimes is the product σ -algebra and \mathcal{B} is the Borel σ -algebra. On the other hand, X is \mathcal{F}_t -adapted if $X(t)$ is \mathcal{F}_t -measurable for every $t \geq 0$. It is a direct consequence of the definition that if X is progressively measurable, then X is \mathcal{F}_t -adapted.

Let us define the quotient space:

$$S_d^0[0, T] \stackrel{\text{def}}{=} \left\{ X : \Omega \times [0, T] \rightarrow \mathbb{R}^d : X \text{ is a prog. measurable cont. stochastic process} \right\} / \mathcal{R}.$$

We set $S_d^0 \stackrel{\text{def}}{=} \bigcap_{T \geq 0} S_d^0[0, T]$. Furthermore, for $\nu > 0$, we define $S_d^\nu[0, T]$ as the subset of processes $X(t)$ in $S_d^0[0, T]$ such that

$$S_d^\nu[0, T] \stackrel{\text{def}}{=} \left\{ X \in S_d^0[0, T] : \mathbb{E} \left(\sup_{t \in [0, T]} \|X_t\|^\nu \right) < +\infty \right\}.$$

We define $S_d^\nu \stackrel{\text{def}}{=} \bigcap_{T \geq 0} S_d^\nu[0, T]$.

Theorem A.7 in the appendix provides us with sufficient conditions to ensure the existence and uniqueness of the solution to (SDE). These conditions are met in our case under assumptions (H₀) and (H).

Let us now present Itô's formula which plays a central role in the theory of stochastic differential equations.

Proposition 2.3. [33, Chapter 4] Consider X the solution of (SDE), $\phi : \mathbb{R}_+ \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\phi(\cdot, x) \in C^1(\mathbb{R}_+)$ for every $x \in \mathbb{R}^d$ and $\phi(t, \cdot) \in C^2(\mathbb{R}^d)$ for every $t \geq 0$. Then the process

$$Y(t) = \phi(t, X(t)),$$

is an Itô Process such that for all $t \geq 0$

$$Y(t) = Y(0) + \int_0^t \frac{\partial \phi}{\partial t}(s, X(s)) ds - \int_0^t \langle \nabla \phi(s, X(s)), \nabla f(X(s)) \rangle ds \\ + \int_0^t \left\langle \sigma^\top(s, X(s)) \nabla \phi(s, X(s)), dW(s) \right\rangle + \frac{1}{2} \int_0^t \text{tr} \left(\sigma(s, X(s)) \sigma^\top(s, X(s)) \nabla^2 \phi(s, X(s)) \right) ds. \quad (2.1)$$

Moreover, if for all $T > 0$

$$\mathbb{E} \left(\int_0^T \|\sigma^\top(s, X(s)) \nabla \phi(s, X(s))\|^2 ds \right) < +\infty,$$

then $\int_0^t \left\langle \sigma^\top(s, X(s)) \nabla \phi(s, X(s)), dW(s) \right\rangle$ is a square-integrable continuous martingale and

$$\mathbb{E}[Y(t)] = Y(0) + \mathbb{E} \left(\int_0^t \frac{\partial \phi}{\partial t}(s, X(s)) ds \right) - \mathbb{E} \left(\int_0^t \langle \nabla \phi(s, X(s)), \nabla f(X(s)) \rangle ds \right) \\ + \frac{1}{2} \mathbb{E} \left(\int_0^t \text{tr} \left(\sigma(s, X(s)) \sigma^\top(s, X(s)) \nabla^2 \phi(s, X(s)) \right) ds \right). \quad (2.2)$$

The C^2 assumption on $\phi(t, \cdot)$ in Itô's formula is crucial. This can be weakened in certain cases leading to the following inequality that will be useful in our context.

Proposition 2.4. Consider X the solution of (SDE), $\phi_1 \in C^1(\mathbb{R}_+)$, $\phi_2 \in C_L^{1,1}(\mathbb{R}^d)$ and $\phi(t, x) = \phi_1(t)\phi_2(x)$. Then the process

$$Y(t) = \phi(t, X(t)) = \phi_1(t)\phi_2(X(t)),$$

is an Itô Process such that

$$Y(t) \leq Y(0) + \int_0^t \phi_1'(s)\phi_2(X(s)) ds - \int_0^t \phi_1(s) \langle \nabla \phi_2(X(s)), \nabla f(X(s)) \rangle ds \\ + \int_0^t \left\langle \sigma^\top(s, X(s)) \phi_1(s) \nabla \phi_2(X(s)), dW(s) \right\rangle + \frac{L}{2} \int_0^t \phi_1(s) \text{tr} \left(\sigma(s, X(s)) \sigma^\top(s, X(s)) \right) ds. \quad (2.3)$$

Moreover, if for all $T > 0$

$$\mathbb{E} \left(\int_0^T \|\sigma^\top(s, X(s)) \phi_1(s) \nabla \phi_2(X(s))\|^2 ds \right) < +\infty,$$

then

$$\mathbb{E}[Y(t)] \leq Y(0) + \mathbb{E} \left(\int_0^t \phi_1'(s)\phi_2(X(s)) ds \right) - \mathbb{E} \left(\int_0^t \phi_1(s) \langle \nabla \phi_2(X(s)), \nabla f(X(s)) \rangle ds \right) \\ + \frac{L}{2} \mathbb{E} \left(\int_0^t \phi_1(s) \text{tr} \left(\sigma(s, X(s)) \sigma^\top(s, X(s)) \right) ds \right). \quad (2.4)$$

Proof. Proof. Analogous to the proof of [28, Proposition C.2]. □

3 Convergence properties for convex differentiable functions

We consider f (called the potential) and study the dynamic (SDE) under hypotheses (\mathbf{H}_0) (i.e. $f \in C_L^{1,1}(\mathbb{R}^d) \cap \Gamma_0(\mathbb{R}^d)$) and (\mathbf{H}) . Recall the definitions of σ_* and $\sigma_\infty(t)$ from (1.1). Observe that from (\mathbf{H}) one can take $\sigma_*^2 = md \sup_{t \geq 0, x \in \mathbb{R}^d} |\sigma_{ik}(t, x)|^2$. Throughout the rest of the paper, we will use the shorthand notation

$$\Sigma(t, x) \stackrel{\text{def}}{=} \sigma(t, x)\sigma(t, x)^\top.$$

3.1 Almost sure convergence of trajectory

Our first main result establish almost convergence of $X(t)$ to an \mathcal{S} -valued random variable as $t \rightarrow +\infty$.

Theorem 3.1. *Consider the dynamic (SDE) where f and σ satisfy the assumptions (\mathbf{H}_0) and (\mathbf{H}) . Then, there exists a unique solution $X \in S_d^\nu$ of (SDE), for every $\nu \geq 2$. Additionally, if $\sigma_\infty \in L^2(\mathbb{R}_+)$, then:*

- (i) $\sup_{t \geq 0} \mathbb{E}[\|X(t)\|^2] < +\infty$.
- (ii) $\forall x^* \in \mathcal{S}$, $\lim_{t \rightarrow \infty} \|X(t) - x^*\|$ exists a.s. and $\sup_{t \geq 0} \|X(t)\| < +\infty$ a.s.
- (iii) $\lim_{t \rightarrow \infty} \|\nabla f(X(t))\| = 0$ a.s. As a result, $\lim_{t \rightarrow \infty} f(X(t)) = \min f$ a.s.
- (iv) In addition to (iii), there exists an \mathcal{S} -valued random variable x^* such that $\lim_{t \rightarrow \infty} X(t) = x^*$ a.s.

Proof. The existence and uniqueness of a solution follows directly from the fact that the conditions of Theorem A.7 are satisfied under (\mathbf{H}_0) and (\mathbf{H}) . The architecture of the proof of Theorem 3.1 consists of three steps that we briefly describe:

- The first step is based on Itô's formula (Proposition 2.3). Theorem A.9 then allows us to conclude that for all $x^* \in \mathcal{S}$, $\lim_{t \rightarrow \infty} \|X(t) - x^*\|$ exists a.s. Then, a separability argument is used to conclude that almost surely, for every $x^* \in \mathcal{S}$, $\lim_{t \rightarrow \infty} \|X(t) - x^*\|$ exists.
- The second step consists in using another conclusion of Theorem A.9 to conclude that $\|\nabla f(X(\cdot))\|^2 \in L^1(\mathbb{R}_+)$ a.s. After proving that this function is eventually uniformly continuous, we proceed according to Barbalat's Lemma (see [36]) to conclude that $\lim_{t \rightarrow \infty} \|\nabla f(X(t))\| = 0$ a.s. As a consequence of the convexity of f we deduce that $\lim_{t \rightarrow \infty} f(X(t)) = \min f$ a.s.
- Finally, the third step consists in using Opial's Lemma to conclude that there exists an \mathcal{S} -valued random variable x^* such that $\lim_{t \rightarrow \infty} X(t) = x^*$ a.s.

- (i) Let x^* be taken arbitrarily in \mathcal{S} . Let us define the corresponding anchor function $\phi(x) = \frac{\|x - x^*\|^2}{2}$. Using Itô's formula we obtain

$$\begin{aligned} \phi(X(t)) &= \underbrace{\frac{\|X_0 - x^*\|^2}{2}}_{\xi} + \underbrace{\frac{1}{2} \int_0^t \text{tr}(\Sigma(s, X(s))) ds}_{A_t} - \underbrace{\int_0^t \langle \nabla f(X(s)), X(s) - x^* \rangle ds}_{U_t} \\ &\quad + \underbrace{\int_0^t \langle \sigma^\top(s, X(s))(X(s) - x^*), dW(s) \rangle}_{M_t}. \end{aligned} \tag{3.1}$$

Since $X \in S_d^2$ by Proposition 2.3, we have for every $T > 0$, that

$$\mathbb{E} \left(\int_0^T \|\sigma^\top(s, X(s))(X(s) - x^*)\|^2 ds \right) \leq \mathbb{E} \left(\sup_{t \in [0, T]} \|X(t) - x^*\|^2 \right) \int_0^T \sigma_\infty^2(s) ds < +\infty.$$

Therefore M_t is a square-integrable continuous martingale. It is also a continuous local martingale (see [35, Theorem 1.3.3]), which implies that $\mathbb{E}(M_t) = 0$.

Let us now take the expectation of (3.1). Using that

$$0 \leq \text{tr}(\Sigma(s, X(s))) \leq \sigma_\infty^2(s) \text{ and } \langle \nabla f(X(s)), X(s) - x^* \rangle \geq 0,$$

and taking the supremum over $t \geq 0$, we obtain that

$$\sup_{t \geq 0} \mathbb{E} \left(\frac{\|X(t) - x^*\|^2}{2} \right) \leq \frac{\|X_0 - x^*\|^2}{2} + \frac{1}{2} \int_0^\infty \sigma_\infty^2(s) ds < +\infty.$$

This shows the first claim.

- (ii) A_t and U_t are two continuous adapted increasing processes with $A_0 = U_0 = 0$ a.s. Since $\phi(X(t))$ is nonnegative and $\sup_{x \in \mathbb{R}^d} \|\sigma(\cdot, x)\|_F \in L^2(\mathbb{R}_+)$, we deduce that $\lim_{t \rightarrow \infty} A_t < +\infty$. Then, we can use Theorem A.9 to conclude that

$$\int_0^\infty \langle \nabla f(X(s)), X(s) - x^* \rangle ds < +\infty \quad a.s. \quad (3.2)$$

and

$$\forall x^* \in \mathcal{S}, \exists \Omega_{x^*} \in \mathcal{F}, \text{ such that } \mathbb{P}(\Omega_{x^*}) = 1 \text{ and } \lim_{t \rightarrow \infty} \|X(\omega, t) - x^*\| \text{ exists } \forall \omega \in \Omega_{x^*}. \quad (3.3)$$

Since \mathbb{R}^d is separable, there exists a countable set $Z \subseteq \mathcal{S}$, such that $\text{cl}(Z) = \mathcal{S}$. Let $\tilde{\Omega} = \bigcap_{z \in Z} \Omega_z$. Since Z is countable

$$\mathbb{P}(\tilde{\Omega}) = 1 - \mathbb{P} \left(\bigcup_{z \in Z} \Omega_z^c \right) \geq 1 - \sum_{z \in Z} \mathbb{P}(\Omega_z^c) = 1.$$

For arbitrary $x^* \in \mathcal{S}$, there exists a sequence $(z_k)_{k \in \mathbb{N}} \subseteq Z$ such that $z_k \rightarrow x^*$. In view of (3.3), for every $k \in \mathbb{N}$ there exists $\tau_k : \Omega_{z_k} \rightarrow \mathbb{R}_+$ such that

$$\lim_{t \rightarrow \infty} \|X(\omega, t) - z_k\| = \tau_k(\omega), \quad \forall \omega \in \Omega_{z_k}.$$

Moreover, $\lim_{k \rightarrow \infty} \tau_k(\omega)$ exists since $(z_k)_{k \in \mathbb{N}}$ is convergent. Now, let $\omega \in \tilde{\Omega}$. Using the triangle inequality, we obtain that

$$\left| \|X(\omega, t) - z_k\| - \|X(\omega, t) - x^*\| \right| \leq \|z_k - x^*\|.$$

Taking $\limsup_{t \rightarrow \infty}$ over the previous inequality, we conclude that

$$\left| \tau_k(\omega) - \limsup_{t \rightarrow \infty} \|X(\omega, t) - x^*\| \right| \leq \|z_k - x^*\|.$$

A similar conclusion holds for the $\liminf_{t \rightarrow \infty}$. Then, taking the limit over k , we deduce

$$\lim_{t \rightarrow \infty} \|X(\omega, t) - x^*\| = \lim_{k \rightarrow \infty} \tau_k(\omega), \quad \forall \omega \in \tilde{\Omega},$$

whence we obtain that the previous limit exists on a set of probability 1 independently of x^* .

Let us recall that there exists $\Omega_c \in \mathcal{F}$ such that $\mathbb{P}(\Omega_c) = 1$ and $X(\omega, \cdot)$ is continuous for every $\omega \in \Omega_c$. Now let $x^* \in \mathcal{S}$ arbitrary, since the limit exists, for every $\omega \in \tilde{\Omega} \cap \Omega_c$ there exists $T(\omega)$ such that $\|X(\omega, t) - x^*\| \leq 1$ for every $t \geq T(\omega)$. Besides, since $X(\omega, \cdot)$ is continuous, by Bolzano's theorem $\sup_{t \in [0, T(\omega)]} \|X(\omega, t)\| = \max_{t \in [0, T(\omega)]} \|X(\omega, t)\| \stackrel{\text{def}}{=} h(\omega) < +\infty$. Therefore, $\sup_{t \geq 0} \|X(t)\| < \max\{h(\omega), 1 + \|x^*\|\} < \infty$.

