



HAL
open science

An SDE Perspective on Stochastic Inertial Gradient Dynamics with Time-Dependent Viscosity and Geometric Damping

Rodrigo Maulen-Soto, Jalal Fadili, Hedy Attouch, Peter Ochs

► **To cite this version:**

Rodrigo Maulen-Soto, Jalal Fadili, Hedy Attouch, Peter Ochs. An SDE Perspective on Stochastic Inertial Gradient Dynamics with Time-Dependent Viscosity and Geometric Damping. 2024. hal-04637159

HAL Id: hal-04637159

<https://hal.science/hal-04637159>

Preprint submitted on 5 Jul 2024

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 Inertial Gradient Dynamics with Time-Dependent Viscosity and Geometric Damping

Rodrigo Maulen-Soto^{*} Jalal Fadili[†] Hedy Attouch[‡] Peter Ochs[§]

Abstract. Our approach is part of the close link between continuous dissipative dynamical systems and optimization algorithms. We aim to solve convex minimization problems by means of stochastic inertial differential equations which are driven by the gradient of the objective function. This will provide a general mathematical framework for analyzing fast optimization algorithms with stochastic gradient input. Our study is a natural extension of our previous work devoted to the first-order in time stochastic steepest descent. Our goal is to develop these results further by considering second-order stochastic differential equations in time, incorporating a viscous time-dependent damping and a Hessian-driven damping. To develop this program, we rely on stochastic Lyapunov analysis. Assuming a square-integrability condition on the diffusion term times a function dependant on the viscous damping, and that the Hessian-driven damping is a positive constant, our first main result shows that almost surely, there is convergence of the values, and states fast convergence of the values in expectation. Besides, in the case where the Hessian-driven damping is zero, we conclude with the fast convergence of the values in expectation and in almost sure sense, we also managed to prove almost sure weak convergence of the trajectory. We provide a comprehensive complexity analysis by establishing several new pointwise and ergodic convergence rates in expectation for the convex and strongly convex case.

Key words. Stochastic optimization, Inertial gradient system, Convex optimization, Stochastic Differential Equation, Time-dependent viscosity, Convergence rate, Asymptotic behavior.

AMS subject classifications. 37N40, 46N10, 49M99, 65B99, 65K05, 65K10, 90B50, 90C25, 60H10, 90C53, 60G12

1 Introduction

1.1 Problem Statement

Let us fix the framework of our study. We consider the minimization problem

$$\min_{x \in \mathbb{H}} f(x), \tag{P}$$

where \mathbb{H} is a real Hilbert space and the objective function $f : \mathbb{H} \rightarrow \mathbb{R}$ satisfies the following standing assumptions:

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

To solve (P), a fundamental dynamic is the gradient flow system:

$$\begin{cases} \dot{x}(t) + \nabla f(x(t)) = 0, & t > t_0; \\ x(t_0) = x_0. \end{cases} \tag{GF}$$

This dynamic is known to yield a convergence rate of $\mathcal{O}(t^{-1})$ (in fact even $o(t^{-1})$) on the values. Second-order inertial dynamical systems have been introduced to provably accelerate the convergence behaviour. Among them, the Inertial

^{*}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

[§]Department of Mathematics and Computer Science, Saarland University, Germany, E-mail: ochs@math.uni-saarland.de

System with Implicit Hessian Damping (ISIHD) is the following differential equation starting at $t_0 > 0$ with initial condition $x_0, v_0 \in \mathbb{H}$:

$$\begin{cases} \ddot{x}(t) + \gamma(t)\dot{x}(t) + \nabla f(x(t) + \beta(t)\dot{x}(t)) = 0, & t > t_0; \\ x(t_0) = x_0, \quad \dot{x}(t_0) = v_0, \end{cases} \quad (\text{ISIHD})$$

where $\gamma, \beta : [t_0, +\infty[\rightarrow \mathbb{R}_+$. (ISIHD) was first considered by [1]; see also [2, 3]. Following the physical interpretation of this ODE, we call the non-negative parameters γ and β as the viscous and geometric damping parameters, respectively. Also, this ODE was found to have a smoothing effect on the energy error and oscillations [1, 2, 3]. The use of the term ‘‘implicit’’ comes from the fact that by Taylor expansion (as $t \rightarrow +\infty$) one has

$$\nabla f(x(t) + \beta(t)\dot{x}(t)) \approx \nabla f(x(t)) + \beta(t)\nabla^2 f(x(t))\dot{x}(t), \quad (1.1)$$

thus making appear the Hessian-driven damping with coefficient $\beta(t)$.

In many practical situations, the gradient evaluation is subject to stochastic errors. This is for example the case if the cost per iteration is very high and thus cheap and random approximations of the gradient are necessary. These errors can also be due to some other exogenous factor. The continuous-time approach through stochastic differential equations (SDE) is a powerful way to model these errors in a unified way, and stochastic algorithms can then be viewed as time-discretizations. In fact, several recent works have used the dynamic

$$\begin{cases} dX(t) = -\nabla f(X(t))dt + \sigma(t, X(t))dW(t), \\ X(t_0) = X_0. \end{cases} \quad (1.2)$$

to model SGD-type algorithms; (see *e.g.* [4, 5, 6, 7, 8, 9, 10, 11, 12, 13]). In fact, the continuous-time perspective offers a deep insight and unveils the key properties of the dynamic without being tied to a specific discretization.

In this setting, we can model the associated errors using a stochastic integral with respect to the measure defined by a continuous Itô martingale. This entails the following stochastic differential equation (SDE for short), which is the stochastic counterpart of (ISIHD):

$$\begin{cases} dX(t) &= V(t)dt, \\ dV(t) &= -\gamma(t)V(t)dt - \nabla f(X(t) + \beta(t)V(t))dt + \sigma(t, X(t) + \beta(t)V(t))dW(t), \\ X(t_0) &= X_0, \quad V(t_0) = V_0. \end{cases} \quad (\text{S} - \text{ISIHD})$$

They are defined over a filtered probability space $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$, where the diffusion (volatility) term $\sigma : [t_0, +\infty[\times \mathbb{H} \rightarrow \mathcal{L}_2(\mathbb{K}; \mathbb{H})$ is a measurable function, and W is a \mathbb{K} -valued Brownian motion. Where $X_0, V_0 \in L^\nu(\Omega; \mathbb{H})$ for some $\nu \geq 2$ are given initial data. Besides, γ and β are parameters called viscous damping and geometric damping, respectively. They are explained and discussed in more detail below.

Our goal is to provide a general mathematical framework for analyzing fast gradient-based optimization algorithms with stochastic gradient input. For this, we will study second-order stochastic differential equations in time, *i.e.*, also involving acceleration, and whose drift term is the gradient of the function to be minimized. In this context, considering inertial dynamics with a time-dependent viscosity coefficient is a key property to obtain fast convergent methods. Our study is related to two recent works:

- On the one hand, it is a natural extension of the article by Maulen-Soto, Fadili and Attouch [11] devoted to the first-order in time stochastic steepest descent.
- On the other hand, we will rely on the Lyapunov analysis for the dynamic (ISIHD) done by Attouch, Fadili, and Kungurtsev in [3]; and the one for (IGS $_\gamma$) done by Attouch and Cabot in [14].

More precisely, our goal is to study the dynamic (S – ISIHD) and its long-time behavior in order to solve (P). We develop an independent analysis and do not use the results of the deterministic case, which will serve as a useful comparison since we will be able to recover the known results for (ISIHD) when there is no noise (*i.e.* $\sigma = 0$) in the dynamic (S – ISIHD).

Let us first identify the assumptions needed to expect that the position state of **(S – ISIHD)** approaches \mathcal{S} in the long run. In the case where $\mathbb{H} = \mathbb{K} = \mathbb{R}^d$, $\gamma(\cdot) \equiv \gamma > 0$, $\beta \equiv 0$, and $\sigma(\cdot, \cdot) \equiv \sigma I_d$ ($\sigma > 0$), under mild assumptions one can show that **(S – ISIHD)** has a unique invariant distribution π_σ in (x, v) with density $\propto \exp\left(-\frac{2\gamma}{\sigma^2}\left(f(x) + \frac{\|v\|^2}{2}\right)\right)$, see e.g., [15, Proposition 6.1]. Clearly, as $\sigma \rightarrow 0^+$, π_σ gets concentrated around $\operatorname{argmin} f \times \{0_d\}$ as σ tends to 0^+ , with $\lim_{\sigma \rightarrow 0^+} \pi_\sigma(\operatorname{argmin} f \times \{0_d\}) = 1$, see e.g. [16]. Motivated by these observations and the fact that we aim to exactly solve **(P)**, our paper will then mainly focus on the case where $\sigma(\cdot, x)$ vanishes fast enough as $t \rightarrow +\infty$ uniformly in x .

Our objectives are largely motivated by recent analysis in the deterministic setting. In fact, the dynamic **(S – ISIHD)** comes naturally as a stochastic version of **(ISIHD)**. **(ISIHD)** is one of the most recent developments regarding the use of gradient-based dynamic systems for optimization. Let us briefly recall the steps that led to its emergence. In this regard, let us stress the importance of working with a time-dependent viscosity coefficient $\gamma(t)$. It is with the introduction of the non-autonomous inertial dynamics

$$\ddot{x}(t) + \frac{\alpha}{t}\dot{x}(t) + \nabla f(x(t)) = 0, \quad t > t_0, \quad (1.3)$$

that Su, Boyd, and Candès [17] showed the rate of convergence $1/t^2$ of the values, thus making the link with the accelerated gradient method of Nesterov [18]. Since then, abundant literature has been devoted to the study of inertial dynamics with time-dependent viscosity coefficient

$$\ddot{x}(t) + \gamma(t)\dot{x}(t) + \nabla f(x(t)) = 0, \quad t > t_0, \quad (\text{IGS}_\gamma)$$

where the importance of the case $\gamma(t) = \frac{\alpha}{t}$, and the subtle tuning of the parameter α is elucidated. Indeed, α must be taken greater than or equal to 3 for getting the rate of convergence $\mathcal{O}(1/t^2)$ of the values, and $\alpha > 3$ provides an even better rate of convergence with little o instead of big \mathcal{O} ; see Attouch-Cabot [19] and Attouch-Peypouquet [20].

However, because of the inertial aspects, and the asymptotic vanishing viscous damping coefficient, **(IGS $_\gamma$)** may exhibit many small oscillations which are not desirable from an optimization point of view. To remedy this, a powerful tool consists in introducing into the dynamic a geometric damping driven by the Hessian of f . This gives the Inertial System with Explicit Hessian-driven Damping

$$\ddot{x}(t) + \gamma(t)\dot{x}(t) + \beta(t)\nabla^2 f(x(t))\dot{x}(t) + b(t)\nabla f(x(t)) = 0, \quad (\text{ISEHD})$$

where γ and β are, the already presented, damping parameters, and b is a time scale parameter. This dynamic is the explicit version of **(ISIHD)**. The time discretization of this system has been studied by Attouch, Chbani, Fadili, and Riahi [21]. It provides a rich family of first-order methods for minimizing f . At first glance, the presence of the Hessian may seem to entail numerical difficulties. However, this is not the case as the Hessian intervenes in the above ODE in the form $\nabla^2 f(x(t))\dot{x}(t)$, which is nothing but the time derivative of $t \mapsto \nabla f(x(t))$. This explains why the time discretization of this dynamic provides first-order algorithms. On the contrary, the time-continuous dynamics can be argued to be truly of second-order nature, i.e., close to Newton’s and Levenberg-Marquardt’s dynamics [22]. This understanding suggests that **(ISIHD)** may represent the nature of first-order algorithms better than **(ISEHD)**. Let us recall that in our stochastic setting, we do not have direct access to evaluate the gradient of f . Instead, we model the associated errors with a continuous Itô martingale (denoted as $M(t)$). Therefore, it is meaningless to ask for the time derivative of $\nabla f(X(t)) + M(t)$ because (non-constant) martingales are not differentiable a.s.. This is why we are going to focus on the implicit form of the Hessian-driven damping **(S – ISIHD)**.

1.2 Contributions

Our main contributions are the following:

- Given that the Lyapunov analysis has already been done in the case of the first-order in-time stochastic gradient system from which our inertial system is derived (see [11, 13]), our analysis is greatly simplified. This allows us to show almost sure convergence of the trajectory and convergence rates in expectation for the case with time-dependent coefficients $\gamma(t)$ and a particular choice of $\beta(t)$.
- We will develop a Lyapunov analysis to obtain convergence rates, integral estimates, and almost sure results in the general case of coefficients $\gamma(t)$ and $\beta(t)$.

- In the case where the coefficient $\beta(t)$ is zero, we show that under some hypotheses, we have almost sure convergence of the trajectory, convergence rates, and integral estimates. As a special case, we focus on viscous damping coefficient $\gamma(t) = \frac{\alpha}{t^r}$, $r \in [0, 1]$, $\alpha \geq 1 - r$.

1.3 Relation to prior work

Kinetic diffusion dynamics for sampling Let us consider (S – ISIHD) in the case where $\mathbb{H} = \mathbb{K} = \mathbb{R}^d$, $\gamma(t) = \gamma > 0$, $\beta(t) = 0$ and $\sigma = \sqrt{2\gamma}I_d$. Then one recovers the kinetic Langevin diffusion (or second-order Langevin process). In this case, the continuous-time Markov process $(X(t), V(t))$ is positive recurrent and has a unique invariant distribution which has the density $\propto \exp\left(-f(x) - \frac{\|v\|^2}{2}\right)$ with respect to the Lebesgue measure on \mathbb{R}^{2d} . Time-discretized versions of this Langevin diffusion process have been studied in the literature to (approximately) sample from $\propto \exp(-f(x))$ with asymptotic and non-asymptotic convergence guarantees in various topologies and under various conditions have been studied; see [23, 24, 25] and references therein.

Inexact inertial gradient systems There is an abundant literature regarding the dynamics (ISIHD) and (ISEHD), either in the exact case or with errors but only deterministic ones; see [1, 3, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35]). We are not aware of any such work in the stochastic case. Only a few papers have been devoted to studying the second-order in-time inertial stochastic gradient systems with viscous damping, *i.e.* stochastic versions of (IGS $_\gamma$), either with vanishing damping $\gamma(t) = \alpha/t$ or constant damping $\gamma(t)$ (stochastic HBF); see *e.g.* [12, 36, 37]. For instance, [12] provide asymptotic $\mathcal{O}(1/t^2)$ convergence rate on the objective values in expectation under integrability conditions on the diffusion term as well as other rates under additional geometrical properties of the objective. The corresponding stochastic algorithms for these two choices of γ , whose mathematical formulation and analysis is simpler, have been the subject of active research work; see *e.g.* [38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51].

Time scaling and averaging An SDE to solve (P) has been thoroughly studied in [11]; see also [13] for the non-smooth setting. This SDE has the form

$$\begin{cases} dX(t) = -\nabla f(X(t))dt + \sigma(t, X(t))dW(t), & t \geq t_0, \\ X(t_0) = X_0. \end{cases} \quad (1.4)$$

The authors in [52] proposed time scaling and averaging to link (GF) and (ISIHD) with a general viscous damping function γ and a properly adjusted geometric damping function β (related to γ). Moreover, in [53] we extended the results of [52] to the stochastic case. Leveraging the techniques shown in [52] with a general function γ and an appropriate β , we were able to transfer all the results we obtained in [11] for (1.4) to (S – ISIHD). This avoids, in particular, going through an intricate and dedicated Lyapunov analysis for (S – ISIHD) at the cost of having a particular β (related to γ). A local convergence analysis also became easily accessible through those lenses while it is barely possible otherwise. We also specialized those results to a standard case where $\gamma(t) = \frac{\alpha}{t}$ and $\beta(t) = \frac{t}{\alpha-1}$. However, a distinguished analysis of the advantages of adding a Hessian-driven damping term requires a fine choice of β that is independent of γ . Therefore, we go through that dedicated Lyapunov analysis for (S – ISIHD) and consider the more general case for γ and β .

2 Notation and Preliminaries

We will use the following shorthand notations: Given $n \in \mathbb{N}$, $[n] \stackrel{\text{def}}{=} \{1, \dots, n\}$. Consider \mathbb{H}, \mathbb{K} real separable Hilbert spaces endowed with the inner product $\langle \cdot, \cdot \rangle_{\mathbb{H}}$ and $\langle \cdot, \cdot \rangle_{\mathbb{K}}$, respectively, and norm $\|\cdot\|_{\mathbb{H}} = \sqrt{\langle \cdot, \cdot \rangle_{\mathbb{H}}}$ and $\|\cdot\|_{\mathbb{K}} = \sqrt{\langle \cdot, \cdot \rangle_{\mathbb{K}}}$, respectively (we omit the subscripts \mathbb{H} and \mathbb{K} for the sake of clarity). $I_{\mathbb{H}}$ is the identity operator on \mathbb{H} . $\mathcal{L}(\mathbb{K}; \mathbb{H})$ is the space of bounded linear operators from \mathbb{K} to \mathbb{H} , $\mathcal{L}_1(\mathbb{K})$ is the space of trace-class operators, and $\mathcal{L}_2(\mathbb{K}; \mathbb{H})$ is the space of bounded linear Hilbert-Schmidt operators from \mathbb{K} to \mathbb{H} . For $M \in \mathcal{L}_1(\mathbb{K})$, its trace is defined by

$$\text{tr}(M) \stackrel{\text{def}}{=} \sum_{i \in I} \langle M e_i, e_i \rangle < +\infty,$$

where $I \subseteq \mathbb{N}$ and $(e_i)_{i \in I}$ is an orthonormal basis of \mathbb{K} . Besides, for $M \in \mathcal{L}(\mathbb{K}; \mathbb{H})$, $M^* \in \mathcal{L}(\mathbb{H}; \mathbb{K})$ is the adjoint operator of M , and for $M \in \mathcal{L}_2(\mathbb{K}; \mathbb{H})$,

$$\|M\|_{\text{HS}} \stackrel{\text{def}}{=} \sqrt{\text{tr}(MM^*)} < +\infty$$

is its Hilbert-Schmidt norm (in the finite-dimensional case is equivalent to the Frobenius norm). We denote by $w\text{-lim}$ the limit for the weak topology of \mathbb{H} . The notation $A : \mathbb{H} \rightrightarrows \mathbb{H}$ means that A is a set-valued operator from \mathbb{H} to \mathbb{H} . Consider $f : \mathbb{H} \rightarrow \mathbb{R}$, the sublevel of f at height $r \in \mathbb{R}$ is denoted $[f \leq r] \stackrel{\text{def}}{=} \{x \in \mathbb{H} : f(x) \leq r\}$. For $1 \leq p \leq +\infty$, $L^p([a, b])$ is the space of measurable functions $g : \mathbb{R} \rightarrow \mathbb{R}$ such that $\int_a^b |g(t)|^p dt < +\infty$, with the usual adaptation when $p = +\infty$. For two functions $f, g : \mathbb{R} \rightarrow \mathbb{R}$ we will denote $f \sim g$ as $t \rightarrow +\infty$, if $\lim_{t \rightarrow +\infty} \frac{f(t)}{g(t)} = 1$. On the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, $L^p(\Omega; \mathbb{H})$ denotes the (Bochner) space of \mathbb{H} -valued random variables whose p -th moment (with respect to the measure \mathbb{P}) is finite. Other notations will be explained when they first appear.

Let us recall some important definitions and results from convex analysis; for a comprehensive coverage, we refer the reader to [54].

We denote by $\Gamma_0(\mathbb{H})$ the class of proper lsc and convex functions on \mathbb{H} taking values in $\mathbb{R} \cup \{+\infty\}$. For $\mu > 0$, $\Gamma_\mu(\mathbb{H}) \subset \Gamma_0(\mathbb{H})$ is the class of μ -strongly convex functions, roughly speaking, this means that there exists a quadratic lower bound on the growth of these functions. We denote by $C^s(\mathbb{H})$ the class of s -times continuously differentiable functions on \mathbb{H} . For $L \geq 0$, $C_L^{1,1}(\mathbb{H}) \subset C^1(\mathbb{H})$ is the set of functions on \mathbb{H} whose gradient is L -Lipschitz continuous, and $C_L^2(\mathbb{H})$ is the subset of $C_L^{1,1}(\mathbb{H})$ whose functions are twice differentiable.

The class of $C_L^{1,1}(\mathbb{H})$ functions enjoys the well-known *descent lemma* which plays a central role in the analysis of optimization dynamics.

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

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

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

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

On stochastic differential equations For the necessary notation and preliminaries on stochastic processes, see [13, Section A.2]. Moreover, the existence and uniqueness of a solution of **(S – ISIHD)** is discussed in Proposition A.9.

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

Proposition 2.3. [55, Section 2.3] *Consider (X, V) a solution of **(S – ISIHD)** and W a K -valued Brownian motion, let $\phi : [t_0, +\infty[\times \mathbb{H} \times \mathbb{H} \rightarrow \mathbb{R}$ be such that $\phi(\cdot, x, v) \in C^1([t_0, +\infty[)$ for every $x, v \in \mathbb{H}$, $\phi(t, \cdot, \cdot) \in C^2(\mathbb{H} \times \mathbb{H})$ for every $t \geq t_0$. Then the process*

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

is an Itô Process, such that for all $t \geq t_0$

$$\begin{aligned} Y(t) = & Y(t_0) + \int_{t_0}^t \frac{\partial \phi}{\partial t}(s, X(s), V(s)) ds + \int_{t_0}^t \langle \nabla_x \phi(s, X(s), V(s)), V(s) \rangle ds \\ & - \int_{t_0}^t \langle \nabla_v \phi(s, X(s), V(s)), \gamma(s)V(s) + \nabla f(X(s) + \beta(s)V(s)) \rangle ds \\ & + \int_{t_0}^t \langle \sigma^*(s, X(s) + \beta(s)V(s)) \nabla_v \phi(s, X(s), V(s)), dW(s) \rangle \\ & + \frac{1}{2} \int_{t_0}^t \text{tr}[\sigma(s, X(s) + \beta(s)V(s)) \sigma^*(s, X(s) + \beta(s)V(s)) \nabla_v^2 \phi(s, X(s), V(s))] ds, \end{aligned} \tag{2.1}$$

where ∇_v^2 is the Hessian with respect to the double differentiation of v and σ^* is the adjoint operator of σ . Moreover, if for all $T > t_0$

$$\mathbb{E} \left(\int_{t_0}^T \|\sigma^*(s, X(s) + \beta(s)V(s))\nabla_v \phi(s, X(s), V(s))\|^2 ds \right) < +\infty,$$

then $\int_{t_0}^t \langle \sigma^*(s, X(s) + \beta(s)V(s))\nabla_v \phi(s, X(s), V(s)), dW(s) \rangle$ is a square-integrable continuous martingale and

$$\mathbb{E} \left(\int_{t_0}^t \langle \sigma^*(s, X(s) + \beta(s)V(s))\nabla_v \phi(s, X(s), V(s)), dW(s) \rangle \right) = 0 \quad (2.2)$$

3 (S – ISIHD) with general γ and β

In this section, we will develop a Lyapunov analysis based on [3] to study almost sure, and in expectation properties of the dynamic (S – ISIHD), when the parameters γ and β are general functions. This will allow to go much further and consider parameters not covered in [53] which exploits the relationship between first-order and second-order systems. We will also apply our results to two special cases: (i) γ is a differentiable, decreasing and vanishing function, with vanishing derivative, and β is a positive constant; and (ii) $\gamma(t) = \frac{\alpha}{t}$, and $\beta(t) = \gamma_0 + \frac{\beta}{t}$ (with $\gamma_0, \beta > 0$). These cases are again not covered by results in [53].

Recall that our focus in this paper is on an optimization perspective, and as we argued in the introduction, we will study the long time behavior of (S – ISIHD) as the diffusion term vanishes when $t \rightarrow +\infty$. Therefore, throughout the paper, we assume that the diffusion (volatility) term σ satisfies:

$$\begin{cases} \sup_{t \geq t_0, x \in \mathbb{H}} \|\sigma(t, x)\|_{\text{HS}} < +\infty, \\ \|\sigma(t, x') - \sigma(t, x)\|_{\text{HS}} \leq l_0 \|x' - x\|, \end{cases} \quad (\text{H}_\sigma)$$

for some $l_0 > 0$ and for all $t \geq t_0, x, x' \in \mathbb{H}$. The Lipschitz continuity assumption is mild and required to ensure the well-posedness of (S – ISIHD).

Remark 3.1. Under the hypothesis (H_σ) we have that there exists $\sigma_*^2 > 0$ such that

$$\|\sigma(t, x)\|_{\text{HS}}^2 \leq \sigma_*^2, \quad \forall t \geq t_0, \forall x \in \mathbb{H}.$$

Let us also define $\sigma_\infty : [t_0, +\infty[\rightarrow \mathbb{R}_+$ as: $\sigma_\infty(t) \stackrel{\text{def}}{=} \sup_{x \in \mathbb{H}} \|\sigma(t, x)\|_{\text{HS}}$.

Now, we follow with the hypotheses we will require over the damping parameters.

For $t_0 > 0$, let $\gamma : [t_0, +\infty[\rightarrow \mathbb{R}_+$ be a viscous damping, we denote

$$p(t) \stackrel{\text{def}}{=} \exp \left(\int_{t_0}^t \gamma(s) ds \right). \quad (3.1)$$

Besides, if

$$\int_{t_0}^{\infty} \frac{ds}{p(s)} < +\infty, \quad (\text{H}_\gamma)$$

we define $\Gamma : [t_0, +\infty[\rightarrow \mathbb{R}_+$ by

$$\Gamma(t) \stackrel{\text{def}}{=} p(t) \int_t^{\infty} \frac{ds}{p(s)}. \quad (3.2)$$

Remark 3.2. Let us notice that Γ satisfies the relation

$$\Gamma' = \gamma\Gamma - 1.$$

For $t_0 > 0$, let $\beta : [t_0, +\infty[\rightarrow \mathbb{R}_+$ be a geometric damping that we will assume to be a differentiable function. We will occasionally need to impose the additional assumption:

there exists $c_1, c_2 > 0$, and $t_1 > t_0$ such that

$$\begin{aligned} \beta(t) &\leq c_1, \\ \left| \frac{\beta'(t) - \gamma(t)\beta(t) + 1}{\beta(t)} \right| &\leq c_2, \forall t \geq t_1. \end{aligned} \tag{H_\beta}$$

We recall also that $\mathcal{S} \stackrel{\text{def}}{=} \text{argmin}(f)$.

3.1 Reformulation of (S – ISIHD)

The formulation of the dynamic (S – ISIHD) is known as the Hamiltonian formulation. However, it is not the only one. In the deterministic case, an alternative equivalent and more flexible first-order reformulation of (ISIHD) was proposed in [3]. The motivation there was that this equivalent reformulation can handle the case where f is non-smooth. Although we will not consider the non-smooth case here, we will still extend and use that equivalent reformulation to the stochastic case.