(iii) By convexity of f and (3.2), we have that there exists $\Omega_f \in \mathcal{F}$ such that $\mathbb{P}(\Omega_f) = 1$ and $f(X(\omega, \cdot)) - \min f \in L^1(\mathbb{R}_+)$ for every $\omega \in \Omega_f$. By Corollary 2.2, we obtain that $\|\nabla f(X(\omega, \cdot))\| \in L^2(\mathbb{R}_+)$ for every $\omega \in \Omega_f$. Let $\omega \in \Omega_f$ arbitrary, then $\liminf_{t \rightarrow \infty} \|\nabla f(X(\omega, t))\| = 0$. If $\limsup_{t \rightarrow \infty} \|\nabla f(X(\omega, t))\| = 0$ then we conclude. Suppose by contradiction that $\limsup_{t \rightarrow \infty} \|\nabla f(X(\omega, t))\| > 0$. Then, by Lemma A.4, there exists $\delta > 0$ satisfying

$$0 = \liminf_{t \rightarrow \infty} \|\nabla f(X(\omega, t))\| < \delta < \limsup_{t \rightarrow \infty} \|\nabla f(X(\omega, t))\|,$$

and there exists $(t_k)_{k \in \mathbb{N}} \subset \mathbb{R}_+$ such that $\lim_{k \rightarrow \infty} t_k = \infty$,

$$\|\nabla f(X(\omega, t_k))\| > \delta \text{ and } t_{k+1} - t_k > 1, \quad \forall k \in \mathbb{N}.$$

Let $M_t = \int_0^t \sigma(s, X(s)) dW(s)$. This is a continuous martingale (w.r.t. the filtration \mathcal{F}_t), which verifies

$$\mathbb{E}(|M_t|^2) = \mathbb{E} \left(\int_0^t \|\sigma(s, X(s))\|_F^2 ds \right) \leq \mathbb{E} \left(\int_0^\infty \sigma_\infty^2(s) ds \right) < \infty, \forall t \geq 0.$$

According to Theorem A.8, we deduce that there exists a random variable M_∞ w.r.t. \mathcal{F}_∞ , and which verifies: $\mathbb{E}(|M_\infty|^2) < +\infty$, and there exists $\Omega_M \in \mathcal{F}$ such that $\mathbb{P}(\Omega_M) = 1$ and

$$\lim_{t \rightarrow \infty} M_t(\omega) = M_\infty(\omega) \text{ for every } \omega \in \Omega_M.$$

Let $\Omega_{\text{conv}} \stackrel{\text{def}}{=} \tilde{\Omega} \cap \Omega_c \cap \Omega_f \cap \Omega_M$, hence $\mathbb{P}(\Omega_{\text{conv}}) = 1$. Take any $\omega_0 \in \Omega_{\text{conv}}$. We allow ourselves the abuse of notation $X(t) \stackrel{\text{def}}{=} X(\omega_0, t)$ during the rest of the proof from this point.

Let $\varepsilon \in]0, \min\left(\frac{\delta^2}{4L^2}, L\right)[$. Note that $([t_k, t_k + \frac{\varepsilon}{2L}] : k \in \mathbb{N})$ are disjoint intervals. On the other hand, according to the convergence property of M_t and the fact that $\|\nabla f(X(\cdot))\| \in L^2(\mathbb{R}_+)$, there exists $k' > 0$ such that for every $k \geq k'$

$$\sup_{t \geq t_k} |M_t - M_{t_k}|^2 < \frac{\varepsilon}{4} \text{ and } \int_{t_k}^\infty \|\nabla f(X(s))\|^2 ds \leq \frac{L}{2}.$$

Besides, for every $k \geq k'$, $t \in [t_k, t_k + \frac{\varepsilon}{2L}]$

$$\|X(t) - X(t_k)\|^2 \leq 2(t - t_k) \int_{t_k}^t \|\nabla f(X(s))\|^2 ds + 2|M_t - M_{t_k}|^2 \leq 2(t - t_k) \frac{L}{2} + \frac{\varepsilon}{2} \leq \varepsilon.$$

Since $f \in C_L^{1,1}(\mathbb{R}^d)$, we have that for every $k \geq k'$ and $t \in [t_k, t_k + \frac{\varepsilon}{2L}]$

$$\|\nabla f(X(t)) - \nabla f(X(t_k))\|^2 \leq L^2 \|X(t) - X(t_k)\|^2 \leq \left(\frac{\delta}{2}\right)^2.$$

Therefore, for every $k \geq k'$, $t \in [t_k, t_k + \frac{\varepsilon}{2L}]$

$$\|\nabla f(X(t))\| \geq \|\nabla f(X(t_k))\| - \underbrace{\|\nabla f(X(t)) - \nabla f(X(t_k))\|}_{\leq \frac{\delta}{2}} \geq \frac{\delta}{2}.$$

Finally,

$$\int_0^\infty \|\nabla f(X(s))\|^2 ds \geq \sum_{k \geq k'} \int_{t_k}^{t_k + \frac{\varepsilon}{2L}} \|\nabla f(X(s))\|^2 ds \geq \sum_{k \geq k'} \frac{\delta^2 \varepsilon}{8L} = \infty,$$

which contradicts $\|\nabla f(X(\cdot))\| \in L^2(\mathbb{R}_+)$. So,

$$\limsup_{t \rightarrow \infty} \|\nabla f(X(\omega, t))\| = \liminf_{t \rightarrow \infty} \|\nabla f(X(\omega, t))\| = \lim_{t \rightarrow \infty} \|\nabla f(X(\omega, t))\| = 0, \quad \forall \omega \in \Omega_{\text{conv}}.$$

Let $x^* \in \mathcal{S}$ and $\omega \in \Omega_{\text{conv}}$ taken arbitrary. By convexity and Cauchy-Schwarz inequality:

$$0 \leq f(X(\omega, t)) - \min f \leq \|\nabla f(X(\omega, t))\| \|X(\omega, t) - x^*\|.$$

The claim then follows as we have already obtained that $\lim_{t \rightarrow \infty} \|X(\omega, t) - x^*\|$ exists, and

$$\lim_{t \rightarrow \infty} \|\nabla f(X(\omega, t))\| = 0.$$

- (iv) Let $\omega \in \Omega_{\text{conv}}$ and $\bar{x}(\omega)$ be a sequential limit point of $X(\omega, t)$. Equivalently, there exists an increasing sequence $(t_k)_{k \in \mathbb{N}} \subset \mathbb{R}_+$ such that $\lim_{k \rightarrow \infty} t_k = \infty$ and

$$\lim_{k \rightarrow \infty} X(\omega, t_k) = \bar{x}(\omega).$$

Since $\lim_{t \rightarrow \infty} f(X(\omega, t)) = \min f$ and by continuity of f , we obtain directly that $\bar{x}(\omega) \in \mathcal{S}$. Finally by Opial's Lemma (see [37]) we conclude that there exists $x^*(\omega) \in \mathcal{S}$ such that $\lim_{t \rightarrow \infty} X(\omega, t) = x^*(\omega)$. In other words, since $\omega \in \Omega_{\text{conv}}$ was arbitrary, there exists an \mathcal{S} -valued random variable x^* such that $\lim_{t \rightarrow \infty} X(t) = x^*$ a.s. □

3.2 Convergence rates of the objective

Our first result, stated below, summarizes the global convergence rates in expectation satisfied by the trajectories of (SDE).

Theorem 3.2. *Consider the dynamic (SDE) where f and σ satisfy the assumptions (H₀) and (H). The following statements are satisfied by the solution trajectory $X \in S_a^2$ of (SDE):*

- (i) Let $\overline{f \circ X}(t) \stackrel{\text{def}}{=} t^{-1} \int_0^t f(X(s)) ds$ and $\overline{X}(t) = t^{-1} \int_0^t X(s) ds$. Then

$$\mathbb{E}(f(\overline{X}(t)) - \min f) \leq \mathbb{E}(\overline{f \circ X}(t) - \min f) \leq \frac{\text{dist}(X_0, \mathcal{S})^2}{2t} + \frac{\sigma_*^2}{2}, \quad \forall t > 0. \quad (3.4)$$

Besides, if σ_∞ is $L^2(\mathbb{R}_+)$, then

$$\mathbb{E}(f(\overline{X}(t)) - \min f) \leq \mathbb{E}(\overline{f \circ X}(t) - \min f) = \mathcal{O}\left(\frac{1}{t}\right). \quad (3.5)$$

- (ii) If moreover $f \in \Gamma_\mu(\mathbb{R}^d)$ with $\mu > 0$, then $\mathcal{S} = \{x^*\}$ and

$$\mathbb{E}(\|X(t) - x^*\|^2) \leq \|X_0 - x^*\|^2 e^{-2\mu t} + \frac{\sigma_*^2}{2\mu}, \quad \forall t \geq 0. \quad (3.6)$$

Besides, if σ_∞ is decreasing and vanishes at infinity, then for every $\lambda \in]0, 1[$:

$$\mathbb{E}(\|X(t) - x^*\|^2) \leq \|X_0 - x^*\|^2 e^{-2\mu t} + \frac{\sigma_*^2}{2\mu} e^{-2\mu(1-\lambda)t} + \sigma_\infty^2(\lambda t), \quad \forall t \geq 0. \quad (3.7)$$

Proof. (i) Let $x^* \in \mathcal{S}$. Let $g(t) = \phi(X(t)) = \frac{\|X(t) - x^*\|^2}{2}$ and $G(t) = \mathbb{E}(g(t))$. By applying Proposition 2.3 with ϕ , and using the convexity of f , we obtain

$$\begin{aligned} G(t) - G(0) &= \mathbb{E} \left(\int_0^t \langle \nabla f(X(s)), x^* - X(s) \rangle ds \right) + \frac{1}{2} \mathbb{E} \left(\int_0^t \text{tr}[\Sigma(s, X(s))] ds \right) \\ &\leq -\mathbb{E} \left(\int_0^t (f(X(s)) - \min f) ds \right) + \frac{1}{2} \mathbb{E} \left(\int_0^t \text{tr}[\Sigma(s, X(s))] ds \right) \\ &\leq -\mathbb{E} \left(\int_0^t (f(X(s)) - \min f) ds \right) + \frac{\sigma_*^2}{2} t. \end{aligned} \quad (3.8)$$

Then rearranging the terms in (3.8), using $G(t) \geq 0$, and dividing by $t > 0$, we obtain

$$\frac{1}{t} \mathbb{E} \left(\int_0^t (f(X(s)) - \min f) ds \right) \leq \frac{\|X_0 - x^*\|^2}{2t} + \frac{\sigma_*^2}{2}, \quad \forall t > 0. \quad (3.9)$$

Since x^* is arbitrary, by taking the infimum with respect to $x^* \in \mathcal{S}$ in (3.9), we obtain

$$\frac{1}{t} \mathbb{E} \left(\int_0^t (f(X(s)) - \min f) ds \right) \leq \frac{\text{dist}(X_0, \mathcal{S})^2}{2t} + \frac{\sigma_*^2}{2}, \quad \forall t > 0. \quad (3.10)$$

Moreover, if $\sigma_\infty \in L^2(\mathbb{R}_+)$, then using inequality (3.8), we have

$$G(t) - G(0) \leq -\mathbb{E} \left(\int_0^t (f(X(s)) - \min f) ds \right) + \frac{1}{2} \left(\int_0^{+\infty} \sigma_\infty^2(s) ds \right).$$

Rearranging as before, we conclude that

$$\frac{1}{t} \mathbb{E} \left(\int_0^t (f(X(s)) - \min f) ds \right) \leq \frac{\text{dist}(X_0, \mathcal{S})^2}{2t} + \frac{1}{2t} \int_0^{+\infty} \sigma_\infty^2(s) ds, \quad \forall t > 0. \quad (3.11)$$

Then complete the result with the inequality

$$\mathbb{E} (f(\overline{X}(t)) - \min f) \leq \mathbb{E} (\overline{f \circ X}(t) - \min f)$$

which follows from convexity of f and Jensen's inequality.

(ii) Let $g(t) = \phi(X(t)) = \frac{\|X(t) - x^*\|^2}{2}$, $G(t) = \mathbb{E}(g(t))$. By Proposition 2.3 with ϕ , we obtain

$$G(t) - G(0) = \mathbb{E} \left(\int_0^t \langle -\nabla f(X(s)), X(s) - x^* \rangle ds \right) + \frac{1}{2} \mathbb{E} \left(\int_0^t \text{tr}[\Sigma(s, X(s))] ds \right). \quad (3.12)$$

Using that $f \in \Gamma_\mu(\mathbb{R}^d)$, we deduce that

$$G(t) \leq G(0) - 2\mu \int_0^t G(s) ds + \int_0^t \frac{\sigma_*^2}{2} ds, \quad \forall t \geq 0.$$

In order to invoke Lemma A.2, we solve the ODE

$$\begin{cases} y'(t) &= -2\mu y(t) + \frac{\sigma_*^2}{2}, \quad t > 0 \\ y(0) &= \frac{\|X_0 - x^*\|^2}{2}. \end{cases}$$

Solving it by the integrating factor method, we conclude that

$$G(t) \leq \frac{\|X_0 - x^*\|^2}{2} e^{-2\mu t} + \frac{\sigma_*^2}{4\mu}, \quad \forall t \geq 0.$$

Combining this inequality with the gradient descent Lemma 2.1, we obtain

$$\mathbb{E}[f(X(t)) - \min f] \leq L \left(\frac{\|X_0 - x^*\|^2}{2} e^{-2\mu t} + \frac{\sigma_*^2}{4\mu} \right), \quad \forall t \geq 0. \quad (3.13)$$

Suppose now that σ_∞ is decreasing and vanishes at infinity. We can bound the trace term by σ_∞^2 in (3.12). To use Lemma A.2, we need to solve

$$\begin{cases} y'(t) &= -2\mu y(t) + \frac{\sigma_\infty^2(t)}{2}, \quad t > 0 \\ y(0) &= \frac{\|X_0 - x^*\|^2}{2}. \end{cases}$$

Let $\lambda \in]0, 1[$, using the integrating factor method, we get

$$\begin{aligned} y(t) &\leq y(0)e^{-2\mu t} + e^{-2\mu t} \int_0^t \frac{\sigma_\infty^2(s)}{2} e^{2\mu s} ds \\ &\leq y(0)e^{-2\mu t} + e^{-2\mu t} \left(\int_0^{\lambda t} \frac{\sigma_\infty^2(s)}{2} e^{2\mu s} ds + \int_{\lambda t}^t \frac{\sigma_\infty^2(s)}{2} e^{2\mu s} ds \right) \\ &\leq y(0)e^{-2\mu t} + e^{-2\mu t} \left(\frac{\sigma_*^2}{2} \int_0^{\lambda t} e^{2\mu s} ds + \frac{\sigma_\infty^2(\lambda t)}{2} \int_{\lambda t}^t e^{2\mu s} ds \right) \\ &\leq y(0)e^{-2\mu t} + e^{-2\mu t} \left(\frac{\sigma_*^2}{4\mu} e^{2\mu \lambda t} + \frac{\sigma_\infty^2(\lambda t)}{2} e^{2\mu t} \right), \quad \forall t \geq 0. \end{aligned}$$

According to Lemma A.2, we deduce that

$$G(t) \leq \frac{\|X_0 - x^*\|^2}{2} e^{-2\mu t} + \frac{\sigma_*^2}{4\mu} e^{-2\mu(1-\lambda)t} + \frac{\sigma_\infty^2(\lambda t)}{2}, \quad \forall t \geq 0,$$

which is our claim (3.7). □

Under a stronger assumption on σ_∞ , we also have the following pointwise sublinear convergence rate in expectation.