Consider the dynamic (S – ISIHD), and let us define the auxiliary variable

$$Y(t) = X(t) + \beta(t)V(t), \quad t > t_0.$$

We have that

$$\begin{aligned} dY(t) &= dX(t) + \beta'(t)V(t) + \beta(t)dV(t) \\ &= -\beta(t)\nabla f(Y(t))dt - (\beta'(t) - \gamma(t)\beta(t) + 1) \left(\frac{X(t) - Y(t)}{\beta(t)} \right) dt + \beta(t)\sigma(t, Y(t))dW(t). \end{aligned}$$

So we can reformulate (S – ISIHD) in terms of X, Y in the following way:

$$\begin{cases} dX(t) &= - \left(\frac{X(t) - Y(t)}{\beta(t)} \right) dt, \quad t > t_0, \\ dY(t) &= -\beta(t)\nabla f(Y(t))dt - (\beta'(t) - \gamma(t)\beta(t) + 1) \left(\frac{X(t) - Y(t)}{\beta(t)} \right) dt + \beta(t)\sigma(t, Y(t))dW(t), \quad t > t_0, \\ X(t_0) &= X_0, \quad Y(t_0) = X_0 + \beta(t_0)V_0, \end{cases} \tag{ISIHD – S_R}$$

where the subscript 'R' indicates that this is a reformulation. Moreover, we can reformulate (ISIHD – S_R) in the product space $\mathbb{H} \times \mathbb{H}$ by setting $Z(t) = (X(t), Y(t)) \in \mathbb{H} \times \mathbb{H}$, and thus (ISIHD – S_R) can be equivalently written as

$$\begin{cases} dZ(t) &= -\beta(t)\nabla \mathcal{G}(Z(t))dt - \mathcal{D}(t, Z(t))dt + \hat{\sigma}(t, Z(t))dW(t), \quad t > t_0, \\ Z(t_0) &= (X_0, X_0 + \beta(t_0)V_0), \end{cases} \tag{3.3}$$

where $\mathcal{G} : \mathbb{H} \times \mathbb{H} \rightarrow \mathbb{R}$ is the convex function defined as $\mathcal{G}(Z) = f(Y)$, and the time-dependent operator $\mathcal{D} : [t_0, +\infty[\times \mathbb{H} \times \mathbb{H} \rightarrow \mathbb{H} \times \mathbb{H}$ is given by

$$\mathcal{D}(t, Z) = \left(\frac{1}{\beta(t)}(X - Y), \frac{\beta'(t) - \gamma(t)\beta(t) + 1}{\beta(t)}(X - Y) \right), \tag{3.4}$$

and the stochastic noise $\hat{\sigma} \in \mathcal{M}_{2 \times 2}(\mathcal{L}_2(\mathbb{K}; \mathbb{H}))$ defined by $\hat{\sigma}(t, Z) = \begin{pmatrix} 0 & 0 \\ 0 & \beta(t)\sigma(t, Y) \end{pmatrix}$, and $W(t) = (W_1(t), W_2(t))$, where W_1, W_2 are two independent \mathbb{K} -valued Brownian motions.

3.2 Fast convergence properties: convex case

To obtain properties in almost sure sense and in expectation of **(S – ISIHD)**, we are going to adapt the Lyapunov analysis shown on [3] for the dynamic **(ISIHD)**.

To that purpose, let us consider $t_1 > t_0$, $\gamma, \beta : [t_0, +\infty[\rightarrow \mathbb{R}_+$ be fixed functions and let $a, b, c, d : [t_0, +\infty[\rightarrow \mathbb{R}$ be differentiable functions (on $]t_0, +\infty[$) satisfying the following system for all $t > t_1$:

$$\begin{cases} a'(t) - b(t)c(t) & \leq 0 \\ -a(t)\beta(t) & \leq 0 \\ -a(t)\gamma(t)\beta(t) + a(t)\beta'(t) + a(t) - c(t)^2 + b(t)c(t)\beta(t) & = 0 \\ b'(t)b(t) + \frac{d'(t)}{2} & \leq 0 \\ b'(t)c(t) + b(t)(b(t) + c'(t) - c(t)\gamma(t)) + d(t) & = 0 \\ c(t)(b(t) + c'(t) - c(t)\gamma(t)) & \leq 0. \end{cases} \quad (\mathbf{S}_{a,b,c,d})$$

Given $x^* \in \mathcal{S}$, we consider

$$\mathcal{E}(t, x, v) = a(t)(f(x + \beta(t)v) - \min(f)) + \frac{1}{2}\|b(t)(x - x^*) + c(t)v\|^2 + \frac{d(t)}{2}\|x - x^*\|^2. \quad (3.5)$$

Remark 3.3. It was shown in [1, Section 3.1] and [3, Lemma 1] that energy function \mathcal{E} with a, b, c, d satisfying the system **(S_{a,b,c,d})** is a Lyapunov function for **(ISIHD)** when $\gamma(t) = \frac{\alpha}{t}$ (with $\alpha > 3$) and $\beta(t) = \gamma_0 + \frac{\beta}{t}$ (with $\gamma_0, \beta \geq 0$), hence, useful to obtain convergence guarantees of that dynamic. We will see that the same system **(S_{a,b,c,d})** also covers the case of general coefficients γ and β , hence providing insights on the convergence properties of **(S – ISIHD)** when one can find the corresponding functions a, b, c, d .

In the following proposition, we state an abstract integral bound, as well as almost sure and in expectation convergence properties for **(S – ISIHD)**.

Proposition 3.4. Assume that f, σ satisfy **(H₀)** and **(H_σ)**, respectively. Let $\nu \geq 2$, and consider the dynamic **(S – ISIHD)** with initial data $X_0, V_0 \in L^\nu(\Omega; \mathbb{H})$. Consider also γ, β from **(S – ISIHD)** satisfying **(H_γ)** and **(H_β)**. Then, there exists a unique solution $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$ of **(S – ISIHD)**. Moreover, if $t \mapsto m(t)\sigma_\infty^2(t) \in L^1([t_0, +\infty[)$, where $m(t) \stackrel{\text{def}}{=} \max\{1, a(t), c^2(t)\}$, then the following statements hold:

(i) If $b(t)c(t) - a'(t) = \mathcal{O}(c(t)(\gamma(t)c(t) - c'(t) - b(t)))$, then

$$\int_{t_0}^{\infty} (b(s)c(s) - a'(s))(f(X(s)) - \min f + \|V(s)\|^2) ds < +\infty \quad a.s..$$

(ii) If there exists $\eta > 0, \hat{t} > t_0$ such that

$$\eta \leq c(t)(\gamma(t)c(t) - c'(t) - b(t)), \quad \eta \leq a(t)\beta(t), \quad \gamma(t) \leq \eta, \quad \forall t > \hat{t},$$

then $\lim_{t \rightarrow \infty} \|V(t)\| = 0$ a.s., $\lim_{t \rightarrow \infty} \|\nabla f(X(t) + \beta(t)V(t))\| = 0$ a.s., and $\lim_{t \rightarrow \infty} \|\nabla f(X(t))\| = 0$ a.s.

(iii) If there exists $D > 0, \tilde{t} > t_0$ such that $d(t) \geq D$ for $t > \tilde{t}$, then :

$$\mathbb{E}(\|V(t)\|^2) = \mathcal{O}\left(\frac{1 + b^2(t)}{c^2(t)}\right),$$

and

$$\mathbb{E}(f(X(t)) - \min f) = \mathcal{O}\left(\max\left\{\frac{1}{a(t)}, \frac{\beta(t)\sqrt{1 + b^2(t)}}{\sqrt{a(t)c(t)}}, \frac{\beta^2(t)(1 + b^2(t))}{c^2(t)}\right\}\right).$$

This is a compact version that extracts only the most important points from the more detailed and complete Propositions **A.13, A.14** which are proved in the appendix.

The complete version of the previous proposition (*i.e.* Propositions A.13 and A.14) generalizes the results proved in [3] to the stochastic setting. However, they lack practical use if we cannot exhibit a, b, c, d functions that satisfy $(S_{a,b,c,d})$. Although we are not able to solve this system in general, in Corollaries 3.5 and 3.7 we will specify some particular cases for γ and β where such functions a, b, c, d can be exhibited to satisfy the system $(S_{a,b,c,d})$.

The following corollary provides a specific case where a solution to the system $(S_{a,b,c,d})$ can be exhibited, which was not discussed in [1, 3]. Moreover, we show the implications it has on the stochastic setting.

Corollary 3.5 (Decreasing and vanishing γ , with vanishing γ' and positive constant β). *Consider the context of Proposition 3.4 in the case where $\beta(t) \equiv \beta > 0$, γ satisfying (H_γ) , such that it is a differentiable, decreasing, and vanishing function, with $\lim_{t \rightarrow \infty} \gamma'(t) = 0$, and satisfying that:*

$$\text{there exists } t_2 \geq t_0 \text{ and } m < \frac{3}{2} \text{ such that } \gamma(t)\Gamma(t) \leq m \text{ for every } t \geq t_2. \quad (H'_\gamma)$$

Let $b \in]2(m-1), 1[$, then choosing

$$\begin{aligned} a(t) &= \frac{\Gamma(t)(\Gamma(t) - \beta b)}{1 - \beta\gamma(t)}, \\ b(t) &= b, \\ c(t) &= \Gamma(t), \\ d(t) &= b(1 - b), \end{aligned}$$

there exists $\hat{t} > t_0$ such that the system $(S_{a,b,c,d})$ is satisfied for every $t \geq \hat{t}$.

Given $x^* \in \mathcal{S}$ and σ_∞ be such that $t \mapsto \Gamma(t)\sigma_\infty(t) \in L^2([t_0, +\infty[)$, then the following statements hold:

- (i) $\int_{t_0}^{\infty} \Gamma(s) (f(X(s)) - \min f + \|V(s)\|^2) ds < +\infty$ a.s..
- (ii) $\lim_{t \rightarrow \infty} \|\nabla f(X(t))\| + \|V(t)\| = 0$ a.s.,
- (iii) $\mathbb{E}(f(X(t)) - \min f + \|V(t)\|^2) = \mathcal{O}\left(\frac{1}{\Gamma^2(t)}\right)$.

Remark 3.6. When $\gamma(t) = \frac{\alpha}{t}$ with $\alpha > 3$ and $t\sigma_\infty(t) \in L^2([t_0, +\infty[)$, the previous corollary ensures fast convergence of the values, *i.e.*, $\mathcal{O}(t^{-2})$. Besides, by Corollary A.7, when $\gamma(t) = \frac{\alpha}{t^r}$ with $r \in]0, 1[$, $\alpha \geq 1 - r$, and $t^r\sigma_\infty(t) \in L^2([t_0, +\infty[)$, the previous corollary ensures convergence of the objective at a rate $\mathcal{O}(t^{-2r})$. The latter choice indicates that one can require a weaker integrability condition on the noise, compared to the case $\gamma(t) = \frac{\alpha}{t}$ ($\alpha > 3$), but at the price of a slower convergence rate.

Proof. We start by noticing that since γ is decreasing, by [14, Corollary 2.3] we have that $\Gamma(t)$ is increasing and $\gamma(t)\Gamma(t) \geq 1$, for every $t \geq t_0$. Also, it is direct that with a fixed $\beta > 0$ we satisfy (H_β) .

Letting $b \in]2(m-1), 1[$ and $t_1 > t_0$ such that $\beta \leq \frac{1}{\gamma(t_1)}$, this t_1 exists since $t \mapsto \frac{1}{\gamma(t)}$ is an increasing function. We choose $c(t) = \Gamma(t)$, by the fifth equation of $(S_{a,b,c,d})$, we get that $d = b(1 - b)$, and the fourth equation is trivial. The third equation implies that $a(t) = \frac{\Gamma(t)(\Gamma(t) - \beta b)}{1 - \beta\gamma(t)}$ and the choice of β implies that the second equation is satisfied for $t \geq t_1$, since $\beta \leq \frac{1}{\gamma(t_1)} \leq \frac{1}{\gamma(t)} \leq \Gamma(t)$ for every $t > t_1$. By the definition of $c(t)$ and the fact that $b < 1$, we directly have that the sixth equation also holds. We just need to check the first equation, to do that we can see that this equation is equivalent to

$$\frac{\Gamma'(t)(2\Gamma(t) - \beta b)(1 - \beta\gamma(t)) + \beta\Gamma(t)(\Gamma(t) - \beta b)\gamma'(t)}{(1 - \beta\gamma(t))^2} \leq b\Gamma(t),$$

which in turn is equivalent to the following:

$$\begin{aligned} 2\Gamma(t)\Gamma'(t) - \beta b\Gamma'(t) - 2\beta\gamma(t)\Gamma(t)\Gamma'(t) + b\beta^2\gamma(t)\Gamma'(t) + \beta\Gamma^2(t)\gamma'(t) - b\beta^2\Gamma(t)\gamma'(t) \\ \leq b\Gamma(t) - 2b\beta\gamma(t)\Gamma(t) + b\beta^2\gamma^2\Gamma(t). \end{aligned} \quad (3.6)$$

By (H'_γ) , there exists $t_2 > t_0$ such that $\gamma(t)\Gamma(t) \leq m$ and $\Gamma'(t) \leq m - 1$ for every $t \geq t_2$. Note that the terms $b\beta^2\gamma(t)\Gamma'(t)$, $-2b\beta\gamma(t)\Gamma(t)$ are upper and lower bounded by constants. Since the terms

$$-2\beta\gamma(t)\Gamma(t)\Gamma'(t), \quad -\beta b\Gamma'(t), \quad \beta\Gamma^2(t)\gamma'(t)$$

are negative, and $b\beta^2\gamma^2\Gamma(t)$ is positive, if we could prove that there exists $t_3 \geq \max\{t_0, t_1, t_2\}$ such that

$$-b\beta^2\gamma'(t) \leq b - 2\Gamma'(t)$$

for $t \geq t_3$, this would imply that there exists $\hat{t} \geq t_3$ such that (3.6) holds for every $t \geq \hat{t}$. In fact, we see that the previous inequality holds for t large enough (*i.e.* there exists such a t_3) since $\lim_{t \rightarrow \infty} -\gamma'(t) = 0$ and the fact that $2(m-1) < b$ implies that $b - 2\Gamma'(t) > 0$. Thus, we have checked that the proposed a, b, c, d satisfy the system $(S_{a,b,c,d})$ for $t > \hat{t}$.

The rest of the proof is direct from replacing the specified $a, b, c, d, \gamma, \beta$ functions in Proposition 3.4, and the fact that for t large enough, $b\Gamma(t) - a'(t) \geq (b - 2(m-1))\Gamma(t) - C_b$ for some $C_b > 0$, that $\lim_{t \rightarrow \infty} \Gamma(t) = +\infty$, and also $a(t) \geq \Gamma^2(t)$. □

The following result gives us another case in which we can satisfy the system $(S_{a,b,c,d})$. This generalizes to the stochastic setting the results presented in [1, Section 3.1] and [3, Lemma 1]. Besides, it ensures fast convergence of the values whenever $t \mapsto t\sigma_\infty(t) \in L^2([t_0, +\infty[)$.