Proposition 3.3. *Consider the dynamic (SDE) where f and σ satisfy the assumptions (H₀) and (H). Assume that there exists $K \geq 0, \beta \in [0, 1[$ such that*

$$\int_0^t (s+1)\sigma_\infty^2(s) ds \leq Kt^\beta, \quad \forall t \geq 0. \quad (3.14)$$

Then the solution trajectory $X \in S_d^2$ of (SDE) satisfies

$$\mathbb{E}(f(X(t)) - \min f) = \mathcal{O}(t^{\beta-1}).$$

Proof. Given $x^* \in \mathcal{S}$, let us apply Proposition 2.4 successively with $V_1(t, x) = t(f(x) - \min f)$, then with $V_2(x) = \frac{1}{2}\|x - x^*\|^2$. Taking the expectation and adding the two results, we get

$$\begin{aligned} \mathbb{E}(V_1(t, X(t)) + V_2(X(t))) &\leq \frac{1}{2}\|X_0 - x^*\|^2 + \frac{L}{2} \int_0^t s\sigma_\infty^2(s)ds + \frac{1}{2} \int_0^t \sigma_\infty^2(s)ds \\ &\leq \frac{1}{2}\|X_0 - x^*\|^2 + \frac{\max\{1, L\}}{2} \left(\int_0^t (s+1)\sigma_\infty^2(s)ds \right), \end{aligned}$$

where we have used the convexity of f in the first inequality. Then we conclude that

$$\mathbb{E}(f(X(t)) - \min f) \leq \frac{\|X_0 - x^*\|^2}{2t} + \frac{K \max\{1, L\}}{2} t^{\beta-1} = \mathcal{O}(t^{\beta-1}).$$

□

When f is also C^2 , we get an improved $o(t^{-1})$ global convergence rate on the objective in almost sure sense.

Theorem 3.4. Consider the dynamic (SDE). Assume that $f \in C^2(\mathbb{R}^d)$ such that $\nabla^2 f \preceq LI_d$ and satisfies assumption (H₀), and that σ satisfies assumption (H) and that $t \mapsto t\sigma_\infty^2(t) \in L^1(\mathbb{R}_+)$. Then, the solution trajectory $X \in S_d^2$ of (SDE) obeys:

- (i) $t \mapsto t\|\nabla f(X(t))\|^2 \in L^1(\mathbb{R}_+)$ a.s.
- (ii) $f(X(t)) - \min f = o(t^{-1})$ a.s.

Proof. By applying Itô's formula in Proposition 2.3 with $\phi(t, x) = t(f(x) - \min f)$ we get

$$\begin{aligned} t(f(X(t)) - \min f) &= \int_0^t f(X(s)) - \min f ds + \frac{1}{2} \int_0^t \text{str}[\Sigma(s, X(s)) \nabla^2 f(X(s))] ds \\ &\quad - \int_0^t s \|\nabla f(X(s))\|^2 ds + \int_0^t \langle s\sigma^\top(s, X(s)) \nabla f(X(s)), dW(s) \rangle. \end{aligned}$$

By (3.2) and convexity of f , we deduce that $f(X(\cdot)) - \min f \in L^1(\mathbb{R}_+)$ a.s. Moreover,

$$\int_0^\infty \text{str}[\Sigma(s, X(s)) \nabla^2 f(X(s))] ds \leq L \int_0^\infty s\sigma_\infty^2(s) ds < +\infty.$$

Then by Theorem A.9, we have that $\lim_{t \rightarrow \infty} t(f(X(t)) - \min f)$ exists a.s. and $\int_0^\infty t\|\nabla f(X(t))\|^2 dt < +\infty$ a.s. Finally, by Lemma A.1, we conclude that $\lim_{t \rightarrow \infty} t(f(X(t)) - \min f) = 0$ a.s. □

4 Convergence rates under Łojasiewicz inequality

The local convergence rate of the first-order descent methods can be understood using the Łojasiewicz property and the associated Łojasiewicz exponent, see [38, 39]. The Łojasiewicz property has its roots in algebraic geometry, and it essentially describes a relationship between the objective value and its gradient (or subgradient).

Definition 4.1 (Łojasiewicz inequality). Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a differentiable with $\mathcal{S} = \operatorname{argmin}(f) \neq \emptyset$ and $q \in [0, 1[$. f satisfies the Łojasiewicz inequality with exponent q at $\bar{x} \in \mathcal{S}$ if there exists a neighborhood $\mathcal{V}_{\bar{x}}$ of \bar{x} , $r > \min f$ and $\mu > 0$ such that

$$\mu(f(x) - \min f)^q \leq \|\nabla f(x)\|, \quad \forall x \in \mathcal{V}_{\bar{x}} \cap [\min f < f < r]. \quad (4.1)$$

The function f has the Łojasiewicz property on \mathcal{S} if it obeys (4.1) at each point of \mathcal{S} with the same constant μ and exponent q , and we will write $f \in \mathcal{L}^q(\mathcal{S})$.

Error bounds have also been successfully applied to various branches of optimization, and in particular to complexity analysis, see [40]. Of particular interest in our setting is the Hölderian error bound.

Definition 4.2 (Hölderian error bound). Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a proper function such that $\mathcal{S} = \operatorname{argmin}(f) \neq \emptyset$. f satisfies a Hölderian (or power-type) error bound inequality with exponent $p \geq 1$, and we write $f \in \text{EB}^p$, if there exists $\gamma > 0$ and $r > \min f$ such that

$$f(x) - \min f \geq \gamma \operatorname{dist}(x, \mathcal{S})^p, \quad \forall x \in [\min f \leq f \leq r]. \quad (4.2)$$

For a given $r > \min f$ such that (4.2) holds, we will use the shorthand notation $f \in \text{EB}^p([f \leq r])$.

A deep result due Łojasiewicz states that for arbitrary continuous semi-algebraic functions, the Hölderian error bound inequality holds on any compact set, and the Łojasiewicz inequality holds at each point; see [3, 4]. In fact, for convex functions, the Łojasiewicz property and Hölderian error bound are actually equivalent.

Proposition 4.3. Assume that $f \in \Gamma_0(\mathbb{R}^d) \cap C^1(\mathbb{R}^d)$ with $\mathcal{S} = \operatorname{argmin}(f) \neq \emptyset$. Let $q \in [0, 1[$, $p \stackrel{\text{def}}{=} \frac{1}{1-q} \geq 1$ and $r > \min f$. Then f verifies the Łojasiewicz inequality (4.1) at $\bar{x} \in \mathcal{S}$ if and only if the Hölderian error bound (4.2) holds on $\mathcal{V}_{\bar{x}} \cap [\min f < f < r]$.

Proof. Combine [10, Lemma 4 and Theorem 5]. □

We are now ready to state the following ergodic local convergence rate.

Proposition 4.4. Consider the hypotheses of Theorem 3.2 and let $\varepsilon > 0$. If $f \in \text{EB}^p([f \leq r_\varepsilon])$ for $r_\varepsilon > \min f + \frac{\sigma_*^2}{2} + \varepsilon$, then $\exists t_\varepsilon > 0$ such that

$$\operatorname{dist}(\mathbb{E}[\bar{X}(t)], \mathcal{S}) = \mathcal{O}(t^{-\frac{1}{p}}) + \mathcal{O}\left(\sigma_*^{\frac{2}{p}}\right), \quad \forall t \geq t_\varepsilon.$$

Proof. There exists $t_\varepsilon > 0$ such that for all $t \geq t_\varepsilon$, $\frac{\operatorname{dist}(X_0, \mathcal{S})^\varepsilon}{2t} < \varepsilon$. Thus, from (3.4) and Jensen's Inequality, we have

$$f(\mathbb{E}[\bar{X}(t)]) \leq \mathbb{E}[f(\bar{X}(t))] \leq \min f + \frac{\sigma_*^2}{2} + \varepsilon \leq r_\varepsilon, \quad \forall t \geq t_\varepsilon.$$

Clearly, $\mathbb{E}[\bar{X}(t)] \in [f \leq r_\varepsilon]$ for $t \geq t_\varepsilon$. Using Theorem 3.2 and that $f \in \text{EB}^p([f \leq r_\varepsilon])$, letting $\gamma > 0$ the coefficient of the error bound, we have

$$\gamma \operatorname{dist}(\mathbb{E}[\bar{X}(t)], \mathcal{S})^p \leq f(\mathbb{E}[\bar{X}(t)]) - \min f \leq \frac{\operatorname{dist}(X_0, \mathcal{S})^2}{2t} + \frac{\sigma_*^2}{2}, \quad \forall t \geq t_\varepsilon.$$

Since $p \geq 1$, Jensen's inequality yields

$$\operatorname{dist}(\mathbb{E}[\bar{X}(t)], \mathcal{S}) \leq \left(\frac{\operatorname{dist}(X_0, \mathcal{S})^2}{2\gamma r}\right)^{\frac{1}{p}} t^{-\frac{1}{p}} + \left(\frac{\sigma_*^2}{2\gamma r}\right)^{\frac{1}{p}}, \quad \forall t \geq t_\varepsilon.$$

□

4.1 Discussion on the localization of the process

Let us take a moment to elaborate on the localization of the process $X(t)$ generated by (SDE) when $f \in C_L^{1,1}(\mathbb{R}^d) \cap \Gamma_0(\mathbb{R}^d)$ and $\sigma_\infty \in L^2(\mathbb{R}_+)$. This discussion is essential to understand the challenges underlying the analysis of the local convergence properties and rates in a stochastic setting under (local) error bounds. First, observe that the hypothesis of Lipschitz continuity of the gradient is incompatible with a global hypothesis of error bound or Łojasiewicz inequality unless the exponent is $p = 2$ or $q = \frac{1}{2}$, respectively. Therefore, we can only ask for these inequalities to be locally satisfied. Even though, thanks to convexity, we could introduce a global desingularizing function (see [10, Theorem 3]), this function would not be concave nor convex, a fundamental property usually at the heart of the local analysis. In recent literature on stochastic processes and local properties, it is usual to find hypotheses about the almost sure localization of the process or that it is essentially bounded. Nevertheless, these assumptions are unrealistic or outright false due to the behavior of the Brownian Motion. Hence, we will avoid making these kinds of assumptions.

What we will do is to consider that by Theorem 3.1 we have that $\lim_{t \rightarrow \infty} f(X(t)) = \min f$ a.s., which means that there exists $\Omega_{\text{conv}} \in \mathcal{F}$ such that $\mathbb{P}(\Omega_{\text{conv}}) = 1$, and $(\forall r > \min f, \forall \omega \in \Omega_{\text{conv}}, (\exists t_r(\omega) > 0)$ such that $(\forall t > t_r(\omega)), X(\omega, t) \in [f \leq r]$. However, one should not infer from this that $X(t) \in [f \leq r]$ a.s. for t large enough. Indeed, t_r is a random variable which cannot be in general bounded uniformly on Ω_{conv} . Unfortunately, this flawed argument appears quite regularly in the literature. Rather, in this paper, we will invoke measure theoretic arguments to pass from a.s. convergence to almost uniform convergence thanks to Egorov's theorem (see Theorem A.3). More precisely, we will show that

$$(\forall \delta > 0, \forall r > \min f), (\exists \Omega_\delta \in \mathcal{F} \text{ s.t. } \mathbb{P}(\Omega_\delta) \geq 1 - \delta \text{ and } \exists \hat{t}_{r,\delta} > 0), (\forall \omega \in \Omega_\delta, \forall t > \hat{t}_{r,\delta}), \\ X(\omega, t) \in [f \leq r].$$

Hence, this property will allow us to localize $X(t)$ in the sublevel set of f at r for t large enough with probability at least $1 - \delta$. In turn, we will be able to invoke the error bound (or Łojasiewicz) inequality.

4.2 Convergence rates under Łojasiewicz Inequality

Let $\sigma_\infty \in L^2(\mathbb{R}_+)$, $L > 0, \delta > 0, \beta \in [0, 1[$ and some positive constants C_*, C_{**}, C_K . Consider the functions $h_\delta, l_\delta, k_\delta : \mathbb{R}_+ \rightarrow \mathbb{R}$ defined by:

$$h_\delta(t) = \sigma_\infty^2(t) + C_* \sqrt{\delta} \frac{\sigma_\infty^2(t)}{2 \sqrt{\int_{\hat{t}_\delta}^t \sigma_\infty^2(u) du}}, \quad (4.3)$$

$$l_\delta(t) = \frac{L}{2} \sigma_\infty^2(t) + C_{**} \sqrt{\delta} \frac{\sigma_\infty^2(t)}{2 \sqrt{\int_{\hat{t}_\delta}^t \sigma_\infty^2(u) du}}, \quad (4.4)$$

$$k_\delta(t) = \frac{L}{2} \sigma_\infty^2(t) + C_K \sqrt{\delta} \frac{\sigma_\infty^2(t) t^{\beta-1}}{2 \sqrt{\int_{\hat{t}_\delta}^t \sigma_\infty^2(u) u^{\beta-1} du}}. \quad (4.5)$$

We are now ready to state our main local convergence result.

Theorem 4.5. *Consider $X \in S_d^2$ the solution trajectory of (SDE) where f and σ satisfy the assumptions (\mathbf{H}_0) and (\mathbf{H}) , and suppose that $\sigma_\infty \in L^2(\mathbb{R}_+)$ ($C_\infty \stackrel{\text{def}}{=} \|\sigma_\infty\|_{L^2(\mathbb{R}_+)}$). Let $p \geq 2$ and $q \stackrel{\text{def}}{=} 1 - \frac{1}{p} \in [\frac{1}{2}, 1[$, and assume that $f \in \mathcal{L}^q(\mathcal{S})$. Consider also the positive constants $C_*, C_{**}, C_K, C_d, C_f$ (detailed in the proof).*

Then, for all $\delta > 0$, there exists a measurable set Ω_δ such that $\mathbb{P}(\Omega_\delta) \geq 1 - \delta$ and $\hat{t}_\delta > 0$ such that the following statements hold.

(i) If $p = 2$ and σ_∞ is decreasing, then σ_∞ vanishes at infinity and

(a) there exists $\gamma > 0$ such that for every $\lambda \in]0, 1[$,

$$\begin{aligned} \mathbb{E} \left(\frac{\text{dist}(X(t), \mathcal{S})^2}{2} \right) &\leq e^{-2\gamma(t-\hat{t}_\delta)} \mathbb{E} \left(\frac{\text{dist}(X(\hat{t}_\delta), \mathcal{S})^2}{2} \right) \\ &\quad + e^{-2\gamma(1-\lambda)(t-\hat{t}_\delta)} (C_\infty^2 + C_* C_\infty \sqrt{\delta}) \\ &\quad + \frac{h_\delta(\hat{t}_\delta + \lambda(t - \hat{t}_\delta))}{2\gamma} + C_d \sqrt{\delta}, \quad \forall t > \hat{t}_\delta; \end{aligned} \quad (4.6)$$

(b) there exists $\mu > 0$ such that for every $\lambda \in]0, 1[$,

$$\begin{aligned} \mathbb{E} (f(X(t)) - \min f) &\leq e^{-\mu^2(t-\hat{t}_\delta)} \mathbb{E} (f(X(\hat{t}_\delta)) - \min f) \\ &\quad + e^{-\mu^2(1-\lambda)(t-\hat{t}_\delta)} \left(\frac{LC_\infty^2}{2} + C_{**} C_\infty \sqrt{\delta} \right) \\ &\quad + \frac{l_\delta(\hat{t}_\delta + \lambda(t - \hat{t}_\delta))}{\mu^2} + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta. \end{aligned} \quad (4.7)$$

Moreover, if (3.14) holds, then

$$\begin{aligned} \mathbb{E} (f(X(t)) - \min f) &\leq e^{-\mu^2(t-\hat{t}_\delta)} \mathbb{E} (f(X(\hat{t}_\delta)) - \min f) \\ &\quad + e^{-\mu^2(1-\lambda)(t-\hat{t}_\delta)} \left(\frac{LC_\infty^2}{2} + C_K C_\infty \sqrt{\hat{t}_\delta^{\beta-1} \sqrt{\delta}} \right) \\ &\quad + \frac{k_\delta(\hat{t}_\delta + \lambda(t - \hat{t}_\delta))}{\mu^2} + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta. \end{aligned} \quad (4.8)$$

(ii) If $p > 2$:

(a) There exists $\gamma > 0$ such that

$$\mathbb{E} \left(\frac{\text{dist}(X(t), \mathcal{S})^2}{2} \right) \leq y_\delta^*(t) + C_d \sqrt{\delta}, \quad \forall t > \hat{t}_\delta, \quad (4.9)$$

where y_δ^* is the solution of the Cauchy problem

$$(C.1) \quad \begin{cases} y'(t) &= -2^{\frac{p}{2}} \gamma y^{\frac{p}{2}} + h_\delta(t), \quad t > \hat{t}_\delta \\ y(\hat{t}_\delta) &= \mathbb{E} \left(\frac{\text{dist}(X(\hat{t}_\delta), \mathcal{S})^2}{2} \mathbb{1}_{\Omega_\delta} \right). \end{cases}$$

(b) There exists $\mu > 0$ such that

$$\mathbb{E} [f(X(t)) - \min f] \leq w_\delta^*(t) + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta, \quad (4.10)$$

where w_δ^* is the solution of the Cauchy problem

$$(C.2) \quad \begin{cases} y'(t) &= -\mu^2 y(t)^{2q} + l_\delta(t), \quad t > \hat{t}_\delta \\ y(\hat{t}_\delta) &= \mathbb{E} ([f(X(\hat{t}_\delta)) - \min f] \mathbb{1}_{\Omega_\delta}). \end{cases}$$

Moreover, if (3.14) holds, then

$$\mathbb{E} [f(X(t)) - \min f] \leq z_\delta^*(t) + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta, \quad (4.11)$$

where z_δ^* is the solution of the Cauchy problem

$$(C.3) \begin{cases} y'(t) &= -\mu^2 y(t)^{2q} + k_\delta(t), \quad t > \hat{t}_\delta \\ y(\hat{t}_\delta) &= \mathbb{E}([f(X(\hat{t}_\delta) - \min f] \mathbb{1}_{\Omega_\delta}). \end{cases}$$

Before proceeding with the proof, a few remarks are in order.

Remark 4.6. The hypothesis that f has a Lipschitz continuous gradient restricts the Łojasiewicz exponent q to be in $[\frac{1}{2}, 1]$.

Remark 4.7. If we have a *global* error bound (or Łojasiewicz inequality), then as noted in the discussion of Section 4.1, one necessarily has $p = 2$ (or $q = \frac{1}{2}$). In this case, the statements (i) of Theorem 4.5 will hold if we replace $\sigma_\infty \in L^2(\mathbb{R}_+)$ by σ_∞ decreasing and vanishing at infinity, δ by 0 and \hat{t}_δ by 0. Clearly, one recovers (3.7).

Remark 4.8. It is important to highlight the trade-off in the selection of δ . Although δ can be arbitrarily small, the time from which the inequalities are satisfied, \hat{t}_δ , surely increases when δ approaches 0^+ . Besides, let $q_{\delta, \hat{t}_\delta} : \mathbb{R}_+ \rightarrow \mathbb{R}$ be a decreasing function. Our convergence rates in Theorem 4.5 are of the form $\mathbb{E}[m(X(t))] \leq q_{\delta, \hat{t}_\delta}(t) + C\sqrt{\delta}$, $\forall t > t_\delta$, where $m(x) = f(x) - \min f$ or $m(x) = \text{dist}(x, \mathcal{S})^2/2$. Let $\varepsilon \in]0, 2C[$ and $\delta^* = \frac{\varepsilon^2}{4C^2}$. Then one gets an ε -optimal solution for $t > \max\{q^*(\varepsilon), \hat{t}_{\delta^*}\}$.

Remark 4.9. Referring again to the discussion of Section 4.1, we have that there exists $\delta > 0$ and $\Omega_\delta \in \mathcal{F}$ with $\mathbb{P}(\Omega_\delta) \geq 1 - \delta$ over which we have uniform convergence of the objective. If δ could be 0 (a.s. uniform convergence), there would be a $\hat{t} > 0$ such that $X(t) \in [f \leq r], \forall t > \hat{t}$ a.s. Thus, the statements in Theorem 4.5 would hold if we replace δ by 0 and \hat{t}_δ by \hat{t} . The proof is far easier in this case. It is however not easy to ensure the existence of such \hat{t} in general.

Remark 4.10. In order to find explicit convergence rates in Theorem 4.5 we have to solve or bound the solution of the Cauchy problems (C.1), (C.2) and (C.3). We can generalize these problems as follows: Let $a > 0, b > 1, \hat{t}_\delta > 0, \delta > 0, y_0(\hat{t}_\delta, \delta) > 0$ and p_δ a nonnegative integrable function. Consider

$$(C.0) \begin{cases} y'(t) &= -ay^b(t) + p_\delta(t), \quad t > \hat{t}_\delta \\ y(\hat{t}_\delta) &= y_0(\hat{t}_\delta, \delta). \end{cases}$$

Although one could give an explicit ad-hoc p_δ in order to find a particular solution of (C.0), the dependence of this function on \hat{t}_δ is unavoidable, which is a problem, since p_δ is explicitly related to σ_∞ , and this in turn is the one that defines \hat{t}_δ in the first place.

To the best of our knowledge, there is no way to arithmetically solve this non linear ODE, not even a sharp bound of the solution.

Nevertheless, if $y(t) = \mathcal{O}\left((t+1)^{-\frac{1}{b-1}}\right)$, then $p_\delta(t) = \mathcal{O}\left((t+1)^{-\frac{b}{b-1}}\right)$. Which leads us to make the following conjecture:

Conjecture 4.11. *If $p_\delta = \mathcal{O}(\sigma_\infty^2)$ and $\sigma_\infty^2(t) = \mathcal{O}\left((t+1)^{-\frac{b}{b-1}}\right)$ (for constants independent of δ and \hat{t}_δ), then $y(t) = \mathcal{O}\left((t+1)^{-\frac{1}{b-1}}\right)$.*

Proof. Proof of Theorem 4.5. Given that $\sigma_\infty \in L^2(\mathbb{R}_+)$, if it is decreasing, we have immediately that it vanishes at infinity. Let $x^* \in \mathcal{S}$. Let us recall that by claim (i) of Theorem 3.1, there exists $C^* > 0$ such that

$$\sup_{t \geq 0} \mathbb{E}(\text{dist}(X(t), \mathcal{S})^2) \leq \sup_{t \geq 0} \mathbb{E}(\|X(t) - x^*\|^2) \leq C^*.$$

On the other hand, by Theorem 3.1(iii), there exists a set $\Omega_{\text{conv}} \in \mathcal{F}$ such that $\mathbb{P}(\Omega_{\text{conv}}) = 1$ where, for all $\omega \in \Omega_{\text{conv}}$: $\lim_{t \rightarrow \infty} f(X(\omega, t)) = \min f$, $t \mapsto f(X(\omega, t))$ is continuous, and $\lim_{t \rightarrow \infty} \text{dist}(X(\omega, t), \mathcal{S}) = 0$. Then, by Theorem A.3 for every $\delta > 0$ there exists $\Omega_\delta \in \mathcal{F}$ such that $\Omega_\delta \subset \Omega_{\text{conv}}$, $\mathbb{P}(\Omega_\delta) > 1 - \delta$ and $f(X(\cdot, t))$ (resp. $\text{dist}(X(\cdot, t), \mathcal{S})$) converges uniformly to $\min f$ (resp. to 0) on Ω_δ . This means that given $r \geq \min f$, and for every $\delta > 0$, there exist $\hat{t}_\delta > 0$ and $\Omega_\delta \in \mathcal{F}$ with $\mathbb{P}(\Omega_\delta) > 1 - \delta$ such that $X(\omega, t) \in [f \leq r] \cap \mathcal{V}_\mathcal{S}$ for all $t \geq \hat{t}_\delta$ and $\omega \in \Omega_\delta$, where $\mathcal{V}_\mathcal{S}$ is a neighbourhood of \mathcal{S} . On the other hand, since $f \in \mathbf{L}^q(\mathcal{S})$, by Proposition 4.3, there exists $r > \min f$ and a neighbourhood $\mathcal{V}_\mathcal{S}$ of \mathcal{S} such that f verifies the p -Hölderian error bound inequality (4.2) on $[\min f < f < r] \cap \mathcal{V}_\mathcal{S}$. Consequently, for any $\delta > 0$, there exists $t \geq \hat{t}_\delta$ large enough such that the p -Hölderian error bound inequality holds at $X(\omega, t)$ for all $t \geq \hat{t}_\delta$ and $\omega \in \Omega_\delta$.

We are now ready to start. Let $x^* \in \mathcal{S}$, $\delta > 0$, and $t \geq \hat{t}_\delta$.

(i) $p = 2$:

(a) Let $\hat{g}(t) = \hat{\phi}(X(t)) = \frac{\text{dist}(X(t), \mathcal{S})^2}{2}$, $\hat{G}(t) = \mathbb{E}(\hat{g}(t) \mathbb{1}_{\Omega_\delta})$, and $\mu > 0$ be the coefficient of the error bound inequality. We have

$$\nabla \hat{\phi}(X(t)) = X(t) - P_\mathcal{S}(X(t)),$$

where $P_\mathcal{S}(x)$ is the projection of x on \mathcal{S} , so $\hat{\phi} \in C_1^{1,1}(\mathbb{R}^d)$. We use Proposition 2.4 to obtain

$$\begin{aligned} \hat{g}(t) - \hat{g}(\hat{t}_\delta) &\leq - \int_{\hat{t}_\delta}^t \langle \nabla f(X(s), X(s) - P_\mathcal{S}(X(s))) \rangle ds \\ &\quad + \int_{\hat{t}_\delta}^t \text{tr}[\Sigma(s, X(s))] ds + \int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s))(X(s) - P_\mathcal{S}(X(s))), dW(s) \right\rangle. \end{aligned} \quad (4.12)$$

We have that $\text{tr}[\Sigma(s, X(s))] \leq \sigma_\infty^2(s)$ and by convexity

$$- \langle \nabla f(X(s), X(s) - P_\mathcal{S}(X(s))) \rangle \leq - (f(X(s)) - \min f).$$

Therefore,

$$\begin{aligned} \hat{g}(t) - \hat{g}(\hat{t}_\delta) &\leq - \int_{\hat{t}_\delta}^t (f(X(s)) - \min f) ds \\ &\quad + \int_{\hat{t}_\delta}^t \sigma_\infty^2(s) ds + \int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s))(X(s) - P_\mathcal{S}(X(s))), dW(s) \right\rangle. \end{aligned}$$

Then, multiplying this inequality by $\mathbb{1}_{\Omega_\delta}$, and taking expectation we obtain

$$\begin{aligned} \hat{G}(t) - \hat{G}(\hat{t}_\delta) &\leq - \mathbb{E} \left[\int_{\hat{t}_\delta}^t (f(X(s)) - \min f) \mathbb{1}_{\Omega_\delta} ds \right] + \int_{\hat{t}_\delta}^t \sigma_\infty^2(s) ds \\ &\quad + \mathbb{E} \left[\mathbb{1}_{\Omega_\delta} \int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s))(X(s) - P_\mathcal{S}(X(s))), dW(s) \right\rangle \right]. \end{aligned}$$

On the other hand, since $\sigma_\infty \in L^2(\mathbb{R}_+)$, we have for all $T > 0$

$$\begin{aligned} \mathbb{E} \left(\int_0^T \|\sigma^\top(s, X(s))(X(s) - P_{\mathcal{S}}(X(s)))\|^2 ds \right) &\leq \mathbb{E} \left(\int_0^T \sigma_\infty^2(s) \|X(s) - P_{\mathcal{S}}(X(s))\|^2 ds \right) \\ &= \int_0^T \sigma_\infty^2(s) \mathbb{E}(\text{dist}(X(t), \mathcal{S})^2) \\ &\leq C^* \int_0^{+\infty} \sigma_\infty^2(s) < \infty. \end{aligned}$$

Letting $Y(s) = \sigma^\top(s, X(s))(X(s) - P_{\mathcal{S}}(X(s)))$, then

$$\mathbb{E} \left[\int_{\hat{t}_\delta}^t \langle Y(s), dW(s) \rangle \right] = 0.$$

This immediately implies

$$\mathbb{E} \left[\mathbb{1}_{\Omega_\delta} \int_{\hat{t}_\delta}^t \langle Y(s), dW(s) \rangle \right] = -\mathbb{E} \left[\mathbb{1}_{\Omega_{\text{conv}} \setminus \Omega_\delta} \int_{\hat{t}_\delta}^t \langle Y(s), dW(s) \rangle \right].$$

The right hand side can be bounded using Cauchy-Schwarz inequality as follows

$$\begin{aligned} &\left| \mathbb{E} \left[\mathbb{1}_{\Omega_{\text{conv}} \setminus \Omega_\delta} \int_{\hat{t}_\delta}^t \langle Y(s), dW(s) \rangle \right] \right| \\ &= \left| \mathbb{E} \left[\mathbb{1}_{\Omega_{\text{conv}} \setminus \Omega_\delta} \int_{\hat{t}_\delta}^t \langle \sigma^\top(s, X(s))(X(s) - P_{\mathcal{S}}(X(s))), dW(s) \rangle \right] \right| \\ &\leq \sqrt{\mathbb{E}(\mathbb{1}_{\Omega_{\text{conv}} \setminus \Omega_\delta})} \sqrt{\mathbb{E} \left[\left(\int_{\hat{t}_\delta}^t \langle \sigma^\top(s, X(s))(X(s) - P_{\mathcal{S}}(X(s))), dW(s) \rangle \right)^2 \right]} \\ &\leq \sqrt{\delta} \sqrt{\mathbb{E} \left[\int_{\hat{t}_\delta}^t \|\sigma^\top(s, X(s))(X(s) - P_{\mathcal{S}}(X(s)))\|^2 ds \right]} \\ &\leq \sqrt{C^* \delta} \sqrt{\int_{\hat{t}_\delta}^t \sigma_\infty^2(s) ds} = \sqrt{C^* \delta} \int_{\hat{t}_\delta}^t \frac{\sigma_\infty^2(s)}{2\sqrt{\int_{\hat{t}_\delta}^s \sigma_\infty^2(u) du}} ds. \end{aligned}$$

Set $C_* = \sqrt{C^*}$, and recall that $C_\infty = \sqrt{\int_0^\infty \sigma_\infty^2(s) ds}$. Thus, for every $t > \hat{t}_\delta$

$$\hat{G}(t) \leq \hat{G}(\hat{t}_\delta) - \int_{\hat{t}_\delta}^t \mathbb{E}[(f(X(s)) - \min f) \mathbb{1}_{\Omega_\delta}] ds + \int_{\hat{t}_\delta}^t \sigma_\infty^2(s) ds + C_* \sqrt{\delta} \int_{\hat{t}_\delta}^t \frac{\sigma_\infty^2(s)}{2\sqrt{\int_{\hat{t}_\delta}^s \sigma_\infty^2(u) du}} ds. \quad (4.13)$$

Recall $h_\delta(t)$ from (4.3). Then, we can rewrite (4.13) as

$$\hat{G}(t) \leq \hat{G}(\hat{t}_\delta) - \int_{\hat{t}_\delta}^t \mathbb{E}[(f(X(s)) - \min f) \mathbb{1}_{\Omega_\delta}] ds + \int_{\hat{t}_\delta}^t h_\delta(s) ds, \quad \forall t > \hat{t}_\delta. \quad (4.14)$$