Corollary 3.7 ($\gamma(t) = \frac{\alpha}{t}$ and $\beta(t) = \gamma_0 + \frac{\beta}{t}$). *Consider the context of Proposition 3.4 in the case where $\gamma(t) = \frac{\alpha}{t}$ and $\beta(t) = \gamma_0 + \frac{\beta}{t}$, where $\alpha > 3, \gamma_0 > 0, \beta \geq 0$. Then choosing*

$$\begin{aligned} a(t) &= t^2 \left(1 + \frac{(\alpha - b)\gamma_0 t - \beta(\alpha + 1 - b)}{t^2 - \alpha\gamma_0 t - \beta(\alpha + 1)} \right), \\ b(t) &= b \in (2, \alpha - 1), \\ c(t) &= t, \\ d(t) &= b(\alpha - 1 - b), \end{aligned}$$

the system $(S_{a,b,c,d})$ is satisfied.

Given $x^* \in \mathcal{S}$ and σ_∞ be such that $t \mapsto t\sigma_\infty(t) \in L^2([t_0, +\infty[)$, then the following statements hold:

- (i) $\int_{t_0}^{\infty} s (f(X(s)) - \min f + \|V(s)\|^2) ds < +\infty$ a.s..
- (ii) $\lim_{t \rightarrow \infty} \|\nabla f(X(t))\| + \|V(t)\| = 0$ a.s..
- (iii) $\mathbb{E}(f(X(t)) - \min f + \|V(t)\|^2) = \mathcal{O}\left(\frac{1}{t^2}\right)$.

Proof. Direct from replacing the specified $a, b, c, d, \gamma, \beta$ functions in Proposition 3.4, and the fact that for t large enough $bt - a'(t) \geq \frac{(\alpha-3)t}{2}$, and also $a(t) \geq t^2, 0 < \gamma_0 < \beta(t) \leq \gamma_0 + \frac{\beta}{t_0}$. □

Remark 3.8. We can use the choices for a, b, c, d presented in Corollaries 3.5 and 3.7 in Propositions A.13 and A.14 to obtain additional integral bounds, almost sure and in expectation properties of $(S - \text{ISIHD})$. We leave this to the reader.

3.3 Strongly convex case

In the following theorem, we consider the case where the objective function is strongly convex and we present a choice of parameters γ and β to obtain a fast linear convergence to a noise dominated region.

Theorem 3.9. *Assume that $f : \mathbb{H} \rightarrow \mathbb{R}$ satisfies (H_0) , and is μ -strongly convex, $\mu > 0$, and denote x^* its unique minimizer. Suppose also that σ obeys (H_σ) . Let $\nu \geq 2$, consider the dynamic $(S - \text{ISIHD})$ with initial data $X_0 \in L^\nu(\Omega; \mathbb{H})$. Consider also $\gamma \equiv 2\sqrt{\mu}$, and a constant β such that $0 \leq \beta \leq \frac{1}{2\sqrt{\mu}}$. Moreover, suppose that σ_∞ is a non-increasing function such that $\sigma_\infty \in L^2([t_0, +\infty[)$. Define the function $\mathcal{E} : [t_0, +\infty[\times \mathbb{H} \times \mathbb{H} \rightarrow \mathbb{R}_+$ as*

$$\mathcal{E}(t, x, v) \stackrel{\text{def}}{=} f(x + \beta v) - \min f + \frac{1}{2} \|\sqrt{\mu}(x - x^*) + v\|^2.$$

Then, $(S - \text{ISIHD})$ has q unique solution $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$. In addition, there exists positive constants M_1, M_2 such that

$$\mathbb{E}[\mathcal{E}(t, X(t), V(t))] \leq \mathcal{E}(t_0, X_0, V_0) e^{-\frac{\sqrt{\mu}}{2}(t-t_0)} + M_1 e^{-\frac{\sqrt{\mu}}{4}(t-t_0)} + M_2 \sigma_\infty \left(\frac{t_0 + t}{2} \right), \quad \forall t > t_0.$$

Let $\Theta : [t_0, +\infty[\rightarrow \mathbb{R}_+$ defined as $\Theta(t) \stackrel{\text{def}}{=} \max\{e^{-\frac{\sqrt{\mu}}{4}(t-t_0)}, \sigma_\infty\left(\frac{t+t_0}{2}\right)\}$. Consequently,

$$\begin{aligned}\mathbb{E}(f(X(t)) - \min f) &= \mathcal{O}(\Theta(t)), \\ \mathbb{E}(\|X(t) - x^*\|^2) &= \mathcal{O}(\Theta(t)), \\ \mathbb{E}(\|V(t)\|^2) &= \mathcal{O}(\Theta(t)), \\ \mathbb{E}(\|\nabla f(X(t))\|^2) &= \mathcal{O}(\Theta(t)).\end{aligned}$$

Proof. Using Itô's formula with \mathcal{E} , taking expectation and denoting $E(t) \stackrel{\text{def}}{=} \mathbb{E}(\mathcal{E}(t, X(t), V(t)))$, we have

$$E(t) \leq E(t_0) - \int_{t_0}^t \frac{\sqrt{\mu}}{2} E(s) ds - \int_{t_0}^t C(s) ds + (L\beta^2 + 1) \int_{t_0}^t \sigma_\infty^2(s) ds,$$

where

$$\begin{aligned}C(t) &\stackrel{\text{def}}{=} \beta \|\nabla f(X(t) + \beta V(t))\|^2 + \beta \sqrt{\mu} \langle \nabla f(X(t) + \beta V(t)), V(t) \rangle + \frac{\sqrt{\mu}}{2} (\beta^2 \mu + 1) \|V(t)\|^2 \\ &\quad + \beta \mu \sqrt{\mu} \langle X(t) - x^*, V(t) \rangle + \frac{\mu \sqrt{\mu}}{4} \|X(t) - x^* + \beta V(t)\|^2.\end{aligned}$$

It was proved in [3, Theorem 4.2] that under the condition $0 \leq \beta \leq \frac{1}{2\sqrt{\mu}}$ we obtain that $C(t)$ is a non-negative function. Therefore, we can write the following

$$E(t) \leq E(t_0) - \int_{t_0}^t \frac{\sqrt{\mu}}{2} E(s) ds + (L\beta^2 + 1) \int_{t_0}^t \sigma_\infty^2(s) ds.$$

We continue by using [13, Lemma A.2], to do this, we need to solve the following Cauchy problem:

$$\begin{cases} Y'(t) &= -\frac{\sqrt{\mu}}{2} Y(t) + (L\beta^2 + 1) \sigma_\infty^2(t) \\ Y(t_0) &= \mathcal{E}(t_0, X_0, V_0). \end{cases}$$

Using the integrating factor method, we deduce that for all $t \geq t_0$:

$$\begin{aligned}Y(t) &= Y(t_0) e^{\frac{\sqrt{\mu}}{2}(t_0-t)} + (L\beta^2 + 1) e^{-\frac{\sqrt{\mu}}{2}t} \int_{t_0}^t e^{\frac{\sqrt{\mu}}{2}s} \sigma_\infty^2(s) ds \\ &\leq Y(t_0) e^{\frac{\sqrt{\mu}}{2}(t_0-t)} + (L\beta^2 + 1) e^{-\frac{\sqrt{\mu}}{2}t} \left(\int_{t_0}^{\frac{t_0+t}{2}} e^{\frac{\sqrt{\mu}}{2}s} \sigma_\infty^2(s) ds + \int_{\frac{t_0+t}{2}}^t e^{\frac{\sqrt{\mu}}{2}s} \sigma_\infty^2(s) ds \right) \\ &\leq Y(t_0) e^{\frac{\sqrt{\mu}}{2}(t_0-t)} + (L\beta^2 + 1) \sigma_\infty^2\left(\frac{t_0+t}{2}\right) + (L\beta^2 + 1) \sigma_\infty^2(t_0) e^{\frac{\sqrt{\mu}}{4}(t_0-t)} \\ &= \mathcal{O}(\Theta(t)).\end{aligned}$$

By [13, Lemma A.2], we conclude that $E(t) = \mathcal{O}(\Theta(t))$, immediately we observe that

$$\begin{aligned}\mathbb{E}(f(X(t) + \beta V(t)) - \min f) &= \mathcal{O}(\Theta(t)) \\ \mathbb{E}(\|\sqrt{\mu}(X(t) - x^*) + V(t)\|^2) &= \mathcal{O}(\Theta(t))\end{aligned}$$

By the strong convexity of f , we have that $\mathbb{E}(\|X(t) - x^* + \beta V(t)\|^2) = \mathcal{O}(\Theta(t))$, since $\beta \neq \frac{1}{\sqrt{\mu}}$ ($\beta \leq \frac{1}{2\sqrt{\mu}}$), then $\mathbb{E}(\|X(t) - x^*\|^2) = \mathbb{E}(\|V(t)\|^2) = \mathcal{O}(\Theta(t))$, on the other hand, using Lemma 2.1 and Lemma 2.2, $\mathbb{E}(f(X(t)) - f(X(t) + \beta V(t))) = \mathcal{O}(\Theta(t))$, thus,

$$\mathbb{E}(f(X(t)) - \min f) = \mathbb{E}(\|\nabla f(X(t))\|^2) = \mathcal{O}(\Theta(t)).$$

□

4 (S – ISIH) with general γ and $\beta \equiv 0$

In this section we are going to study properties of the dynamic (S – ISIH) in expectation and in almost sure sense, when the parameter γ is a general function and $\beta \equiv 0$. The noiseless case and under deterministic noise is well documented in [19].

Consider the dynamic (S – ISIH) when $\beta \equiv 0$. This dynamic will be a stochastic version of the Hamiltonian formulation of (IGS $_\gamma$) and it will be described by:

$$\begin{cases} dX(t) &= V(t)dt, \quad t > t_0, \\ dV(t) &= -\gamma(t)V(t)dt - \nabla f(X(t))dt + \sigma(t, X(t))dW(t), \quad t > t_0, \\ X(t_0) &= X_0, \quad V(t_0) = V_0. \end{cases} \quad (\text{IGS}_\gamma - \text{S})$$

The main motivation for a separate analysis is that, in Section 3 we consider hypothesis (H $_\beta$) to establish the existence and uniqueness of a solution, from which, the rest of the results follow. This hypothesis is incompatible with the case $\beta \equiv 0$.

We will demonstrate almost sure convergence of the velocity to zero and of the objective to its minimum value, under assumptions that are satisfied for $\gamma(t) = \frac{\alpha}{t^r}$, with $r \in [0, 1], \alpha \geq 1 - r$. Additionally, we will show that for this particular choice of β , we can obtain almost sure (weak) convergence of the trajectory.

4.1 Minimization properties

Let us define for $c > 0$,

$$\lambda_c(t) = \frac{p(t)}{c + \int_{t_0}^t p(s)ds}.$$

We can deduce that $\lambda'_c + \lambda_c^2 - \gamma\lambda_c = 0$, besides, since $p \notin L^1([t_0, +\infty[)$, then $\lambda_c \notin L^1([t_0, +\infty[)$.

Theorem 4.1. *Assume that f and σ satisfy assumptions (H $_0$) and (H $_\sigma$), respectively. Let $\nu \geq 2$, and consider the dynamic (IGS $_\gamma - \text{S}$) with initial data $X_0, V_0 \in L^\nu(\Omega; \mathbb{H})$. Then, there exists a unique solution $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$ of (IGS $_\gamma - \text{S}$). Additionally, if $\sigma_\infty \in L^2([t_0, +\infty[)$, then*

$$\int_{t_0}^{\infty} \gamma(s)\|V(s)\|^2 ds < +\infty \quad a.s..$$

Moreover, suppose that

$$\text{there exists } \hat{t} \geq t_0, \text{ and } c > 0 \text{ such that } \gamma(t) \leq \lambda_c(t) \quad \forall t \geq \hat{t}, \quad (\text{H}_a)$$

and

$$\int_{t_0}^{\infty} \lambda_c(s)\|V(s)\|^2 ds < +\infty \quad a.s.. \quad (\text{H}_b)$$

Then the following properties are satisfied:

- (i) $\int_{t_0}^{\infty} \lambda_c(s)(f(X(s)) - \min f)ds < +\infty$ a.s..
- (ii) $\lim_{t \rightarrow \infty} \|V(t)\| = 0$ a.s. and $\lim_{t \rightarrow \infty} f(X(t)) - \min f = 0$ a.s..

Proof. The existence and uniqueness of a solution of (IGS $_\gamma - \text{S}$) is a direct consequence of [13, Theorem 3.3] in the product space $\mathbb{H} \times \mathbb{H}$.