Using that $f \in \text{EB}^2([f \leq r])$, we obtain

$$\hat{G}(t) \leq \hat{G}(\hat{t}_\delta) - 2\gamma \int_{\hat{t}_\delta}^t \hat{G}(s) ds + \int_{\hat{t}_\delta}^t h_\delta(s) ds, \quad \forall t > \hat{t}_\delta.$$

Observe that $h_\delta \in L^1([\hat{t}_\delta, \infty[)$ since

$$\int_{\hat{t}_\delta}^{\infty} h_\delta(s) ds \leq C_\infty^2 + C_* C_\infty \sqrt{\delta}.$$

The goal now is to apply the comparison lemma to $\hat{G}(t)$ (see Lemma A.2) which necessitates to solve the following ODE

$$\begin{cases} y'(t) = -2\gamma y(t) + h_\delta(t) & t > \hat{t}_\delta \\ y(\hat{t}_\delta) = \hat{G}(\hat{t}_\delta). \end{cases}$$

Let $\lambda \in]0, 1[$. Using the integrating factor method, we obtain

$$\begin{aligned} y(t) &= e^{-2\gamma(t-\hat{t}_\delta)} y(\hat{t}_\delta) + e^{-2\gamma t} \int_{\hat{t}_\delta}^{\hat{t}_\delta + \lambda(t-\hat{t}_\delta)} h_\delta(s) e^{2\gamma s} ds + e^{-2\gamma t} \int_{\hat{t}_\delta + \lambda(t-\hat{t}_\delta)}^t h_\delta(s) e^{2\gamma s} ds \\ &\leq e^{-2\gamma(t-\hat{t}_\delta)} \mathbb{E}(\hat{g}(\hat{t}_\delta)) + e^{-2\gamma(1-\lambda)(t-\hat{t}_\delta)} \int_{\hat{t}_\delta}^{\hat{t}_\delta + \lambda(t-\hat{t}_\delta)} h_\delta(s) ds \\ &\quad + h_\delta(\hat{t}_\delta + \lambda(t-\hat{t}_\delta)) e^{-2\gamma t} \int_{\hat{t}_\delta + \lambda(t-\hat{t}_\delta)}^t e^{2\gamma s} ds \\ &\leq e^{-2\gamma(t-\hat{t}_\delta)} \mathbb{E}(\hat{g}(\hat{t}_\delta)) + e^{-2\gamma(1-\lambda)(t-\hat{t}_\delta)} (C_\infty^2 + C_* C_\infty \sqrt{\delta}) + \frac{h_\delta(\hat{t}_\delta + \lambda(t-\hat{t}_\delta))}{2\gamma}. \end{aligned}$$

where in the first inequality, we used that σ^2 is decreasing and so is h_δ . Lemma A.2 then gives

$$\begin{aligned} \mathbb{E} \left(\frac{\text{dist}(X(t), \mathcal{S})^2}{2} \mathbb{1}_{\Omega_\delta} \right) &\leq e^{-2\gamma(t-\hat{t}_\delta)} \mathbb{E} \left(\frac{\text{dist}(X(\hat{t}_\delta), \mathcal{S})^2}{2} \right) + e^{-2\gamma(1-\lambda)(t-\hat{t}_\delta)} (C_\infty^2 + C_* C_\infty \sqrt{\delta}) \\ &\quad + \frac{h_\delta(\hat{t}_\delta + \lambda(t-\hat{t}_\delta))}{2\gamma}. \end{aligned}$$

According to Corollary A.6 we obtain that for all $t > \hat{t}_\delta$

$$\begin{aligned} \mathbb{E} \left(\frac{\text{dist}(X(t), \mathcal{S})^2}{2} \right) &\leq e^{-2\gamma(t-\hat{t}_\delta)} \mathbb{E} \left(\frac{\text{dist}(X(\hat{t}_\delta), \mathcal{S})^2}{2} \right) + e^{-2\gamma(1-\lambda)(t-\hat{t}_\delta)} (C_\infty^2 + C_* C_\infty \sqrt{\delta}) \\ &\quad + \frac{h_\delta(\hat{t}_\delta + \lambda(t-\hat{t}_\delta))}{2\gamma} + C_d \sqrt{\delta}. \end{aligned}$$

(b) Denote $\tilde{g}(t) = \tilde{\phi}(X(t)) = f(X(t)) - \min f$ and $\tilde{G}(t) = \mathbb{E}(\mathbb{1}_{\Omega_\delta} \tilde{g}(t))$. By Proposition 2.4

$$\begin{aligned} \tilde{g}(t) &\leq \tilde{g}(\hat{t}_\delta) - \int_{\hat{t}_\delta}^t \left\langle \nabla f(X(s)), \nabla \tilde{\phi}(X(s)) \right\rangle ds + \frac{L}{2} \int_{\hat{t}_\delta}^t \text{tr}[\Sigma(s, X(s))] ds \\ &\quad + \mathbb{1}_{\Omega_\delta} \int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s)) \nabla f(X(s)), dW(s) \right\rangle. \quad (4.15) \end{aligned}$$

Multiplying both sides by $\mathbb{1}_{\Omega_\delta}$ and taking expectation we obtain

$$\begin{aligned} \tilde{G}(t) - \tilde{G}(\hat{t}_\delta) &\leq -\mathbb{E} \left[\int_{\hat{t}_\delta}^t \|\nabla f(X(s))\|^2 \mathbb{1}_{\Omega_\delta} ds \right] + \frac{L}{2} \mathbb{E} \left[\int_{\hat{t}_\delta}^t \text{tr}[\Sigma(s, X(s))] ds \right] \\ &\quad + \mathbb{E} \left[\mathbb{1}_{\Omega_\delta} \int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s)) \nabla f(X(s)), dW(s) \right\rangle \right]. \end{aligned} \quad (4.16)$$

On the other hand, we have

$$\begin{aligned} \mathbb{E} \left(\int_0^T \|\sigma^\top(s, X(s)) \nabla f(X(s))\|^2 ds \right) &\leq L^2 \mathbb{E} \left(\int_0^T \sigma_\infty^2(s) \|X(s) - x^*\|^2 ds \right) \\ &\leq L^2 C^* \int_0^{+\infty} \sigma_\infty^2(s) < \infty, \quad \forall T > 0. \end{aligned}$$

Since $\mathbb{E} \left[\int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s)) \nabla f(X(s)), dW(s) \right\rangle \right] = 0$, we have

$$\mathbb{E} \left[\mathbb{1}_{\Omega_\delta} \int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s)) (\nabla f(X(s))), dW(s) \right\rangle \right] = -\mathbb{E} \left[\mathbb{1}_{\Omega_{\text{conv}} \setminus \Omega_\delta} \int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s)) (\nabla f(X(s))), dW(s) \right\rangle \right].$$

The last term can be bounded as

$$\begin{aligned} &\left| \mathbb{E} \left[\mathbb{1}_{\Omega_{\text{conv}} \setminus \Omega_\delta} \int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s)) (\nabla f(X(s))), dW(s) \right\rangle \right] \right| \\ &\leq \sqrt{\mathbb{E}(\mathbb{1}_{\Omega_{\text{conv}} \setminus \Omega_\delta})} \sqrt{\mathbb{E} \left[\left(\int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s)) (\nabla f(X(s))), dW(s) \right\rangle \right)^2 \right]} \\ &\leq L\sqrt{\delta} \sqrt{\mathbb{E} \left[\int_{\hat{t}_\delta}^t \sigma_\infty^2(s) \|X(s) - x^*\|^2 ds \right]} \\ &\leq L\sqrt{C^*} \sqrt{\delta} \sqrt{\int_{\hat{t}_\delta}^t \sigma_\infty^2(s) ds} = L\sqrt{C^*} \sqrt{\delta} \int_{\hat{t}_\delta}^t \frac{\sigma_\infty^2(s)}{2\sqrt{\int_{\hat{t}_\delta}^s \sigma_\infty^2(u) du}} ds. \end{aligned}$$

Let us notice that if (3.14) holds, then Proposition 3.3 tells us that $\mathbb{E}(f(X(t)) - \min f) \leq K't^{\beta-1}$ with $\beta \in [0, 1]$, and for some $K' > 0$. In this case

$$\left| \mathbb{E} \left[\mathbb{1}_{\Omega_{\text{conv}} \setminus \Omega_\delta} \int_{\hat{t}_\delta}^t \left\langle \sigma^\top(s, X(s)) (\nabla f(X(s))), dW(s) \right\rangle \right] \right| \leq \sqrt{2LK'} \sqrt{\delta} \int_{\hat{t}_\delta}^t \frac{\sigma_\infty^2(s) s^{\beta-1}}{2\sqrt{\int_{\hat{t}_\delta}^s \sigma_\infty^2(u) u^{\beta-1} du}} ds.$$

Injecting this into (4.16), we have for all $t > \hat{t}_\delta$

$$\begin{aligned} \tilde{G}(t) &\leq \tilde{G}(\hat{t}_\delta) - \mathbb{E} \left[\int_{\hat{t}_\delta}^t \|\nabla f(X(s))\|^2 \mathbb{1}_{\Omega_\delta} ds \right] + \frac{L}{2} \int_{\hat{t}_\delta}^t \sigma_\infty^2(s) ds \\ &\quad + \begin{cases} C_K \sqrt{\delta} \int_{\hat{t}_\delta}^t \frac{\sigma_\infty^2(s) s^{\beta-1}}{2\sqrt{\int_{\hat{t}_\delta}^s \sigma_\infty^2(u) u^{\beta-1} du}} ds, & \forall t > \hat{t}_\delta \text{ if (3.14) holds,} \\ C_{**} \sqrt{\delta} \int_{\hat{t}_\delta}^t \frac{\sigma_\infty^2(s)}{2\sqrt{\int_{\hat{t}_\delta}^s \sigma_\infty^2(u) du}} ds & \text{otherwise,} \end{cases} \end{aligned} \quad (4.17)$$

where $C_{**} = L\sqrt{C^*}$, $C_K = \sqrt{2LK'}$ and recall that $C_\infty = \sqrt{\int_0^\infty \sigma_\infty^2(s) ds}$. Recalling $l_\delta(t)$ and $k_\delta(t)$ from (4.4)-(4.5), and by Fubini's theorem, (4.17) becomes

$$\tilde{G}(t) \leq \tilde{G}(\hat{t}_\delta) - \int_{\hat{t}_\delta}^t \mathbb{E} [\|\nabla f(X(s))\|^2 \mathbb{1}_{\Omega_\delta}] ds + \begin{cases} \int_{\hat{t}_\delta}^t k_\delta(s) ds & \text{if (3.14) holds,} \\ \int_{\hat{t}_\delta}^t l_\delta(s) ds & \text{otherwise.} \end{cases} \quad (4.18)$$

Since $f \in \mathbb{L}^{1/2}(\mathcal{S})$, there exists $\mu > 0$ such that

$$\tilde{G}(t) \leq \tilde{G}(\hat{t}_\delta) - \mu^2 \int_{\hat{t}_\delta}^t \tilde{G}(s) ds + \begin{cases} \int_{\hat{t}_\delta}^t k_\delta(s) ds & \text{if (3.14) holds,} \\ \int_{\hat{t}_\delta}^t l_\delta(s) ds & \text{otherwise.} \end{cases} \quad (4.19)$$

To get an explicit bound in (4.19), we use Lemma A.2, which involves solving

$$(E.2) \quad \begin{cases} y'(t) &= -\mu^2 y(t) + l_\delta(t), & t > \hat{t}_\delta \\ y(\hat{t}_\delta) &= \tilde{G}(\hat{t}_\delta) \end{cases}$$

$$(E.3) \quad \begin{cases} y'(t) &= -\mu^2 y(t) + k_\delta(t), & t > \hat{t}_\delta \\ y(\hat{t}_\delta) &= \tilde{G}(\hat{t}_\delta) \end{cases}$$

Let $\lambda \in]0, 1[$. Using the integrating factor method as in (i), we get for (E.2)

$$y(t) \leq e^{-\mu^2(t-\hat{t}_\delta)} \mathbb{E}(\tilde{g}(\hat{t}_\delta)) + \begin{cases} e^{-\mu^2(1-\lambda)(t-\hat{t}_\delta)} \left(\frac{LC_\infty^2}{2} + C_{**}C_\infty\sqrt{\delta} \right) + \frac{l_\delta(\hat{t}_\delta + \lambda(t-\hat{t}_\delta))}{\mu^2} & \text{for (E.2)} \\ e^{-\mu^2(1-\lambda)(t-\hat{t}_\delta)} \left(\frac{LC_\infty^2}{2} + C_KC_\infty\sqrt{\hat{t}_\delta^{\beta-1}\sqrt{\delta}} \right) + \frac{k_\delta(\hat{t}_\delta + \lambda(t-\hat{t}_\delta))}{\mu^2} & \text{for (E.3).} \end{cases}$$

Using Lemma A.2 and Corollary A.6

$$\begin{aligned} \mathbb{E}[f(X(t)) - \min f] &\leq y(t) + C_f\sqrt{\delta} \\ &\leq e^{-\mu^2(t-\hat{t}_\delta)} \mathbb{E}[f(X(\hat{t}_\delta)) - \min f] + C_f\sqrt{\delta} \\ &\quad + \begin{cases} e^{-\mu^2(1-\lambda)(t-\hat{t}_\delta)} \left(\frac{LC_\infty^2}{2} + C_{**}C_\infty\sqrt{\delta} \right) + \frac{l_\delta(\hat{t}_\delta + \lambda(t-\hat{t}_\delta))}{\mu^2} & \text{for (E.2)} \\ e^{-\mu^2(1-\lambda)(t-\hat{t}_\delta)} \left(\frac{LC_\infty^2}{2} + C_KC_\infty\sqrt{\hat{t}_\delta^{\beta-1}\sqrt{\delta}} \right) + \frac{k_\delta(\hat{t}_\delta + \lambda(t-\hat{t}_\delta))}{\mu^2} & \text{for (E.3).} \end{cases} \end{aligned}$$

(ii) $p > 2$:

(a) We embark from inequality (4.14) and we now use that $f \in \text{EB}^p([f \leq r])$ with $p > 2$, to get

$$\begin{aligned} \hat{G}(t) &\leq \hat{G}(\hat{t}_\delta) - \int_{\hat{t}_\delta}^t \mathbb{E} [(f(X(s)) - \min f) \mathbb{1}_{\Omega_\delta}] ds + \int_{\hat{t}_\delta}^t h_\delta(s) ds \\ &\leq \hat{G}(\hat{t}_\delta) - 2^{p/2}\gamma \int_{\hat{t}_\delta}^t \hat{G}(s)^{p/2} + \int_{\hat{t}_\delta}^t h_\delta(s) ds. \end{aligned} \quad (4.20)$$

In the last inequality, we used that $p > 2$ and Jensen's inequality.

The idea is again to use the comparison lemma (Lemma A.2), which will now involve solving the Cauchy problem (C.1), and finally invoke Corollary A.6.

(b) The reasoning is similar to the previous point using now that $f \in \mathbb{L}^q(\mathcal{S})$ and the computations of (i)(b). We omit the details for the sake of brevity. \square

5 SDE for nonsmooth structured convex optimization

In this section, we turn to the composite convex minimization problem with additive structure

$$\min_{x \in \mathbb{R}^d} f(x) + g(x), \quad (5.1)$$

where

$$\begin{cases} f \in C_L^{1,1}(\mathbb{R}^d) \cap \Gamma_0(\mathbb{R}^d) \text{ and } g \in \Gamma_0(\mathbb{R}^d); \\ \mathcal{S} = \operatorname{argmin}(f + g) \neq \emptyset. \end{cases} \quad (\text{H}'_0)$$

The importance of this class of problems comes from its wide spectrum of applications ranging from data processing, to machine learning and statistics to name a few.