Let $x^* \in \mathcal{S}$ and $\phi_0 : (x, v) \mapsto \mathbb{R}$ defined by $\phi_0(x, v) = f(x) - \min f + \frac{\|v\|^2}{2}$, by Itô's formula and Theorem A.12 we obtain that $\int_{t_0}^{\infty} \gamma(s)\|V(s)\|^2 ds < +\infty$ a.s.. Moreover, if we assume the hypotheses (H $_a$) and (H $_b$), then:

- (i) Let $x^* \in \mathcal{S}$ and $\phi : (t, x, v) \mapsto \mathbb{R}$ defined by $\phi(t, x, v) = \frac{\|\lambda_c(t)(x-x^*)+v\|^2}{2} + (f(x) - \min f)$. Let \hat{t} defined in the statement, by Itô's formula from \hat{t} to t , we have

$$\begin{aligned} f(X(t)) - \min f + \frac{\|\lambda_c(t)(X(t) - x^*) + V(t)\|^2}{2} &= f(X(\hat{t})) - \min f + \frac{\|\lambda_c(t)(X(\hat{t}) - x^*) + V(\hat{t})\|^2}{2} \\ &+ \int_{\hat{t}}^t \lambda_c(s) \lambda'_c(s) \|X(s) - x^*\|^2 - \gamma(t) \|V(s)\|^2 - \lambda_c(t) \langle \nabla f(X(s)), X(s) - x^* \rangle ds \\ &+ \int_{\hat{t}}^t \lambda_c(s) \|V(s)\|^2 + \text{tr}[\Sigma(s, X(s))] ds + \underbrace{\int_{\hat{t}}^t \langle [\lambda_c(s)(X(s) - x^*) + V(s)] \sigma^*(s, X(s)), dW(s) \rangle}_{M_t}. \end{aligned}$$

By the hypotheses, we have that

$$\int_{\hat{t}}^{\infty} (\lambda_c(s) \|V(s)\|^2 + \text{tr}[\Sigma(s, X(s))]) ds \leq \int_{\hat{t}}^{\infty} (\lambda_c(s) \|V(s)\|^2 + \sigma_{\infty}^2(s)) ds < +\infty \quad a.s..$$

Besides $(M_t)_{t \geq \hat{t}}$ is a continuous martingale. Moreover, by convexity of f and the fact that $\lambda'_c(t) \leq 0 \forall t \geq \hat{t}$,

$$\begin{aligned} &\int_{\hat{t}}^t \lambda_c(s) \lambda'_c(s) \|X(s) - x^*\|^2 - \gamma(t) \|V(s)\|^2 - \lambda_c(t) \langle \nabla f(X(s)), X(s) - x^* \rangle ds \\ &\leq - \int_{\hat{t}}^t \lambda_c(s) (f(X(s)) - \min f) ds. \end{aligned}$$

Then, by Theorem A.12,

$$\int_{\hat{t}}^{\infty} \lambda_c(s) (f(X(s)) - \min f) ds < +\infty \quad a.s., \quad (4.1)$$

and $\lim_{t \rightarrow \infty} f(X(t)) - \min f + \frac{\|\lambda_c(t)(X(t) - x^*) + V(t)\|^2}{2}$ exists a.s..

- (ii) By Lemma A.2 and (4.1), we conclude that $\lim_{t \rightarrow \infty} \frac{\|V(t)\|^2}{2} + f(X(t)) - \min f = 0$ a.s. \square

Corollary 4.2. Consider the context of Theorem 4.1 with $\gamma(t) = \frac{\alpha}{t^r}$, where $r \in [0, 1]$ and $\alpha > 1 - r$. Then (H_a) and (H_b) are satisfied and thus the conclusions of Theorem 4.1 hold.

Proof. • We will prove the case $r = 1$ first, since it is direct, in such case, letting $c = \frac{t_0}{\alpha+1}$ we have that $\lambda_c(t) = \frac{\alpha+1}{t}$, which satisfies (H_a) , moreover, $\lambda(t) = \frac{\alpha+1}{\alpha} \gamma(t)$, so (H_b) is also satisfied.

- Let $r \in]0, 1[$, $c = \frac{\int_0^{t_0} e^{\alpha s^{1-r}} ds}{e^{\alpha t_0^{1-r}}}$. Instead of proving $\gamma(t) \leq \lambda_c(t)$, we will prove the equivalent inequality $\frac{1}{t\lambda_c(t)} \leq \frac{1}{t\gamma(t)}$. In fact, by a change of variable we have that (see notation of I_p in Lemma A.8):

$$\frac{1}{\lambda_c(t)t} = \frac{(1-r)^{\frac{r}{1-r}}}{\alpha^{\frac{1}{1-r}}} I_{\frac{r}{1-r}} \left(\frac{\alpha}{1-r} t^{1-r} \right),$$

Moreover, by the first result of Lemma A.8 we have that

$$\frac{1}{t\lambda_c(t)} \leq \left(\frac{1-r}{\alpha} \right)^{\frac{1}{1-r}} \frac{t^{r-1}}{\alpha} \leq \frac{1}{t\gamma(t)}.$$

where the last inequality comes from the fact that $1-r \leq \alpha$. Moreover, by the second result of Lemma A.8, we obtain that:

$$\left(\frac{\alpha}{1-r} \right)^{\frac{1}{1-r}} \frac{1}{t\lambda_c(t)} \sim \frac{1}{t\gamma(t)}, \quad \text{as } t \rightarrow +\infty.$$

This implies that for every $\varepsilon \in]0, 1[$ there exists $\hat{t} > t_0, \Lambda_\varepsilon \geq 1$ such that $\lambda_c(t) \leq \Lambda_\varepsilon \gamma(t)$ for every $t > \hat{t}$ ($\Lambda_\varepsilon = \left(\frac{\alpha}{1-r}\right)^{\frac{1}{1-r}} \frac{1}{(1-\varepsilon)}$), this implies (H_b) . □

Remark 4.3. Finding all (or at least a larger class of) continuous functions γ that satisfy (H_a) and for which one can prove (H_b) in general is an open problem.

4.2 Fast convergence of the values

In order to illustrate the context of the following result, it is useful to mention that if $\gamma(t) = \frac{\alpha}{t}$, then Theorem 4.1 gives us minimization properties in the case $\alpha > 0$. However, it is widely known in the continuous deterministic setting (1.3) that if $\alpha > 3$, then the values converge at the rate $o(1/t^2)$ (see [19, 20]). Based on [19], we will depict that effect for a general γ in the continuous stochastic setting.

We will rephrase assumption (H_0) on the objective f to:

$$\begin{cases} f \text{ is convex and continuously differentiable with } L\text{-Lipschitz continuous gradient;} \\ f \in C^2(\mathbb{H}) \text{ or } \mathbb{H} \text{ is finite-dimensional;} \\ \mathcal{S} \stackrel{\text{def}}{=} \text{argmin}(f) \neq \emptyset. \end{cases} \quad (H_0^*)$$

(H_0^*) coincides with (H_0) in the infinite-dimensional case, but is weaker than (H_0) when \mathbb{H} is finite-dimensional.

Theorem 4.4. Assume that f, σ and γ satisfy assumptions $(H_0^*), (H_\sigma)$ and (H_γ) - (H'_γ) , respectively. Let $\nu \geq 2$, and consider the dynamic $(IGS_\gamma - S)$ with initial data $X_0, V_0 \in L^\nu(\Omega; \mathbb{H})$. Then, there exists a unique solution $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$ of $(IGS_\gamma - S)$, for every $\nu \geq 2$. Additionally, if $\Gamma \sigma_\infty \in L^2([t_0, +\infty[)$, then:

(i) $\int_{t_2}^\infty \Gamma(t)(f(X(t)) - \min f + \|V(t)\|^2) dt < +\infty$ a.s..

(ii) $f(X(t)) - \min f + \|V(t)\|^2 = o\left(\frac{1}{\Gamma^2(t)}\right)$ a.s..

(iii) $\mathbb{E}(f(X(t)) - \min f + \|V(t)\|^2) = \mathcal{O}\left(\frac{1}{\Gamma^2(t)}\right)$.

Moreover, assume that $\Gamma \notin L^1([t_0, +\infty[)$, and let $\theta(t) \stackrel{\text{def}}{=} \int_{t_0}^t \Gamma(s) ds$. If also $\theta \sigma_\infty^2 \in L^1([t_0, +\infty[)$, then:

(iv) $f(X(t)) - \min f + \|V(t)\|^2 = o\left(\frac{1}{\theta(t)}\right)$ a.s..

(v) $\mathbb{E}(f(X(t)) - \min f + \|V(t)\|^2) = \mathcal{O}\left(\min\left\{\frac{\int_{t_0}^t \Gamma(s) ds}{\theta(t)}, \frac{1}{\Gamma^2(t)}\right\}\right)$.

Remark 4.5. The claim (ii) is new even in the deterministic case. According to the first three items of the previous theorem, the conclusions of Remark 3.6 are also valid. Regarding the fourth item, this can be seen as the extension of [19, Theorem 3.6] to the stochastic setting.

Proof. (i) Let $m < \frac{3}{2}$ and t_2 defined in (H'_γ) , let also $b \in]2(m-1), 1[$ and $x^* \in \mathcal{S}$. Based on $(S_{a,b,c,d})$ with $\beta \equiv 0$, we introduce $\phi_1 : (t, x, v) \mapsto \mathbb{R}$ defined by

$$\phi_1(t, x, v) = \Gamma^2(t)(f(x) - \min f) + \frac{\|b(x - x^*) + \Gamma(t)v\|^2}{2} + \frac{b(1-b)}{2} \|x - x^*\|^2.$$

Since $f \in C^2(\mathbb{H})$, we use Itô's formula from t_2 to t to get

$$\begin{aligned} \phi_1(t, X(t), V(t)) &= \phi_1(t_2, X(t_2), V(t_2)) + \int_{t_2}^t \Gamma(s)[2\Gamma'(s)(f(X(s)) - \min f) - b\langle \nabla f(X(s)), X(s) - x^* \rangle] ds \\ &\quad + (b-1) \int_{t_2}^t \Gamma(s) \|V(s)\|^2 ds + \int_{t_2}^t \Gamma^2(s) \text{tr}[\Sigma(s, X(s))] ds \\ &\quad + \underbrace{\int_{t_2}^t \langle [\Gamma^2(s)V(s) + b\Gamma(s)(X(s) - x^*)] \sigma^*(s, X(s)), dW(s) \rangle}_{M_t}. \end{aligned} \quad (4.2)$$

When \mathbb{H} is finite-dimensional but f is not $C^2(\mathbb{H})$, we can use mollifiers as in [10, Proposition C.2], and get (4.2) as an inequality in this case.

Besides, we have that

$$\int_{t_2}^{\infty} \Gamma^2(s) \text{tr}[\Sigma(s, X(s))] ds \leq \int_{t_2}^{\infty} \Gamma^2(s) \sigma_{\infty}^2(s) ds < +\infty.$$

Besides $(M_t)_{t \geq t_2}$ is a continuous martingale. Moreover, by convexity of f , we have that

$$\begin{aligned} & \int_{t_2}^t \Gamma(s) [2\Gamma'(s)(f(X(s)) - \min f) - b \langle \nabla f(X(s)), X(s) - x^* \rangle] ds \\ & \leq \int_{t_2}^t \Gamma(s) (2\Gamma'(s) - b)(f(X(s)) - \min f) ds. \end{aligned}$$

Since $b - 1 < 0$, and

$$2\Gamma'(t) - b = 2\gamma(t)\Gamma(t) - 2 - b \leq 2(m - 1) - b < 0, \quad \forall t > t_2.$$

By Theorem A.12,

$$\int_{t_2}^{\infty} \Gamma(s) (f(X(s)) - \min f + \|V(s)\|^2) ds < +\infty \quad a.s., \quad (4.3)$$

and

$$\lim_{t \rightarrow \infty} \Gamma^2(t) (f(X(t)) - \min f) + \frac{\|b(X(t) - x^*) + \Gamma(t)V(t)\|^2}{2} + \frac{b(1-b)}{2} \|X(t) - x^*\|^2 \text{ exists a.s..}$$

(ii) On the other hand, let $\phi_2 : (t, x, v) \mapsto \mathbb{R}$ defined by $\phi_2(t, x, v) = \Gamma^2(t) \left(f(x) - \min f + \frac{\|v\|^2}{2} \right)$. Recalling the discussion for ϕ_1 , we get that by Itô's formula from t_2 to t , we have

$$\begin{aligned} \phi_2(t, X(t), V(t)) &= \phi_2(t_2, X(t_2), V(t_2)) + \int_{t_2}^t 2\Gamma(s)\Gamma'(s)(f(X(s)) - \min f) ds \\ &\quad - \int_{t_2}^t \Gamma(s) \|V(s)\|^2 ds + \int_{t_2}^t \Gamma^2(s) \text{tr}[\Sigma(s, X(s))] ds \\ &\quad + \underbrace{\int_{t_2}^t \Gamma^2(s) \langle V(s) \sigma^*(s, X(s)), dW(s) \rangle}_{M_t}. \end{aligned} \quad (4.4)$$

And also, that

$$\begin{aligned} & \int_{t_2}^{\infty} 2\Gamma(s)\Gamma'(s)(f(X(s)) - \min f) + \Gamma^2(s) \text{tr}[\Sigma(s, X(s))] ds \\ & \leq \int_{t_2}^{\infty} \Gamma(s) (f(X(s)) - \min f) + \Gamma^2(s) \sigma_{\infty}^2(s) ds < +\infty \quad a.s.. \end{aligned}$$

Besides $(M_t)_{t \geq t_2}$ is a continuous martingale. By Theorem A.12, we get again that $\int_{t_2}^{\infty} \Gamma(s) \|V(s)\|^2 ds < +\infty$ a.s. and that

$$\lim_{t \rightarrow \infty} \Gamma^2(t) \left(f(X(t)) - \min f + \frac{\|V(t)\|^2}{2} \right) \text{ exists a.s.} \quad (4.5)$$

Let us recall that $\frac{1}{\Gamma} \notin L^1([t_0, +\infty[)$ by Lemma A.3. Therefore, by (4.3) and (4.5), we can use Lemma A.2 to obtain that

$$\lim_{t \rightarrow \infty} \Gamma^2(t) \left(f(X(t)) - \min f + \frac{\|V(t)\|^2}{2} \right) = 0 \quad a.s..$$

(iii) Taking expectation on (4.2) and denoting

$$K_1 \stackrel{\text{def}}{=} \Gamma^2(t_2) \mathbb{E}((f(X(t_2)) - \min f)) + \frac{1}{2} \mathbb{E}(\|b(X(t_2) - x^*) + \Gamma(t_2)V(t_2)\|^2) + \frac{b(1-b)}{2} \|X(t_2) - x^*\|^2,$$

$$K_\Gamma \stackrel{\text{def}}{=} \int_{t_2}^{\infty} \Gamma^2(s) \sigma_\infty^2(s) ds,$$

we obtain directly that

$$\mathbb{E} \left(\Gamma^2(t) (f(X(t)) - \min f) + \frac{\|b(X(t) - x^*) + \Gamma(t)V(t)\|^2}{2} + \frac{b(1-b)}{2} \|X(t) - x^*\|^2 \right) \leq K_1 + K_\Gamma.$$