We consider two different approaches leading to different SDE's. The first is based on a fixed point argument and the use of the notion of cocoercive monotone operator. The second approach is based on a regularization/smoothing argument, for instance the Moreau envelope.

5.1 Fixed point approach via cocoercive monotone operators

Let us start with some classical definitions concerning monotone operators.

Definition 5.1. An operator $A : \mathbb{R}^d \rightarrow \mathcal{P}(\mathbb{R}^d)$ is monotone if

$$\langle u - v, x - y \rangle \geq 0, \quad \forall (x, u) \in \operatorname{graph}(A), (y, v) \in \operatorname{graph}(A).$$

It is maximally monotone if there exists no monotone operator whose graph properly contains $\operatorname{graph}(A)$. Moreover, A is γ -strongly monotone with modulus $\gamma > 0$ if

$$\langle u - v, x - y \rangle \geq \gamma \|x - y\|^2, \quad \forall (x, u) \in \operatorname{graph}(A), (y, v) \in \operatorname{graph}(A).$$

Remark 5.2. If A is maximally monotone and strongly monotone, then $A^{-1}(0) \stackrel{\text{def}}{=} \{x \in \mathbb{R}^d : A(x) = 0\}$ is non-empty and reduced to a singleton.

Remark 5.3. The subdifferential operator ∂g of $g \in \Gamma_0(\mathbb{R}^d)$ is maximally monotone.

Definition 5.4. A single-valued operator $M : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is cocoercive with constant $\rho > 0$ if

$$\langle M(x) - M(y), x - y \rangle \geq \rho \|M(x) - M(y)\|^2, \quad \forall x, y \in \mathbb{R}^d.$$

Remark 5.5. It is clear that a cocoercive operator is ρ^{-1} -Lipschitz continuous. In turn, a cocoercive operator is maximally monotone.

Remark 5.6. If $f \in C_L^{1,1}(\mathbb{R}^d) \cap \Gamma_0(\mathbb{R}^d)$, then the operator ∇f is L^{-1} -cocoercive.

Our interest now is to solve the structured monotone inclusion problem

$$0 \in A(x) + B(x),$$

where A is maximally monotone, and B is cocoercive with $(A + B)^{-1}(0) \neq \emptyset$. This is of course a generalization of (5.1) by taking $A = \partial g$ and $B = \nabla f$.

A favorable situation occurs when one can compute the resolvent operator of A

$$J_{\mu A} = (I + \mu A)^{-1}, \quad \mu > 0.$$

In this case, we can develop a strategy parallel to the one which consists in replacing a maximally monotone operator by its Yosida approximation. Indeed, given $\mu > 0$, we have

$$(A + B)(x) \ni 0 \iff x - J_{\mu A}(x - \mu B(x)) = 0 \iff M_{A,B,\mu}(x) = 0, \quad (5.2)$$

where $M_{A,B,\mu} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the single-valued operator defined by

$$M_{A,B,\mu}(x) = \frac{1}{\mu} (x - J_{\mu A}(x - \mu B(x))). \quad (5.3)$$

$M_{A,B,\mu}$ is closely tied to the well-known forward-backward fixed point operator. Moreover, when $B = 0$, $M_{A,B,\mu} = \frac{1}{\mu} (I - J_{\mu A})$ which is nothing but the Yosida regularization of A with index μ . As a remarkable property, for the μ parameter properly set, the operator $M_{A,B,\mu}$ is cocoercive. This is made precise in the following result.

Proposition 5.7. [5, Lemma B.1] *Let $A : \mathbb{R}^d \rightarrow \mathcal{P}(\mathbb{R}^d)$ be a general maximally monotone operator, and let $B : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a monotone operator which is λ -cocoercive. Assume that $\mu \in]0, 2\lambda[$. Then, $M_{A,B,\mu}$ is ρ -cocoercive with*

$$\rho = \mu \left(1 - \frac{\mu}{4\lambda} \right).$$

We first focus on finding the zeros of M , where

$$M : \mathbb{R}^d \rightarrow \mathbb{R}^d \text{ is cocoercive and } M^{-1}(0) \neq \emptyset. \quad (\text{H}_0^M)$$

We will then specialize our results to the case of a structured operator of the form $M_{A,B,\mu}$.

Our goal is to handle the situation where M can be evaluated up to a stochastic error. We therefore consider the following SDE, defined for (deterministic) initial data $X_0 \in \mathbb{R}^d$,

$$\begin{cases} dX(t) = -M(X(t))dt + \sigma(t, X(t))dW(t), & t \geq 0 \\ X(0) = X_0. \end{cases} \quad (\text{SDE}^M)$$

As in Section 1.1, we will assume that W is a \mathcal{F}_t -adapted m -dimensional Brownian motion, and the volatility matrix $\sigma : \mathbb{R}_+ \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$ satisfies (H).

Let us now state the natural extensions of our main results to this situation.

Theorem 5.8. *Let $M : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a cocoercive operator. Consider the stochastic differential equation (SDE^M) under the hypotheses (H₀^M) and (H). Then, there exists a unique solution $X \in S_d^\nu$, for every $\nu \geq 2$. Moreover, if $\sigma_\infty \in L^2(\mathbb{R}_+)$, then:*

- (i) $\sup_{t \geq 0} \mathbb{E}[\|X(t)\|^2] < \infty$.
- (ii) $\forall x^* \in M^{-1}(0)$, $\lim_{t \rightarrow \infty} \|X(t) - x^*\|$ exists a.s. and $\sup_{t \geq 0} \|X(t)\| < \infty$ a.s.
- (iii) $\lim_{t \rightarrow \infty} \|M(X(t))\| = 0$ a.s.
- (iv) There exists an $M^{-1}(0)$ -valued random variable x^* such that $\lim_{t \rightarrow \infty} X(t) = x^*$ a.s.

Proof. Existence and uniqueness follow from Theorem A.7 since M is Lipschitz continuous and σ verifies (H). The proof of the first two items remains the same as for Theorem 3.1. For the third item, we use the cocoercivity of M instead of the convexity of f and Corollary 2.2 to prove that $\lim_{t \rightarrow \infty} \|M(X(t))\| = 0$ a.s. For the last item, it suffices to use that the operator M is continuous (since it is Lipschitz continuous) to conclude with Opial's Lemma. \square

Theorem 5.9. Let $M : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a ρ -cocoercive operator. Let us make the assumptions (\mathbf{H}_0^M) and (H). Let $X \in S_d^2$ be the solution of (SDE^M) with initial condition X_0 . Then the following properties are satisfied:

(i) Let $\overline{M \circ X}(t) \stackrel{\text{def}}{=} t^{-1} \int_0^t M(X(s)) ds$ and $\overline{\|M(X(t))\|^2} \stackrel{\text{def}}{=} t^{-1} \int_0^t \|M(X(s))\|^2 ds$. We have

$$\mathbb{E} [\|\overline{M \circ X}(t)\|^2] \leq \mathbb{E} [\overline{\|M(X(t))\|^2}] \leq \frac{\text{dist}(X_0, M^{-1}(0))^2}{2\rho t} + \frac{\sigma_*^2}{2\rho}, \quad \forall t > 0. \quad (5.4)$$

Besides, if σ_∞ is $L^2(\mathbb{R}_+)$, then

$$\mathbb{E} [\|\overline{M \circ X}(t)\|^2] \leq \mathbb{E} [\overline{\|M(X(t))\|^2}] = \mathcal{O}\left(\frac{1}{t}\right), \quad \forall t > 0. \quad (5.5)$$

(ii) If M is γ -strongly monotone, then $M^{-1}(0) = \{x^*\}$ and

$$\mathbb{E} \left(\frac{\|X(t) - x^*\|^2}{2} \right) \leq \frac{\|X_0 - x^*\|^2}{2} e^{-2\gamma t} + \frac{\sigma_*^2}{4\gamma}, \quad \forall t \geq 0. \quad (5.6)$$

If, moreover, σ_∞ is decreasing and vanishes at infinity, then for every $\lambda \in]0, 1[$

$$\mathbb{E} \left(\frac{\|X(t) - x^*\|^2}{2} \right) \leq \frac{\|X_0 - x^*\|^2}{2} e^{-2\gamma t} + \frac{\sigma_*^2}{4} e^{-2\gamma t(1-\lambda)} + \frac{\sigma_\infty^2(\lambda t)}{2}, \quad \forall t > 0. \quad (5.7)$$

Proof. Analogous to Theorem 3.2. \square

We now turn to the local convergence properties. To this end, we need an extension of the Hölderian error bound inequality (or Łojasiewicz inequality) to the operator setting. For convex functions, it is known that error bound inequalities are closely related to metric subregularity of the subdifferential [41, 42, 43]. This leads to the following definition.

Definition 5.10. Let $M : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a single-valued operator. We say that M satisfies the Hölder metric subregularity property with exponent $p \geq 2$ at $x^* \in M^{-1}(0)$ if there exists $\gamma > 0$ and a neighbourhood \mathcal{V}_{x^*} such that

$$\|M(x)\|^2 \geq \gamma \text{dist}(x, M^{-1}(0))^p, \quad \forall x \in \mathcal{V}_{x^*}. \quad (5.8)$$

If this inequality holds for any $x^* \in M^{-1}(0)$ with the same γ , we will write $M \in \text{HMS}^p(\mathbb{R}^d)$.

Theorem 5.11. Let M be a ρ -cocoercive operator such that $M \in \text{HMS}^2(\mathbb{R}^d)$. Let $X \in S_d^2$ be the solution of (SDE^M) under the hypotheses (\mathbf{H}_0^M) , (H). Suppose that $\sigma_\infty \in L^2(\mathbb{R}_+)$ ($C_\infty \stackrel{\text{def}}{=} \|\sigma_\infty\|_{L^2(\mathbb{R}_+)}$) and σ_∞ is decreasing. Consider also the positive constants C, C_d, γ . Then, for all $\delta > 0$, there exists $\hat{t}_\delta > 0$ such that for every $\lambda \in (0, 1)$:

$$\begin{aligned} \mathbb{E} \left(\frac{\text{dist}(X(t), M^{-1}(0))^2}{2} \right) &\leq e^{-2\gamma\rho(t-\hat{t}_\delta)} \mathbb{E} \left(\frac{\text{dist}(X(\hat{t}_\delta), M^{-1}(0))^2}{2} \right) \\ &\quad + e^{-2\gamma\rho(1-\lambda)(t-\hat{t}_\delta)} (C_\infty^2 + C_\infty C\sqrt{\delta}) \\ &\quad + \frac{h_\delta(\hat{t}_\delta + \lambda(t - \hat{t}_\delta))}{2\gamma\rho} + C_d\sqrt{\delta}, \quad \forall t > \hat{t}_\delta, \end{aligned} \quad (5.9)$$

where $h_\delta(t) = \sigma_\infty^2(t) + C\sqrt{\delta} \frac{\sigma_\infty^2(t)}{2\sqrt{\int_{t_\delta}^t \sigma_\infty^2(u)du}}$.

Proof. The proof is essentially the same as that of Theorem 4.5(i)(a), where instead of convexity in (4.12), we use cocoercivity of M , and in (4.14) we invoke Theorem 5.8 and Hölder metric subregularity. \square

Remark 5.12. We can naturally extend the previous result for $p > 2$ as in Theorem 4.5(ii). Nevertheless, since that bound is not explicit, we will skip this extension.

As an immediate consequence of the above result, by considering the cocoercive operator $M_{A,B,\mu}$ defined in (5.3), we obtain the following result.

Corollary 5.13. *Let $A : \mathbb{R}^d \rightarrow \mathcal{P}(\mathbb{R}^d)$ be a maximally monotone operator and $B : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a λ -cocoercive operator, $\lambda > 0$. Let $M_{A,B,\mu}$ be the operator defined in (5.3). Assume that $\mu \in]0, 2\lambda[$ and $(A + B)^{-1}(0) \neq \emptyset$. Then, the operator $M_{A,B,\mu}$ is ρ -cocoercive with $\rho = \mu \left(1 - \frac{\mu}{4\lambda}\right)$, and the SDE:*

$$\begin{cases} dX(t) = -M_{A,B,\mu}(X(t))dt + \sigma(t, X(t))dW(t), & t \geq 0 \\ X(0) = X_0, \end{cases}$$

has a unique solution $X \in S_d^\nu$, for every $\nu \geq 2$, that verifies the conclusions of Theorem 5.8 and Theorem 5.9. In particular, if $\sigma_\infty \in L^2(\mathbb{R}_+)$, there exists an $(A + B)^{-1}(0)$ -valued random variable x^* such that $\lim_{t \rightarrow \infty} X(t) = x^*$ a.s.

This result naturally applies to problem (5.1) when $\mathcal{S} = \operatorname{argmin}(f + g) \neq \emptyset$ by taking $A = \partial g$ and $B = \nabla f$. In this case, one has that $X(t)$ converges a.s. to an \mathcal{S} -valued random variable. Moreover, using standard inequalities, see e.g. [44], one can show that

$$\mathbb{E} \left[(f + g) \left(t^{-1} \int_0^t (\operatorname{prox}_{\mu g}(x - \mu \nabla f(x))) ds \right) - \min(f + g) \right] = \mathcal{O} \left(\sqrt{\mathbb{E} \left[\|M(X(t))\|^2 \right]} \right),$$

where $\operatorname{prox}_{\mu g} = (I + \mu \nabla g)^{-1}$ is the proximal mapping of g . From this, one can deduce an $\mathcal{O}(t^{-1/2})$ rate thanks to (5.4) and (5.5).

5.2 Approach via Moreau-Yosida regularization

The previous approach, though it is able to deal with more general setting (that of monotone inclusions), took us out of the framework of convex optimization by considering instead a dynamic governed by a cocoercive operator. In particular, the perturbation/noise is considered on the whole operator evaluation and not on a part of it (i.e. B) as it is standard in many applications. Moreover this approach led to a pessimistic convergence rate estimate when specialized to convex function minimization. By contrast, the following approach will operate directly on problem (5.1) and is based on a standard smoothing approach, replacing the non-smooth part g by its Moreau envelope [45].

5.2.1 Moreau envelope

Let us start by recalling some basic facts concerning the Moreau envelope.

Definition 5.14. Let $g \in \Gamma_0(\mathbb{R}^d)$. Given $\theta > 0$, the Moreau envelope of g of parameter θ is the function

$$g_\theta(x) \stackrel{\text{def}}{=} \inf_{y \in \mathbb{R}^d} \left(g(y) + \frac{1}{2\theta} \|x - y\|^2 \right) = \left(g \square \frac{1}{\theta} q \right) (x)$$

where \square is the infimal convolution operator and $q(x) = \frac{1}{2} \|x\|^2$.

The Moreau envelope has remarkable approximation and regularization properties, as summarized in the following statement.

Proposition 5.15. Let $g \in \Gamma_0(\mathbb{R}^d)$.

- (i) $g_\theta(x) \downarrow \inf g(\mathbb{R}^d)$ as $\theta \uparrow +\infty$.
- (ii) $g_\theta(x) \uparrow g(x)$ as $\theta \downarrow 0$.
- (iii) $g_\theta(x) \leq g(x)$ for any $\theta > 0$ and $x \in \mathbb{R}^d$,
- (iv) $\text{argmin}(g_\theta) = \text{argmin}(g)$ for any $\theta > 0$,
- (v) $g(x) - g_\theta(x) \leq \frac{\theta}{2} \|\partial^0 g(x)\|^2$ for any $\theta > 0$ and $x \in \text{dom}(\partial g)$,
- (vi) $g_\theta \in C_{\frac{1}{\theta}}^{1,1}(\mathbb{R}^d) \cap \Gamma_0(\mathbb{R}^d)$ for any $\theta > 0$.

We use the following notation in the rest of the section: $F \stackrel{\text{def}}{=} f + g$, $\mathcal{S} \stackrel{\text{def}}{=} \text{argmin} F$, $F_\theta \stackrel{\text{def}}{=} f + g_\theta$ and $\mathcal{S}_\theta \stackrel{\text{def}}{=} \text{argmin} F_\theta$.

Note that $F_\theta \in C_{L+\frac{1}{\theta}}^{1,1}(\mathbb{R}^d) \cap \Gamma_0(\mathbb{R}^d)$. Thus we will use F_θ as the potential driving (SDE), that is

$$\begin{cases} dX(t) = -\nabla F_\theta(X(t))dt + \sigma(t, X(t))dW(t), & t \geq 0 \\ X(0) = X_0. \end{cases} \quad (\text{SDE}_\theta)$$

Under (\mathbf{H}'_0) and (\mathbf{H}) , we will show almost sure convergence of the trajectory and corresponding convergence rates.

Remark 5.16. Though we focus here on the Moreau envelope, our convergence results, in particular, Proposition 5.19, still hold with infimal-convolution based smoothing using more general smooth kernels beyond the norm squared; see [45, Section 4.4].

5.2.2 Convergence of the trajectory

Applying Theorem 3.1 to F_θ , we have the following result.

Proposition 5.17. For any $\theta > 0$, let $X_\theta \in \mathcal{S}_\theta^2$ be the solution of the dynamic (SDE $_\theta$) governed by the potential F_θ , and make assumptions (\mathbf{H}'_0) , $\mathcal{S}_\theta \neq \emptyset$, (\mathbf{H}) and $\sigma_\infty \in \mathbb{L}^2(\mathbb{R}_+)$. Then there exists an \mathcal{S}_θ -valued random variable x_θ^* such that

$$\lim_{t \rightarrow \infty} X_\theta(t) = x_\theta^*, \quad a.s.$$

If $f = 0$, then $\mathcal{S}_\theta = \mathcal{S}$ (see Proposition 5.15(iv)), and Proposition 5.17 provides almost sure convergence to a solution of (5.1). On the other hand for $f \neq 0$, $\mathcal{S} \neq \mathcal{S}_\theta$ in general and we only obtain an "approximate" solution of (5.1); see Proposition 5.18(ii) for a quantitative estimate of this approximation when f is strongly convex. To obtain a true solution of the initial problem, a common device consists in using a diagonalization process which combines the dynamic with the approximation. Specifically, one considers

$$\begin{cases} dX(t) = -\nabla F_{\theta(t)}(X(t))dt + \sigma(t, X(t))dW(t), & t \geq 0 \\ X(0) = X_0, \end{cases} \quad (\text{SDE}_{\theta(t)})$$

where $\theta(t) \downarrow 0$ as $t \rightarrow +\infty$. In the deterministic case, an abundant literature has been devoted to the convergence of this type of systems. Note that unlike the cocoercive approach, we are now faced with a non-autonomous stochastic differential equation, making this a difficult problem, a subject for further research.

5.2.3 Convergence rates

We start with the following uniform bound on \mathcal{S}_θ which holds under slightly reinforced, but reasonable assumptions on f and g .

Proposition 5.18. *Consider f, g where f and g are proper lsc and convex, and g is also L_0 -Lipschitz continuous.*

(i) *Assume that $F = f + g$ is coercive. Then for any $\theta \geq 0$ there exists $C > 0$ (independent of θ) such that*

$$\sup_{z \in \mathcal{S}_\theta} \|z\| \leq C. \quad (5.10)$$

(ii) *Assume that $f \in \Gamma_\mu(\mathbb{R}^d)$ for $\mu > 0$, then (5.10) holds, $\mathcal{S} = \{x^*\}$, $\mathcal{S}_\theta = \{x_\theta^*\}$ and*

$$\|x_\theta^* - x^*\|^2 \leq \frac{L_0}{\mu} \theta. \quad (5.11)$$

Proof. (i) Since F is coercive, so is F_θ . Thus both \mathcal{S} and \mathcal{S}_θ are non-empty compact sets. Let $x_\theta^* \in \mathcal{S}_\theta$ and $x^* \in \mathcal{S}$. By Proposition 5.15(v) and Lipschitz continuity of g , we obtain

$$F(x_\theta^*) \leq F_\theta(x_\theta^*) + \frac{L_0^2}{2} \theta.$$

Moreover,

$$F_\theta(x_\theta^*) + \frac{L_0^2}{2} \theta \leq F_\theta(x^*) + \frac{L_0^2}{2} \theta \leq F(x^*) + \frac{L_0^2}{2} \theta \leq \min(F) + \frac{L_0^2}{2} \theta \stackrel{\text{def}}{=} \tilde{C},$$

where the second inequality is given by Proposition 5.15(iv). On the other hand, the coercivity of F implies that there exists $a > 0, b \in \mathbb{R}$ such that for any $x \in \mathbb{R}^d$

$$a\|x\| + b \leq F(x).$$

Therefore, collecting the above inequalities yields

$$a\|x_\theta^*\| + b \leq F(x_\theta^*) \leq \tilde{C}.$$

Taking the supremum over x_θ^* and defining $C \stackrel{\text{def}}{=} \frac{\tilde{C}-b}{a} \geq 0$, we obtain (5.10), or equivalently that the set of approximate minimizers is bounded independently of θ .

(ii) Since f is μ -strongly convex, so are F and F_θ . In turn, F is coercive and thus (5.10) holds by claim (i). Strong convexity implies uniqueness of minimizers of F and F_θ . Moreover,

$$\frac{\mu}{2} \|x_\theta^* - x^*\|^2 \leq F_\theta(x^*) - F_\theta(x_\theta^*). \quad (5.12)$$

From Proposition 5.15(iii)-(v) and Lipschitz continuity of g , we infer that

$$F_\theta(x^*) - F_\theta(x_\theta^*) \leq F(x^*) - F_\theta(x_\theta^*) \leq F(x_\theta^*) - F_\theta(x_\theta^*) = g(x_\theta^*) - g_\theta(x_\theta^*) \leq \frac{L_0}{2} \theta. \quad (5.13)$$

Combining (5.12) and (5.13), we get the claimed bound. □

We are now ready to establish complexity results.

Proposition 5.19. *Suppose that in addition to (\mathbf{H}'_0) and (\mathbf{H}) , $F = f + g$ is coercive and g is L_0 -Lipschitz continuous. Let X_θ be the solution of (\mathbf{SDE}_θ) governed by F_θ with $\theta > 0$. Let $C_0 = \|X_0\| + C$, where C is the constant (independent of θ), defined in (5.10). Then the following statements hold for any $t > 0$.*

(i) Let $\overline{X}_\theta(t) = t^{-1} \int_0^t X_\theta(s) ds$, then

$$\mathbb{E} (F(\overline{X}_\theta(t)) - \min F) \leq \frac{C_0^2}{2t} + \frac{\sigma_*^2}{2} + \theta \frac{L_0^2}{2}.$$

Besides, if $\sigma_\infty \in L^2(\mathbb{R}_+)$, then

$$\mathbb{E} (F(\overline{X}_\theta(t)) - \min F) = \frac{C_0^2 + \int_0^{+\infty} \sigma_\infty^2(s) ds}{2t} + \theta \frac{L_0^2}{2}.$$

(ii) If σ_∞ verifies (3.14) and $\theta \in]0, 1]$, then

$$\mathbb{E} (F(X(t)) - \min F) = \frac{C_0^2}{2t} + \frac{K(1+L)}{2\theta} t^{\beta-1} + \theta \frac{L_0^2}{2}.$$

(iii) If, in addition, $f \in \Gamma_\mu(\mathbb{R}^d)$ for some $\mu > 0$ then $\mathcal{S} = \{x^*\}$, $\mathcal{S}_\theta = \{x_\theta^*\}$, and

$$\mathbb{E} (\|X_\theta(t) - x^*\|^2) \leq 2C_0^2 e^{-2\mu t} + \frac{\sigma_*^2}{\mu} + 2 \frac{L_0}{\mu} \theta.$$

Besides, if σ_∞ is decreasing and vanishes at infinity, then $\forall \lambda \in]0, 1[$:

$$\mathbb{E} (\|X_\theta(t) - x^*\|^2) \leq 2C_0^2 e^{-2\mu t} + \frac{\sigma_*^2}{\mu} e^{-2\mu(1-\lambda)t} + 2\sigma_\infty^2(\lambda t) + 2 \frac{L_0}{\mu} \theta.$$

Remark 5.20. Observe that when $f = 0$, then $\mathcal{S}_\theta = \mathcal{S}$. Therefore in Proposition 5.19 we have $x_\theta^* = x^*$ and the last term in θ can be dropped.

Proof. (i) Combine Theorem 3.2(i) applied to F_θ , Proposition 5.15(iii) and (v), and Proposition 5.18(i) to see that $\text{dist}(X_0, \mathcal{S}_\theta) \leq C_0$.

(ii) Argue as in claim (i) using Proposition 3.3 instead of Theorem 3.2(i), and use the fact that ∇F_θ is Lipschitz continuous with constant

$$L + \frac{1}{\theta} \leq \frac{L+1}{\theta} \quad \text{for } \theta \in]0, 1].$$

(iii) Combine Theorem 3.2(ii) applied to F_θ , Proposition 5.18(ii) and Jensen's inequality. □

A Auxiliary results

A.1 Deterministic results

The following lemma is straightforward to prove. We omit the details.

Lemma A.1. Let $t_0 > 0$ and $g : [t_0, +\infty[\rightarrow \mathbb{R}_+$. Suppose that $\lim_{t \rightarrow \infty} g(t)$ exists and $\int_{t_0}^{\infty} \frac{g(s)}{s} ds < +\infty$. Then $\lim_{t \rightarrow \infty} g(t) = 0$.

The next result is an adaptation of [46, Proposition 2.3] to our specific context but under slightly less stringent assumptions.

Lemma A.2 (Comparison Lemma). Let $t_0 \geq 0$ and $T > t_0$. Assume that $h : [t_0, +\infty[\rightarrow \mathbb{R}_+$ is measurable with $h \in L^1([t_0, T])$, that $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is continuous and nondecreasing, $\varphi_0 > 0$ and the Cauchy problem

$$\begin{cases} \varphi'(t) = -\psi(\varphi(t)) + h(t) & \text{for almost all } t \in [t_0, T] \\ \varphi(t_0) = \varphi_0 \end{cases}$$

has an absolutely continuous solution $\varphi : [t_0, T] \rightarrow \mathbb{R}_+$. If a bounded from below lower semicontinuous function $\omega : [t_0, T] \rightarrow \mathbb{R}_+$ satisfies

$$\omega(t) \leq \omega(s) - \int_s^t \psi(\omega(\tau)) d\tau + \int_s^t h(\tau) d\tau$$

for $t_0 \leq s < t \leq T$ and $\omega(t_0) = \varphi_0$, then

$$\omega(t) \leq \varphi(t) \quad \text{for } t \in [t_0, T].$$

Theorem A.3 (Egorov's Theorem). [47, Chapter 3, Exercise 16] If $\mu(X) < \infty$ and $(f_t)_{t \in \mathbb{R}_+}$ is a family of real functions such that for all $x \in X$:

1. $\lim_{t \rightarrow \infty} f_t(x) = f(x)$ and
2. $t \mapsto f_t(x)$ is continuous.

Then, for every $\delta > 0$, there exists a measurable set $E_\delta \subset X$, with $\mu(X \setminus E_\delta) < \delta$, such that $(f_t)_{t \in \mathbb{R}_+}$ converges uniformly on E_δ .

Lemma A.4. Let $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ and $\liminf_{t \rightarrow \infty} f(t) \neq \limsup_{t \rightarrow \infty} f(t)$. Then there exists a constant α , satisfying $\liminf_{t \rightarrow \infty} f(t) < \alpha < \limsup_{t \rightarrow \infty} f(t)$, such that for every $\beta > 0$, we can define a sequence $(t_k)_{k \in \mathbb{N}} \subset \mathbb{R}$ such that

$$f(t_k) > \alpha, \quad t_{k+1} > t_k + \beta, \quad \forall k \in \mathbb{N}.$$

Proof. Since $\liminf_{t \rightarrow \infty} f(t)$ and $\limsup_{t \rightarrow \infty} f(t)$ are different real numbers, there exists α such that

$$\liminf_{t \rightarrow \infty} f(t) < \alpha < \limsup_{t \rightarrow \infty} f(t).$$

Moreover, by definition of \limsup , there exists a sequence $(t_k)_{k \in \mathbb{N}}$ such that $\lim_{k \rightarrow \infty} t_k = \infty$ and $f(t_k) > \alpha$. Let $\beta > 0$ and $n_0 = 0$, let us define recursively for $j \geq 1$, $n_j = \min\{n > n_{j-1} : t_n - t_{n_{j-1}} > \beta\}$. Let $j' \in \mathbb{N}$ be the first natural such that $n_{j'} = \infty$. This implies that for every $n > n_{j'-1}$, $t_n \leq \beta + t_{n_{j'-1}} < \infty$, a contradiction since $\lim_{n \rightarrow \infty} t_n = \infty$, then for every $j \in \mathbb{N}$, $n_j < \infty$. Thus, we can define $(t_{n_j})_{j \in \mathbb{N}}$ a subsequence of $(t_k)_{k \in \mathbb{N}}$ such that $\lim_{j \rightarrow \infty} t_{n_j} = \infty$ and for every $j \in \mathbb{N}$, $t_{n_{j+1}} - t_{n_j} > \beta$. \square

A.2 Stochastic results

Lemma A.5. Let $\delta > 0$, $\Omega_\delta \in \mathcal{F}$ such that $\mathbb{P}(\Omega_\delta) \geq 1 - \delta$ and $h : \Omega \times \mathbb{R}_+ \rightarrow \mathbb{R}$ a stochastic process such that $\sup_{t \geq 0} \mathbb{E}[h(\omega, t)^2] < \infty$. Then

$$\mathbb{E}[h(\omega, t) \mathbb{1}_{\Omega \setminus \Omega_\delta}] = \mathcal{O}(\sqrt{\delta}).$$

Proof. Note that $\mathbb{P}(\Omega \setminus \Omega_\delta) \leq \delta$ and

$$\mathbb{E}[h(\omega, t) \mathbb{1}_{\Omega \setminus \Omega_\delta}] \leq \sqrt{\delta} \sqrt{\mathbb{E}[h(\omega, t)^2]} \leq \sqrt{\delta} \sqrt{\sup_{t \geq 0} \mathbb{E}[h(\omega, t)^2]},$$

where we have used the Cauchy-Schwarz inequality for our first inequality. \square