From this, is direct that $\sup_{t \geq t_2} \mathbb{E}(\|X(t) - x^*\|^2) < +\infty$, and this in turn imply

$$\mathbb{E} \left(f(X(t)) - \min f + \frac{\|V(t)\|^2}{2} \right) = \mathcal{O} \left(\frac{1}{\Gamma^2(t)} \right).$$

(iv) Moreover, assume that $\Gamma \notin L^1([t_0, +\infty[)$, and let $\theta(t) = \int_{t_0}^t \Gamma(s)$. If also $\theta \sigma_\infty^2 \in L^1([t_0, +\infty[)$, then we consider $\phi_3(t, x, v) = \theta(t) \left(f(x) - \min f + \frac{\|v\|^2}{2} \right)$, by Itô's formula from t_2 to t , we get

$$\begin{aligned} \phi_3(t, X(t), V(t)) &= \phi_3(t_2, X(t_2), V(t_2)) + \int_{t_2}^t \Gamma(s) \left(f(X(s)) - \min f + \frac{\|V(s)\|^2}{2} \right) ds \\ &\quad - \int_{t_2}^t \gamma(s) \theta(s) \|V(s)\|^2 ds + \frac{1}{2} \int_{t_2}^t \theta(s) \text{tr}[\Sigma(s, X(s))] ds \\ &\quad + \underbrace{\int_{t_2}^t \theta(s) \langle V(s) \sigma^*(s, X(s)), dW(s) \rangle}_{M_t}. \end{aligned} \tag{4.6}$$

Also, by the first item and new hypothesis on the diffusion term, we get that

$$\begin{aligned} &\int_{t_2}^t \Gamma(s) \left(f(X(s)) - \min f + \frac{\|V(s)\|^2}{2} \right) + \theta(s) \text{tr}[\Sigma(s, X(s))] ds \\ &\leq \int_{t_2}^{\infty} \Gamma(s) \left(f(X(s)) - \min f + \frac{\|V(s)\|^2}{2} \right) + \theta(s) \sigma_\infty^2(s) ds < +\infty. \end{aligned} \tag{4.7}$$

Besides $(M_t)_{t \geq t_2}$ is a continuous martingale. By Theorem A.12, we get that $\int_{t_2}^{\infty} \gamma(s) \theta(s) \|V(s)\|^2 ds < +\infty$ a.s. and that

$$\lim_{t \rightarrow \infty} \theta(t) \left(f(X(t)) - \min f + \frac{\|V(t)\|^2}{2} \right) \text{ exists a.s.}, \tag{4.8}$$

Using Lemma A.4 with $q(t) = \theta(t)$, we get that $\frac{\Gamma}{\theta} \notin L^1([t_2, +\infty[)$. Besides, recalling that

$$\int_{t_2}^{\infty} \Gamma(s) \left(f(X(s)) - \min f + \frac{\|V(s)\|^2}{2} \right) < +\infty, \quad a.s.,$$

we invoke Lemma A.2 to conclude that $\lim_{t \rightarrow \infty} \theta(t) \left(f(X(t)) - \min f + \frac{\|V(t)\|^2}{2} \right) = 0$ a.s..

(v) Taking expectation in (4.6) and upper bounding we get

$$\begin{aligned} \mathbb{E}(\phi_3(t, X(t), V(t))) &\leq \mathbb{E}(\phi_3(t_2, X(t_2), V(t_2))) + \int_{t_2}^t \Gamma(s) \mathbb{E} \left(f(X(s)) - \min f + \frac{\|V(s)\|^2}{2} \right) ds \\ &\quad + \frac{1}{2} \int_{t_2}^{\infty} \theta(s) \sigma_\infty^2(s) ds. \end{aligned} \tag{4.9}$$

By the third item, we have that $\mathbb{E} \left(f(X(s)) - \min f + \frac{\|V(s)\|^2}{2} \right) = \mathcal{O} \left(\frac{1}{\Gamma^2(s)} \right)$, so we conclude that

$$\mathbb{E} \left(\theta(t) \left(f(X(t)) - \min f + \frac{\|V(t)\|^2}{2} \right) \right) = \mathcal{O} \left(\int_{t_2}^t \frac{ds}{\Gamma(s)} \right). \quad (4.10)$$

Thus,

$$\mathbb{E} \left(f(X(t)) - \min f + \frac{\|V(t)\|^2}{2} \right) = \mathcal{O} \left(\min \left\{ \frac{\int_{t_2}^t \frac{ds}{\Gamma(s)}}{\theta(t)}, \frac{1}{\Gamma^2(t)} \right\} \right).$$

□

4.3 Almost sure weak convergence of trajectories

In the deterministic setting with $\alpha > 3$, it is also well-known that one can obtain weak convergence of the trajectory. Our aim in this section is to show this claim for a general γ in the stochastic setting.

Theorem 4.6. *Consider the setting of Theorem 4.4. Then, if $\Gamma\sigma_\infty \in L^2([t_0, +\infty[)$ we have that:*

- (i) $\mathbb{E} \left[\sup_{t \geq t_2} \|X(t)\|^\nu \right] < +\infty$.
- (ii) $\forall x^* \in \mathcal{S}$, $\lim_{t \rightarrow \infty} \|X(t) - x^*\|$ exists a.s..
- (iii) If γ is non-increasing, there exists an \mathcal{S} -valued random variable X^* such that $w\text{-}\lim_{t \rightarrow \infty} X(t) = X^*$ a.s..

Proof. (i) Analogous to the proof of the first point of [11, Theorem 3.1].

(ii) Recalling the proof of Theorem 4.4, we combine the fact that both

$$\lim_{t \rightarrow \infty} \Gamma^2(t) (f(X(t)) - \min f) + \frac{\|b(X(t) - x^*) + \Gamma(t)V(t)\|^2}{2} + \frac{b(1-b)}{2} \|X(t) - x^*\|^2,$$

and

$$\lim_{t \rightarrow \infty} \Gamma^2(t) \left(f(X(t)) - \min f + \frac{\|V(t)\|^2}{2} \right)$$

exist a.s.. We can subtract both quantities to obtain that

$$\lim_{t \rightarrow \infty} \frac{\|X(t) - x^*\|^2}{2} + \Gamma(t) \langle X(t) - x^*, V(t) \rangle \text{ exists a.s..}$$

Thus, for every $x^* \in \mathcal{S}$ there exists $\Omega_{x^*} \in \mathcal{F}$ with $\mathbb{P}(\Omega_{x^*}) = 1$ and $\exists \ell : \Omega_{x^*} \mapsto \mathbb{R}$ such that

$$\lim_{t \rightarrow \infty} \frac{\|X(\omega, t) - x^*\|^2}{2} + \Gamma(t) \langle V(\omega, t), X(\omega, t) - x^* \rangle = \ell(\omega).$$

Let $Z(\omega, t) = \frac{\|X(\omega, t) - x^*\|^2}{2} - \ell(\omega)$ and $\varepsilon > 0$ arbitrary. There exists $T(\omega) \geq t_0$ such that $\forall t \geq T(\omega)$

$$\left\| Z(\omega, t) + \Gamma(t) \langle V(\omega, t), X(\omega, t) - x^* \rangle \right\| < \varepsilon.$$

Let $g(t) \stackrel{\text{def}}{=} \exp \left(\int_{t_2}^t \frac{ds}{\Gamma(s)} \right)$, multiplying the previous inequality by $\frac{g(t)}{\Gamma(t)}$, there exists $T(\omega) \geq t_0$ such that for every $t \geq T(\omega)$:

$$\left\| \frac{g(t)}{\Gamma(t)} Z(t) + g(t) \langle V(\omega, t), X(\omega, t) - x^* \rangle \right\| < \frac{\varepsilon}{\Gamma(t)} g(t).$$

On the other hand, $dZ(t) = \langle V(t), X(t) - x^* \rangle dt$ and

$$d(g(t)Z(t)) = \left(\frac{g(t)}{\Gamma(t)} Z(t) + g(t) \langle V(t), X(t) - x^* \rangle \right) dt.$$

Thus,

$$\begin{aligned} \|g(t)Z(t) - g(T)z(T)\| &= \left\| \int_T^t d(g(s)Z(s)) \right\| = \left\| \int_T^t \left(\frac{g(s)}{\Gamma(s)} Z(s) + g(s)\langle V(s), X(s) - x^* \rangle \right) ds \right\| \\ &\leq \varepsilon \int_T^t \frac{g(s)}{\Gamma(s)} ds = \varepsilon(g(t) - g(T)). \end{aligned}$$

So,

$$\|Z(t)\| \leq \frac{g(T)}{g(t)} \|z(T)\| + \varepsilon.$$

By Lemma A.3, we obtain that $\lim_{t \rightarrow \infty} g(t) = +\infty$. Hence, $\limsup_{t \rightarrow \infty} \|Z(t)\| \leq \varepsilon$. And we conclude that for every $x^* \in \mathcal{S}$, $\lim_{t \rightarrow \infty} \frac{\|X(t) - x^*\|}{2}$ exists a.s.. By a separability argument (see proof of [11, Theorem 3.1] or [13, Theorem 3.6]) there exists $\tilde{\Omega} \in \mathcal{F}$ (independent of x^*) such that $\mathbb{P}(\tilde{\Omega}) = 1$ and $\lim_{t \rightarrow \infty} \frac{\|X(\omega, t) - x^*\|}{2}$ exists for every $\omega \in \tilde{\Omega}$, $x^* \in \mathcal{S}$.

(iii) If γ is non-increasing, then Γ is non-decreasing (see [19, Corollary 2.3]). Then, by item (ii) of Theorem 4.4, we have that:

$$\lim_{t \rightarrow +\infty} f(X(t)) = \min f \quad a.s..$$

Let $\Omega_f \in \mathcal{F}$ be the set of events on which this limit is satisfied. Thus $\mathbb{P}(\Omega_f) = 1$. Set $\Omega_{\text{conv}} \stackrel{\text{def}}{=} \Omega_f \cap \tilde{\Omega}$. We have $\mathbb{P}(\Omega_{\text{conv}}) = 1$. Now, let $\omega \in \Omega_{\text{conv}}$ and $\tilde{X}(\omega)$ be a weak sequential cluster point of $X(\omega, t)$ (which exists by boundedness). 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

$$\text{w-lim}_{k \rightarrow \infty} X(\omega, t_k) = \tilde{X}(\omega).$$

Since $\lim_{t \rightarrow \infty} f(X(\omega, t)) = \min f$ and the fact that f is weakly lower semicontinuous (since it is convex and continuous), we obtain directly that $\tilde{X}(\omega) \in \mathcal{S}$. Finally by Opial's Lemma (see [56]) we conclude that there exists $X^*(\omega) \in \mathcal{S}$ such that $\text{w-lim}_{t \rightarrow \infty} X(\omega, t) = X^*(\omega)$. In other words, since $\omega \in \Omega_{\text{conv}}$ was arbitrary, there exists an \mathcal{S} -valued random variable X^* such that $\text{w-lim}_{t \rightarrow \infty} X(t) = X^*$ a.s.. □

A Auxiliary results

A.1 Deterministic results

Lemma A.1. *Let $a, b \in \mathbb{R}$ and $x, y \in \mathbb{H}$, then*

$$\|ax - by\| \leq \max\{|a|, |b|\} \|x - y\| + |a - b| \max\{\|x\|, \|y\|\}.$$

Lemma A.2. *Let $t_0 > 0$ and $a, b : [t_0, +\infty[\rightarrow \mathbb{R}_+$. If $\lim_{t \rightarrow \infty} a(t)$ exists, $b \notin L^1([t_0, +\infty[)$ and $\int_{t_0}^{\infty} a(s)b(s)ds < +\infty$, then $\lim_{t \rightarrow \infty} a(t) = 0$.*

Lemma A.3. *Under hypothesis (H_γ) , then*

$$\int_{t_0}^{\infty} \frac{ds}{\Gamma(s)} = +\infty.$$

Proof. Let $q(t) \stackrel{\text{def}}{=} \int_t^{\infty} \frac{ds}{p(s)}$, since $\int_{t_0}^{\infty} \frac{ds}{p(s)} < +\infty$, then $\lim_{t \rightarrow \infty} q(t) = 0$ and $q'(t) = -\frac{1}{p(t)}$. On the other hand

$$\int_{t_0}^{\infty} \frac{ds}{\Gamma(s)} = - \int_{t_0}^{\infty} \frac{q'(s)}{q(s)} ds = \ln(q(t_0)) - \lim_{t \rightarrow \infty} \ln(q(t)) = +\infty.$$

□

Lemma A.4. Let $q : [t_0, +\infty[\rightarrow \mathbb{R}_+$ be a non-decreasing differentiable function, if $q \notin L^1([t_0, +\infty[)$, then $\frac{q'}{q} \notin L^1([t_0, +\infty[)$

Proof. By definition,

$$\int_{t_0}^{\infty} \frac{q'(s)}{q(s)} ds = \lim_{t \rightarrow \infty} \ln(q(t)) - \ln(q(t_0)) = +\infty.$$

□

Lemma A.5. For $a, x > 0$, let us define the upper incomplete Gamma function as:

$$\Gamma_{inc}(a; x) = \int_x^{\infty} s^{a-1} e^{-s} ds.$$

Then, the following holds:

- (i) $x^{1-a} e^x \Gamma_{inc}(a; x) \leq 1$ for $0 < a \leq 1$.
- (ii) $x^{1-a} e^x \Gamma_{inc}(a; x) \geq 1$ for $a \geq 1$.
- (iii) $\lim_{x \rightarrow \infty} x^{1-a} e^x \Gamma_{inc}(a; x) = 1$

Proof. See [57, Section 8].

□

Remark A.6. Do not confuse $\Gamma_{inc}(a; x)$ with $\Gamma(t)$ defined in (3.2).