Corollary A.6. Let X be the solution of (SDE) under hypotheses (H₀), (H) on f and σ , and that $\sigma_\infty \in L^2(\mathbb{R}_+)$. Then $h_1(\omega, t) = \frac{\text{dist}(X(\omega, t), \mathcal{S})^2}{2}$ and $h_2(\omega, t) = f(X(\omega, t)) - \min f$ satisfy the hypothesis of Lemma A.5, this means that there exists $C_d, C_f > 0$:

$$\begin{aligned} \mathbb{E} \left(\frac{\text{dist}(X(t), \mathcal{S})^2}{2} \right) - \mathbb{E} \left[\frac{\text{dist}(X(t), \mathcal{S})^2}{2} \mathbb{1}_{\Omega_\delta} \right] &\leq C_d \sqrt{\delta}, \\ \mathbb{E}(f(X(t)) - \min f) - \mathbb{E}[(f(X(t)) - \min f) \mathbb{1}_{\Omega_\delta}] &\leq C_f \sqrt{\delta}. \end{aligned}$$

Proof. Let $x^* \in \mathcal{S}$ be arbitrary. Using Proposition 2.4 with $\hat{\phi}(x) = \frac{\text{dist}(x, \mathcal{S})^2}{2}$, squaring it, and taking expectation, we obtain

$$\begin{aligned} \mathbb{E} \left[\frac{\text{dist}^4(X(t), \mathcal{S})}{4} \right] &\leq 3 \frac{\text{dist}(X_0, \mathcal{S})^2}{4} + 3 \left(\int_0^t \sigma_\infty^2(s) ds \right)^2 \\ &\quad + 3 \mathbb{E} \left[\left(\int_0^t \langle \sigma^\top(s, X(s))(X(s) - P_{\mathcal{S}}(X(s))), dW(s) \rangle \right)^2 \right] \\ &\leq 3 \frac{\text{dist}(X_0, \mathcal{S})^2}{4} + 3 \left(\int_0^t \sigma_\infty^2(s) ds \right)^2 + 3 \sup_{t \geq 0} \mathbb{E}[\|X(t) - x^*\|^2] \left[\int_0^t \sigma_\infty^2(s) ds \right]. \end{aligned}$$

Taking the supremum over $t \geq 0$, we obtain

$$\begin{aligned} \sup_{t \geq 0} \mathbb{E} \left[\left(\frac{\text{dist}(X(t), \mathcal{S})^2}{2} \right)^2 \right] &\leq 3 \frac{\text{dist}(X_0, \mathcal{S})^2}{4} + 3 \left(\int_0^\infty \sigma_\infty^2(s) ds \right)^2 \\ &\quad + 3 \sup_{t \geq 0} \mathbb{E}[\|X(t) - x^*\|^2] \left[\int_0^\infty \sigma_\infty^2(s) ds \right] \stackrel{\text{def}}{=} C_d < \infty. \end{aligned}$$

In the above estimation we used that $\sigma_\infty \in L^2(\mathbb{R}_+)$ and $\sup_{t \geq 0} \mathbb{E}[\|X(t) - x^*\|^2] < \infty$ by Theorem 3.1(i).

On the other hand, using Proposition 2.4 with $\tilde{\phi}(x) = f(x) - \min f$, squaring it, and taking expectation,

we obtain

$$\begin{aligned}
\mathbb{E} [[f(X(t)) - \min f]^2] &\leq 3[f(X_0) - \min f]^2 + \frac{3L}{2} \left(\int_0^t \sigma_\infty^2(s) ds \right)^2 \\
&\quad + 3\mathbb{E} \left[\left(\int_0^t \langle \sigma^\top(s, X(s)) (\nabla f(X(s))), dW(s) \rangle \right)^2 \right] \\
&\leq 3[f(X_0) - \min f]^2 + \frac{3L}{2} \left(\int_0^t \sigma_\infty^2(s) ds \right)^2 \\
&\quad + 3L^2 \sup_{t \geq 0} \mathbb{E} [\|X(t) - x^*\|^2] \left[\int_0^t \sigma_\infty^2(s) ds \right].
\end{aligned}$$

Taking the supremum over $t \geq 0$, we obtain

$$\begin{aligned}
\sup_{t \geq 0} \mathbb{E} [[f(X(t)) - \min f]^2] &\leq 3[f(X_0) - \min f]^2 + \frac{3L}{2} \left(\int_0^\infty \sigma_\infty^2(s) ds \right)^2 \\
&\quad + 3L^2 \sup_{t \geq 0} \mathbb{E} [\|X(t) - x^*\|^2] \left[\int_0^\infty \sigma_\infty^2(s) ds \right] \stackrel{\text{def}}{=} C_f < \infty.
\end{aligned}$$

□

Let us consider the Stochastic Differential Equation:

$$\begin{cases} dX(t) = F(t, X(t))dt + G(t, X(t))dW(t), & t \geq 0, \\ X(0) = X_0, \end{cases} \quad (\text{A.1})$$

where $F : \mathbb{R}_+ \times \mathbb{R}^d \rightarrow \mathbb{R}^d$, $G : \mathbb{R}_+ \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$ are measurable functions and W is a \mathcal{F}_t -adapted m -dimensional Brownian Motion.

Theorem A.7. (See [33, Theorem 5.2.1], [35, Theorem 2.4.1]) Let $F : \mathbb{R}_+ \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $G : \mathbb{R}_+ \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$ be measurable functions satisfying, for every $T > 0$:

$$\|F(t, x) - F(t, y)\| + \|G(t, x) - G(t, y)\|_F \leq C_1 \|x - y\|, \quad \forall x, y \in \mathbb{R}^d, \forall t \in [0, T], \quad (\text{A.2})$$

for some constant $C_1 \geq 0$. Then (A.1) has a unique solution $X \in S_d^l$, for every $l \geq 2$.

Proof. Condition (A.2) implies that there exists $C_2 \geq 0$ such that

$$\|F(t, x)\| + \|G(t, x)\|_F \leq C_2(1 + \|x\|), \quad \forall x \in \mathbb{R}^d, \forall t \in [0, T].$$

These are the hypotheses of [33, Theorem 5.2.1] to ensure the existence and uniqueness of the solution $X \in S_d^2$ of (A.1). Moreover, condition (A.2) implies the existence of $C_3 \geq 0$ such that

$$\langle x, F(t, x) \rangle + \|G(t, x)\|_F^2 \leq C_3(1 + \|x\|^2) \quad \forall x \in \mathbb{R}^d, \forall t \in [0, T]. \quad (\text{A.3})$$

Thus (A.3) is the necessary inequality to use [48, Lemma 3.2] and deduce that $X \in S_d^l$, for every $l \geq 2$. □

A.3 On martingales

Theorem A.8. [49] Let $(M_t)_{t \geq 0} : \Omega \rightarrow \mathbb{R}$ be a continuous martingale such that $\sup_{t \geq 0} \mathbb{E}(|M_t|^p) < \infty$ for some $p > 1$. Then there exists a random variable M_∞ such that $\mathbb{E}(|M_\infty|^p) < \infty$ and $\lim_{t \rightarrow \infty} M_t = M_\infty$ a.s.

Theorem A.9. [35, Theorem 1.3.9] Let $\{A_t\}_{t \geq 0}$ and $\{U_t\}_{t \geq 0}$ be two continuous adapted increasing processes with $A_0 = U_0 = 0$ a.s. Let $\{M_t\}_{t \geq 0}$ be a real valued continuous local martingale with $M_0 = 0$ a.s. Let ξ be a nonnegative \mathcal{F}_0 -measurable random variable. Define

$$X_t = \xi + A_t - U_t + M_t \quad \text{for } t \geq 0.$$

If X_t is nonnegative and $\lim_{t \rightarrow \infty} A_t < \infty$ a.s., then a.s. $\lim_{t \rightarrow \infty} X_t$ exists and is finite, and $\lim_{t \rightarrow \infty} U_t < \infty$.

References

- [1] R. N. Bhattacharya. Criteria for recurrence and existence of invariant measures for multidimensional diffusions. *Ann. Prob.*, 6(4):541–553, 1978.
- [2] Olivier Catoni. Simulated annealing algorithms and Markov chains with rare transitions. In *Séminaire de Probabilités XXXIII*, volume 1709 of *Lecture Notes in Mathematics*, pages 70–119. Springer, 1999.
- [3] S. Łojasiewicz. Une propriété topologique des sous-ensembles analytiques réels. In *Les Équations aux Dérivées Partielles*, pages 87–89. Editions du Centre National de la Recherche Scientifique, 1963.
- [4] S. Łojasiewicz. Ensembles semi-analytiques. *Lectures Notes IHES (Bures-sur-Yvette)*, 1965.
- [5] Hedy Attouch and Alexandre Cabot. Convergence of a relaxed inertial forward-backward algorithm for structured monotone inclusions. *Applied Mathematics and Optimization, special issue on Games, Dynamics and Optimization*, 80 (3):547–598, 2019.
- [6] Jérôme Bolte. Continuous gradient projection method in Hilbert spaces. *Journal of Optimization Theory and its Applications*, 119(2):235–259, 2003.
- [7] A.S. Antipin. Minimization of convex functions on convex sets by means of differential equations. *Differ. Uravn.*, 30(9):1475–1486, 1994.
- [8] A. Cauchy. Méthode générale pour la résolution des systèmes d’équations simultanées. *Comptes Rendus de l’Académie des Sciences de Paris*, 25:536–538, 1847.
- [9] S. Łojasiewicz. Sur les trajectoires du gradient d’une fonction analytique. *Semin. Geom., Univ. Studi Bologna*, 1982/1983:115–117, 1984.
- [10] Jérôme Bolte, Trong Phong Nguyen, Juan Peypouquet, and Bruce W. Suter. From error bounds to the complexity of first-order descent methods for convex functions. *Mathematical Programming*, 165:471–507, 2016.
- [11] T. Colding and W. Minicozzi H. Łojasiewicz inequalities and applications. *Surveys in Differential Geometry*, XIX:63–82, 2014.
- [12] Herbert Robbins and Sutton Monro. A stochastic approximation method. *Ann. Math. Statist.* 22, pages 400–407, 1951.
- [13] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic gradient algorithms. *arXiv:1511.06251*, 2017.
- [14] Antonio Orvieto and Aurelien Lucchi. Continuous-time models for stochastic optimization algorithms. *33rd Conference on Neural Information Processing Systems (NeurIPS 2019)*, 2019.
- [15] Wenqing Hu, Chris Junchi Li, Lei Li, and Jian-Guo Lui. On the diffusion approximation of nonconvex stochastic gradient descent. *arXiv:1705.07562v2*, 2018.
- [16] Bin Shi, Weijie J. Su, and Michael I. Jordan. On learning rates and Schrödinger operators. *arXiv:2004.06977*, 2020.
- [17] Zhiyuan Li, Sadhika Malladi, and Sanjeev Arora. On the validity of modeling sgd with stochastic differential equations. *arXiv:2102.12470*, 2021.

- [18] S. Soatto and P. Chaudhari. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. *2018 Information Theory and Applications Workshop (ITA)*, pages 1–10, 2018.
- [19] C.W. Gardiner. Handbook of stochastic methods. *Springer*, 3, 1985.
- [20] G. Parisi. Correlation functions and computer simulations. *Nucl. Phys. B*, 180(3):378–384, 1981.
- [21] Arnak S. Dalalyan and Avetik Karagulyan. User-friendly guarantees for the Langevin monte carlo with inaccurate gradient. *arXiv:1710.00095v3*, 2018.
- [22] Arnak S. Dalalyan. Theoretical guarantees for approximate sampling from smooth and log-concave densities. *J.R. Stat. Soc. Series B. Stat. Methodol.*, 79(3):651–676, 2017.
- [23] A. Durmus and E Moulines. High-dimensional bayesian inference via the unadjusted Langevin algorithm. *arXiv:1605.01559*, 2016.
- [24] A. Durmus and E. Moulines. Nonasymptotic convergence analysis for the unadjusted Langevin algorithm. *Ann. Appl. Probab.*, 27(3):1551–1587, 2017.
- [25] X. Cheng, N. S. Chatterji, P. L. Bartlett, and M. I. Jordan. Underdamped Langevin mcmc: A non-asymptotic analysis. *arXiv:1707.03663*, 2017.
- [26] J. Huggins and J. Zou. Quantifying the accuracy of approximate diffusions and markov chains. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, 54 of Proceedings of Machine Learning Research:382–391, 2017.
- [27] Michel Benaïm. Dynamics of stochastic approximation algorithms. In *Séminaire de Probabilités XXXIII*, volume 1709 of *Lecture Notes in Mathematics*, pages 1–69. Springer, 1999.
- [28] Panayotis Mertikopoulos and Mathias Staudigl. On the convergence of gradient-like flows with noisy gradient input. *SIAM Journal on Optimization*, 28(1):163–197, 2018.
- [29] Maxim Raginsky and Jake Bouvrie. Continuous-time stochastic mirror descent on a network: Variance reduction, consensus, convergence. *2012 IEEE 51st IEEE Conference on Decision and Control*, 2012.
- [30] Peter Bartlett and Walid Krichene. Acceleration and averaging in stochastic mirror descent dynamics. *arXiv: 1707.06219*, 2017.
- [31] Steffen Dereich and Sebastian Kassing. Cooling down stochastic differential equations: Almost sure convergence. *arXiv:2106.03510*, 2021.
- [32] R.T. Rockafellar. Convex analysis. *Princeton univeristy press*, 28, 1997.
- [33] Bernt Øksendal. Stochastic differential equations. *Springer*, 2003.
- [34] Etienne Pardoux and Aurel Rascanu. *Stochastic differential equations, backward SDEs, partial differential equations*. Springer, 2014.
- [35] Xuerong Mao. Stochastic differential equations and applications. *Elsevier*, 2007.
- [36] Bálint Farkas and Sven-Ake Wegner. Variations on Barbalat’s lemma. *The American Mathematical Monthly*, 123:8:825–830, 2016.
- [37] Z. Opial. Weak convergence of the sequence of successive approximations for nonexpansive mappings. *Bull. Amer. Math. Soc.*, 73:591–597, 1967.
- [38] H. Attouch, J. Bolte, and B. F. Svaiter. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, Forward–Backward splitting, and regularized Gauss–Seidel methods. *Mathematical Programming*, 137(1-2):91–129, 2013.
- [39] Pierre Frankel, Guillaume Garrigos, and Juan Peypouquet. Splitting methods with variable metric for Kurdyka–Łojasiewicz functions and general convergence rates. *Journal of Optimization Theory and Applications*, 165(3):874–900, 2014.
- [40] Jong-Shi Pang. Error bounds in mathematical programming. *Mathematical Programming*, 79(1):299–332, 1997.
- [41] Francisco Javier Aragón Artacho and Michel Geoffroy. Characterization of metric regularity of subdifferentials. *Journal of Convex Analysis*, 15:365–380, 01 2008.
- [42] Alexander Y. Kruger. Error bounds and hölder metric subregularity. *Set-Valued and Variational Analysis*, 23(4):705–736, 2015.
- [43] Alexander Y. Kruger. Error bounds and metric subregularity. *Optimization*, 64(1):49–79, 2015.

- [44] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [45] Amir Beck and Marc Teboulle. Smoothing and first order methods: A unified framework. *SIAM Journal on Optimization*, 22(2):557–580, 2012.
- [46] Radoslaw Matusik, Andrzej Nowakowski, Slawomir Plaskacz, and Andrzej Rogoswski. Finite-time stability for differential inclusions with applications to neural networks. *arXiv:1804.08440v2*, 2019.
- [47] Walter Rudin. Real and complex analysis. *McGraw-Hill*, 1987.
- [48] Desmond J. Higham, Xuerong Mao, and Andrew M. Stuart. Strong convergence of Euler-type methods for nonlinear stochastic differential equations. *SIAM J. Numer. Anal.*, 40(3), pages 1041–1063, 2006.
- [49] J.L. Doob. Stochastic processes. *Wiley*, 1991.