Corollary A.7. Let us consider the viscous damping function $\gamma : [t_0, +\infty[\rightarrow \mathbb{R}_+$ defined by $\gamma(t) = \frac{\alpha}{t^r}$ with $r \in]0, 1[$ and $\alpha \geq 1 - r$, then:

- (i) γ satisfies (H_γ) .
- (ii) $\Gamma(t) = \mathcal{O}(t^r)$.
- (iii) γ satisfies (H'_γ) .

Proof. (i) Let $c \stackrel{\text{def}}{=} \frac{\alpha}{1-r} \geq 1$, we first notice that after the change of variable $u = cs^{1-r}$, we get

$$\int_0^{\infty} \exp(-cs^{1-r}) ds = \frac{1}{\alpha c^{\frac{r}{1-r}}} \int_0^{\infty} u^{\frac{1}{1-r}-1} e^{-u} du < +\infty,$$

since the last integral is the classical Gamma function (see e.g. [57, Section 5]) evaluated at $\frac{1}{1-r}$, and this function is well defined for positive arguments, then (H_γ) is satisfied.

- (ii) Besides, by definition $\Gamma(t) = \exp(ct^{1-r}) \int_t^{\infty} \exp(-cs^{1-r}) ds$. Using the same change of variable as before, we obtain that

$$\Gamma(t) = \frac{\exp(ct^{1-r})}{\alpha c^{\frac{r}{1-r}}} \Gamma_{inc} \left(\frac{1}{1-r}, ct^{1-r} \right). \quad (\text{A.1})$$

By (iii) of Lemma A.5 with $a = \frac{1}{1-r} > 1$ and $x = ct^{1-r}$, for every $\varepsilon > 0$, there exists $t_1 > t_0$ such that for every $t > t_1$:

$$\Gamma(t) \leq \frac{1+\varepsilon}{\alpha c^{\frac{1}{1-r}}} t^r.$$

- (iii) Moreover, if we restrict $\varepsilon \in]0, \frac{1}{2}[$, there exists $t_1 > t_0$ such that for every $t > t_1$:

$$\gamma(t)\Gamma(t) \leq \frac{1+\varepsilon}{c^{\frac{1}{1-r}}} \leq 1 + \varepsilon.$$

Defining m as $1 + \varepsilon$, we have that $m < \frac{3}{2}$, and we conclude.

□

Lemma A.8. Let us define $p > 0$ and $I_p(t) \stackrel{\text{def}}{=} \int_0^1 e^{-tu} (1-u)^p du$. Then

- (i) $I_p(t) \leq t^{-1}$ for every $t > 0$.
- (ii) $I_p(t) \sim t^{-1}$ as $t \rightarrow +\infty$.

Proof. The first result comes from bounding the term $(1-u)^p$ by 1 in the integral, then we can notice directly that $I_p(t) \leq t^{-1}$ for every $t > 0$. The second result is an application of Watson's Lemma (see [58]).

□

A.2 Stochastic results

A.2.1 On stochastic processes

We refer to the notation and results discussed in [13, Section A.2].

Proposition A.9. Consider $\nu \geq 2$, $X_0, V_0 \in L^\nu(\Omega; \mathbb{H})$, f and σ satisfying (\mathbf{H}_0) and (\mathbf{H}_σ) , respectively. Consider also γ satisfying (\mathbf{H}_γ) , and β satisfying (\mathbf{H}_β) . Then $(\mathbf{S} - \mathbf{ISIHD})$ has a unique solution $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$.

Remark A.10. Hypothesis (\mathbf{H}_β) does not allow us to consider the case $\beta \equiv 0$, nevertheless, this case is well studied in Section 4.

Proof. We rewrite $(\mathbf{S} - \mathbf{ISIHD})$ as in the reformulation $(\mathbf{ISIHD} - \mathbf{S}_R)$, we recall (3.3)

$$\begin{cases} dZ(t) &= -\beta(t)\nabla\mathcal{G}(Z(t))dt - \mathcal{D}(t, Z(t))dt + \hat{\sigma}(t, Z(t))dW(t), \quad t > t_0, \\ Z(t_0) &= (X_0, X_0 + \beta(t_0)V_0), \end{cases}$$

Since $\beta(t) \leq c_1$, we have that $-\beta(t)\nabla\mathcal{G}(Z(t))$ is Lipschitz, besides, since $\left| \frac{\beta'(t) - \gamma(t)\beta(t) + 1}{\beta(t)} \right| \leq c_2$, we have that \mathcal{D} is a Lipschitz operator. Then, using the hypotheses on σ we can use [13, Theorem 3.3] and conclude the existence and uniqueness of a process $Z \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$, this, in turn, implies the existence and uniqueness of a solution $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$ of $(\mathbf{S} - \mathbf{ISIHD})$. \square

Theorem A.11. Let \mathbb{H} be a separable Hilbert space and $(M_t)_{t \geq 0} : \Omega \rightarrow \mathbb{H}$ be a continuous martingale such that $\sup_{t \geq 0} \mathbb{E}(\|M_t\|^2) < +\infty$. Then there exists a \mathbb{H} -valued random variable $M_\infty \in L^2(\Omega; \mathbb{H})$ such that $\lim_{t \rightarrow \infty} M_t = M_\infty$ a.s..

Proof. For the proof, we refer to [13, Theorem A.4]. \square

Theorem A.12. [59, 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 non-negative \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 non-negative and $\lim_{t \rightarrow \infty} A_t < \infty$, then $\lim_{t \rightarrow \infty} X_t$ exists and is finite, and $\lim_{t \rightarrow \infty} U_t < \infty$.

A.3 Abstract integral bounds, almost sure and in expectation properties of $(\mathbf{S} - \mathbf{ISIHD})$

In the following proposition, we state different abstract integral bounds and almost sure properties for $(\mathbf{S} - \mathbf{ISIHD})$, finally concluding with the almost sure convergence of the gradient towards zero.

Proposition A.13. Consider that f, σ satisfy (\mathbf{H}_0) and (\mathbf{H}_σ) , respectively. Let $\nu \geq 2$, and consider the dynamic $(\mathbf{S} - \mathbf{ISIHD})$ with initial data $X_0, V_0 \in L^\nu(\Omega; \mathbb{H})$. Consider also γ, β from $(\mathbf{S} - \mathbf{ISIHD})$ satisfying (\mathbf{H}_γ) and (\mathbf{H}_β) , respectively, and a, b, c, d satisfying $(\mathbf{S}_{a,b,c,d})$. Finally, we consider \mathcal{E} the energy function defined in (3.5).

Then, there exists a unique solution $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$ of $(\mathbf{S} - \mathbf{ISIHD})$. Moreover, if $t \mapsto m(t)\sigma_\infty^2(t) \in L^1([t_0, +\infty[)$, where $m(t) \stackrel{\text{def}}{=} \max\{1, a(t), c^2(t)\}$, then the following properties are satisfied:

(i)

$$\lim_{t \rightarrow +\infty} \mathcal{E}(t, X(t), V(t)) \text{ exists a.s..}$$

(ii) $\int_{t_0}^\infty (b(s)c(s) - a'(s))(f(X(s) + \beta(s)V(s)) - \min f)ds < +\infty$, a.s..

(iii) $\int_{t_0}^\infty a(s)\beta(s)\|\nabla f(X(s)) + \beta(s)V(s)\|^2 ds < +\infty$, a.s..

(iv) $\int_{t_0}^\infty \left(b(s)b'(s) + \frac{d'(s)}{2} \right) \|X(s) - x^*\|^2 ds < +\infty$, a.s..

(v) $\int_{t_0}^\infty c(s)(\gamma(s)c(s) - c'(s) - b(s))\|V(s)\|^2 ds < +\infty$, a.s..

(vi) If $b(t)c(t) - a'(t) = \mathcal{O}(c(t)(\gamma(t)c(t) - c'(t) - b(t)))$, then

$$\int_{t_0}^{\infty} (b(s)c(s) - a'(s))(f(X(s)) - \min f) ds < +\infty \quad a.s..$$

(vii) If there exists $\eta > 0, \hat{t} > t_0$ such that

$$\eta \leq c(t)(\gamma(t)c(t) - c'(t) - b(t)), \quad \eta \leq a(t)\beta(t), \quad \gamma(t) \leq \eta, \quad \forall t > \hat{t},$$

then $\lim_{t \rightarrow \infty} \|V(t)\| = 0$ a.s., $\lim_{t \rightarrow \infty} \|\nabla f(X(t) + \beta(t)V(t))\| = 0$ a.s., and $\lim_{t \rightarrow \infty} \|\nabla f(X(t))\| = 0$ a.s.

Proof. The existence and uniqueness of a solution is a direct consequence of Corollary A.9. Moreover, applying Proposition 2.3 with \mathcal{E} , we can obtain

$$\begin{aligned} \mathcal{E}(t, X(t), V(t)) &\leq \mathcal{E}(t_0, X_0, V_0) - \int_{t_0}^t (b(s)c(s) - a'(s))(f(X(s) + \beta(s)V(s)) - \min f) ds \\ &\quad - \int_{t_0}^t a(s)\beta(s)\|\nabla f(X(s) + \beta(s)V(s))\|^2 ds - \int_{t_0}^t \left(b(s)b'(s) + \frac{d'(s)}{2} \right) \|X(s) - x^*\|^2 ds \\ &\quad - \int_{t_0}^t c(s)(b(s) + c'(s) - c(s)\gamma(s))\|V(s)\|^2 ds + \int_{t_0}^t (La(s)\beta^2(s) + c^2(s))\sigma_{\infty}^2(s) ds + M_t, \end{aligned} \tag{A.2}$$

where $M_t = \int_{t_0}^t \langle \sigma^*(s, X(s) + \beta(s)V(s))(a(s)\beta(s)\nabla f(X(s) + \beta(s)V(s)) + c(s)[b(s)(X(s) - x^*) + c(s)V(s)], dW(s) \rangle$. Since $\sup_{t \in [t_0, T]} \mathbb{E}(\|X(t)\|^2) < +\infty, \sup_{t \in [t_0, T]} \mathbb{E}(\|V(t)\|^2) < +\infty$ for every $T > t_0$, and a, b, c, β are continuous functions, we have that M_t is a continuous martingale, on the other hand, we have that

$$\int_{t_0}^{\infty} (La(s)\beta^2(s) + c^2(s))\sigma_{\infty}^2(s) ds < +\infty.$$

Then, we can apply Theorem A.12 and conclude that $\lim_{t \rightarrow \infty} \mathcal{E}(t, X(t), V(t))$ exists a.s. and

- $\int_{t_0}^{\infty} (b(s)c(s) - a'(s))(f(X(s) + \beta(s)V(s)) - \min f) ds < +\infty$ a.s..
- $\int_{t_0}^{\infty} a(s)\beta(s)\|\nabla f(X(s) + \beta(s)V(s))\|^2 ds < +\infty$ a.s..
- $\int_{t_0}^{\infty} \left(b(s)b'(s) + \frac{d'(s)}{2} \right) \|X(s) - x^*\|^2 ds < +\infty$ a.s..
- $\int_{t_0}^{\infty} c(s)(\gamma(s)c(s) - c'(s) - b(s))\|V(s)\|^2 ds < +\infty$ a.s..

This let us conclude with items (i) to (v).

Let $\tilde{b}(t) = b(t)c(t) - a'(t)$, and

$$\begin{aligned} I_f &\stackrel{\text{def}}{=} \int_{t_0}^{\infty} \tilde{b}(s)(f(X(s)) - \min f) ds \\ &\leq \int_{t_0}^{\infty} \tilde{b}(s)(f(X(s)) - f(X(s) + \beta(s)V(s))) ds + \int_{t_0}^{\infty} \tilde{b}(s)(f(X(s) + \beta(s)V(s)) - \min f) ds \end{aligned}$$

Using Descent Lemma, Cauchy Schwarz Inequality and Corollary 2.2:

$$\begin{aligned} I_f &\leq \sqrt{2L}\beta_0 \left(\int_{t_0}^{\infty} \tilde{b}(s)(f(X(s) + \beta(s)V(s)) - \min f) ds \right)^{\frac{1}{2}} \left(\int_{t_0}^{\infty} \tilde{b}(s)\|V(s)\|^2 ds \right)^{\frac{1}{2}} \\ &\quad + \frac{L\beta_0^2}{2} \int_{t_0}^{\infty} \tilde{b}(s)\|V(s)\|^2 ds + \int_{t_0}^{\infty} \tilde{b}(s)(f(X(s) + \beta(s)V(s)) - \min f) ds. \end{aligned}$$

If $\tilde{b}(t) = b(t)c(t) - a'(t) = \mathcal{O}(c(t)(\gamma(t)c(t) - c'(t) - b(t)))$, we have that

$$\int_{t_0}^{\infty} \tilde{b}(s) \|V(s)\|^2 ds < +\infty \quad a.s..$$

And we conclude with item (vi).

To prove (vii), in particular that $\lim_{t \rightarrow \infty} \|V(t)\| = 0$, we consider that if there exists $\eta > 0, \hat{t} > t_0$ such that $\eta \leq c(t)(\gamma(t)c(t) - c'(t) - b(t)), \forall t > \hat{t}$, then there exists $\Omega_v \in \mathcal{F}$ such that $\mathbb{P}(\Omega_v) = 1$ and

$$\int_{t_0}^{\infty} \|V(\omega, s)\|^2 ds < +\infty, \quad \forall \omega \in \Omega_v.$$

Then, we have $\liminf_{t \rightarrow \infty} \|V(\omega, t)\| = 0, \forall \omega \in \Omega_v$. Let us suppose that $\limsup_{t \rightarrow \infty} \|V(\omega, t)\| > 0, \forall \omega \in \Omega_v$. Then, by [13, Lemma A.3], there exists $\delta > 0$ satisfying

$$0 = \liminf_{t \rightarrow \infty} \|V(\omega, t)\| < \delta < \limsup_{t \rightarrow \infty} \|V(\omega, t)\|, \quad \forall \omega \in \Omega_v.$$

And there exists $(t_k)_{k \in \mathbb{N}} \subset [t_0, +\infty[$ such that $\lim_{k \rightarrow \infty} t_k = +\infty$,

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

Let $N_t \stackrel{\text{def}}{=} \int_{t_0}^t \sigma(s, X(s) + \beta(s)V(s)) dW(s)$. This is a continuous martingale (w.r.t. the filtration \mathcal{F}_t), which verifies

$$\mathbb{E}(\|N_t\|^2) = \mathbb{E} \left(\int_{t_0}^t \|\sigma(s, X(s) + \beta(s)V(s))\|_{\text{HS}}^2 ds \right) \leq \mathbb{E} \left(\int_{t_0}^{\infty} \sigma_{\infty}^2(s) ds \right) < +\infty, \forall t \geq t_0.$$

According to Theorem A.11, we deduce that there exists a \mathbb{H} -valued random variable N_{∞} w.r.t. \mathcal{F}_{∞} , and which verifies: $\mathbb{E}(\|N_{\infty}\|^2) < +\infty$, and there exists $\Omega_N \in \mathcal{F}$ such that $\mathbb{P}(\Omega_N) = 1$ and

$$\lim_{t \rightarrow +\infty} N_t(\omega) = N_{\infty}(\omega) \text{ for every } \omega \in \Omega_N.$$

Let $\omega_0 \in \Omega_{nv} \stackrel{\text{def}}{=} \Omega_N \cap \Omega_v$ ($\mathbb{P}(\Omega_{nv}) = 1$) and the notation $V(t) \stackrel{\text{def}}{=} V(\omega_0, t)$, $\varepsilon \in \left(0, \min\{1, \frac{\delta^2}{4}\}\right)$ arbitrary and recall that $\eta \leq a(t)\beta(t), \gamma(t) \leq \eta$ for every $t > \hat{t}$. Let $k' \in \mathbb{N}$ be such that $t_{k'} > \hat{t}, k > k'$ and $t \in [t_k, t_k + \varepsilon]$, then

$$\begin{aligned} \|V(t) - V(t_k)\|^2 &\leq 3(t - t_k) \int_{t_k}^t \gamma^2(s) \|V(s)\|^2 ds + 3(t - t_k) \int_{t_k}^t \|\nabla f(X(s) + \beta(s)V(s))\|^2 ds \\ &\quad + 3\|N_t - N_{t_k}\|^2 \\ &\leq 3\eta^2(t - t_k) \int_{t_k}^t \|V(s)\|^2 ds + \frac{3}{\eta}(t - t_k) \int_{t_k}^t a(s)\beta(s) \|\nabla f(X(s) + \beta(s)V(s))\|^2 ds \\ &\quad + 3\|N_t - N_{t_k}\|^2 \\ &\leq 3\eta^2 \varepsilon \int_{t_k}^t \|V(s)\|^2 ds + \frac{3}{\eta} \varepsilon \int_{t_k}^t a(s)\beta(s) \|\nabla f(X(s) + \beta(s)V(s))\|^2 ds \\ &\quad + 3\|N_t - N_{t_k}\|^2. \end{aligned}$$

Now let $k'' \in \mathbb{N}$ be such that for every $k > k''$,

$$\int_{t_k}^{\infty} \|V(s)\|^2 ds < \frac{1}{9\eta^2}, \int_{t_k}^{\infty} a(s)\beta(s) \|\nabla f(X(s) + \beta(s)V(s))\|^2 ds < \frac{\eta}{9}, \sup_{t > t_k} \|N_t - N_{t_k}\|^2 < \frac{\varepsilon}{9}.$$

Then, we have that

$$\|V(t) - V(t_k)\|^2 \leq \varepsilon \leq \frac{\delta^2}{4}, \forall t \in [t_k, t_k + \varepsilon], k > \max\{k', k''\}.$$

For such t , we bound using the triangular inequality and obtain

$$\|V(t)\| \geq \|V(t_k)\| - \|V(t) - V(t_k)\| > \frac{\delta}{2}.$$

Now we consider

$$\int_{t_0}^{\infty} \|V(s)\|^2 ds \geq \sum_{k > \max\{k', k''\}} \int_{t_k}^{t_k + \varepsilon} \|V(s)\|^2 ds \geq \sum_{k > \max\{k', k''\}} \frac{\varepsilon \delta^2}{4} = +\infty.$$

Which is a contradiction, then we conclude that $\liminf_{t \rightarrow \infty} \|V(t)\| = \limsup_{t \rightarrow \infty} \|V(t)\| = 0$, a.s..

To prove the second part of (vii), i.e. that $\lim_{t \rightarrow \infty} \|\nabla f(X(t) + \beta(t)V(t))\| = 0$, we recall that there exists $\eta > 0$, $\hat{t} > t_0$ such that $\eta \leq a(t)\beta(t)$, then there exists $\Omega_y \in \mathcal{F}$ such that $\mathbb{P}(\Omega_y) = 1$ and

$$\int_{t_0}^{\infty} \|\nabla f(X(\omega, s) + \beta(s)V(\omega, s))\|^2 ds < +\infty \quad \forall \omega \in \Omega_y$$

So we have that

$$\liminf_{t \rightarrow \infty} \|\nabla f(X(\omega, t) + \beta(t)V(\omega, t))\| = 0, \quad \forall \omega \in \Omega_y.$$

Moreover, if we suppose that

$$\limsup_{t \rightarrow \infty} \|\nabla f(X(\omega, t) + \beta(t)V(\omega, t))\| > 0, \quad \forall \omega \in \Omega_y,$$

by [13, Lemma A.3], there exists $\delta > 0$, $(t_k)_{k \in \mathbb{N}} \subset [t_0, +\infty[$ such that $\lim_{k \rightarrow \infty} t_k = +\infty$,

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

Recall that by (H_β) , there exists β_0 such that $\beta(t) \leq \beta_0$. Let $\varepsilon \in (0, \min\{1, \frac{\delta^2}{4L^2}\})$ arbitrary and consider $Y(t) = X(t) + \beta(t)V(t)$, let also $k \in \mathbb{N}$ arbitrary and $t \in [t_k, t_k + \varepsilon]$. Then, using Lemma A.1 and Jensen's inequality we can bound as follows:

$$\begin{aligned} \|Y(t) - Y(t_k)\|^2 &\leq 2\|X(t) - X(t_k)\|^2 + 2\|\beta(t)V(t) - \beta(t_k)V(t_k)\|^2 \\ &\leq 2\left\|\int_{t_k}^t V(s)ds\right\|^2 + 2(\beta_0\|V(t) - V(t_k)\| + |\beta(t) - \beta(t_k)| \max\{\|V(t)\|, \|V(t_k)\|\})^2 \\ &\leq 2(t - t_k) \int_{t_k}^{\infty} \|V(s)\|^2 ds \\ &\quad + 2(\beta_0\|V(t) - V(t_k)\| + |\beta(t) - \beta(t_k)| \max\{\|V(t)\|, \|V(t_k)\|\})^2 \\ &\leq 2(t - t_k) \int_{t_k}^{\infty} \|V(s)\|^2 ds \\ &\quad + 4\beta_0^2\|V(t) - V(t_k)\|^2 + 4|\beta(t) - \beta(t_k)|^2 \max\{\|V(t)\|, \|V(t_k)\|\}^2. \end{aligned}$$

By the previous point, we have that there exists $\Omega_v \in \mathcal{F}$ such that $\mathbb{P}(\Omega_v) = 1$ such that

$$\int_{t_0}^{\infty} \|V(\omega, s)\|^2 ds < +\infty \quad \forall \omega \in \Omega_v,$$

and $k' \in \mathbb{N}$ such that for every $k > k'$, for all $t \in [t_k, t_k + \varepsilon]$:

$$\int_{t_k}^{\infty} \|V(\omega, s)\|^2 ds < \frac{1}{6}, \quad \max\{\|V(\omega, t)\|, \|V(\omega, t_k)\|\} < 1, \quad \|V(\omega, t) - V(\omega, t_k)\|^2 \leq \frac{\varepsilon}{12\beta_0^2}$$

We consider an arbitrary $\omega_0 \in \Omega_y \cap \Omega_v$ ($\mathbb{P}(\Omega_y \cap \Omega_v) = 1$), and we let us use the abuse of notation $X(t) = X(\omega_0, t)$, $V(t) = V(\omega_0, t)$, and $Y(t) = Y(\omega_0, t)$ for the rest of this proof.

On the other hand, β is continuous, so there exists $\tilde{\delta} > 0$ such that, if $|t - t_k| < \tilde{\delta}$, then $|\beta(t) - \beta(t_k)| < \frac{\sqrt{\varepsilon}}{2\sqrt{3}}$.

Therefore, letting $\varepsilon' = \min\{\varepsilon, \tilde{\delta}\}$, we have that

$$\|\nabla f(Y(t)) - \nabla f(Y(t_k))\|^2 \leq L^2 \|Y(t) - Y(t_k)\|^2 \leq L^2 \varepsilon \leq \frac{\delta^2}{4}, \forall k > k', \forall t \in [t_k, t_k + \varepsilon'].$$

Then, we obtain

$$\|\nabla f(Y(t))\| \geq \|\nabla f(Y(t_k))\| - \|\nabla f(Y(t)) - \nabla f(Y(t_k))\| \geq \frac{\delta}{2}, \forall k > k', \forall t \in [t_k, t_k + \varepsilon'].$$

This implies that

$$\int_{t_0}^{\infty} \|\nabla f(Y(s))\|^2 ds \geq \sum_{k > k'} \int_{t_k}^{t_k + \varepsilon'} \|\nabla f(Y(s))\|^2 ds \geq \sum_{k > k'} \int_{t_k}^{t_k + \varepsilon'} \frac{\delta^2}{4} = \sum_{k > k'} \frac{\varepsilon' \delta^2}{4} = +\infty.$$

Which is a contradiction, then we conclude that

$$\liminf_{t \rightarrow \infty} \|\nabla f(X(t) + \beta(t)V(t))\| = \limsup_{t \rightarrow \infty} \|\nabla f(X(t) + \beta(t)V(t))\| = 0, \quad a.s..$$

To prove the last part of (vii), we consider that $\beta(t) \leq \beta_0$, then

$$\begin{aligned} \|\nabla f(X(t))\| &\leq \|\nabla f(X(t) + \beta(t)V(t))\| + \|\nabla f(X(t)) - \nabla f(X(t) + \beta(t)V(t))\| \\ &\leq \|\nabla f(X(t) + \beta(t)V(t))\| + L\beta_0 \|V(t)\|. \end{aligned}$$

With this bound, we can conclude that $\lim_{t \rightarrow \infty} \|\nabla f(X(t))\| = 0$ a.s.. □

The following proposition states abstract bounds in expectation of (S – ISIHD).

Proposition A.14. *Consider the setting of Proposition A.13, then we have that:*

$$(i) \quad \mathbb{E}(f(X(t) + \beta(t)V(t)) - \min f) = \mathcal{O}\left(\frac{1}{a(t)}\right).$$

Moreover, if there exists $D > 0, \tilde{t} > t_0$ such that $d(t) \geq D$ for $t > \tilde{t}$, then :

$$(ii) \quad \sup_{t > t_0} \mathbb{E}(\|X(t) - x^*\|^2) < +\infty.$$

$$(iii) \quad \mathbb{E}(\|V(t)\|^2) = \mathcal{O}\left(\frac{1 + b^2(t)}{c^2(t)}\right).$$

$$(iv) \quad \mathbb{E}(f(X(t)) - \min f) = \mathcal{O}\left(\max\left\{\frac{1}{a(t)}, \frac{\beta(t)\sqrt{1 + b^2(t)}}{\sqrt{a(t)}c(t)}, \frac{\beta^2(t)(1 + b^2(t))}{c^2(t)}\right\}\right).$$

Proof. To prove this proposition we are going to take expectation in (A.2). First, we are going to bound the negative terms by 0, denoting $E_0 \stackrel{\text{def}}{=} \mathcal{E}(t_0) + \max\{1, L\} \int_{t_0}^{\infty} (a(s)\beta^2(s) + c^2(s))\sigma_{\infty}^2(s) ds$, we obtain that

$$\mathbb{E}(\mathcal{E}(t, X(t), V(t))) \leq E_0.$$

This implies that

- $\mathbb{E}(f(X(t) + \beta(t)V(t)) - \min f) \leq \frac{E_0}{a(t)}$.
- $\mathbb{E}(\|b(t)(X(t) - x^*) + c(t)V(t)\|^2) \leq 2E_0$.

If there exists $D > 0, \tilde{t} > t_0$ such that $d(t) \geq D$ for $t > \tilde{t}$, then for $t > \tilde{t}$:

- $\mathbb{E}(\|X(t) - x^*\|^2) \leq \frac{2E_0}{D}$.

- And also,

$$\begin{aligned}\mathbb{E}(\|V(t)\|^2) &\leq \frac{2}{c^2(t)} [\mathbb{E}(\|b(t)(X(t) - x^*) + c(t)V(t)\|^2) + b^2(t)\mathbb{E}(\|X(t) - x^*\|^2)] \\ &\leq \frac{2}{c^2(t)} \left(2E_0 + \frac{2E_0 b^2(t)}{D} \right) = \frac{4E_0}{c^2(t)} \left(1 + \frac{b^2(t)}{D} \right).\end{aligned}$$

- We bound the following term using the Descent Lemma

$$\begin{aligned}\mathbb{E}(f(X(t)) - f(X(t) + \beta(t)V(t))) &\leq \beta(t)\sqrt{\mathbb{E}(\|\nabla f(X(t) + \beta(t)V(t))\|^2)}\sqrt{\mathbb{E}(\|V(t)\|^2)} \\ &\quad + \frac{L}{2}\beta^2(t)\mathbb{E}(\|V(t)\|^2).\end{aligned}$$

Using Corollary 2.2, we have

$$\begin{aligned}\mathbb{E}(f(X(t)) - f(X(t) + \beta(t)V(t))) &\leq \beta(t)\sqrt{2L\mathbb{E}(f(X(t) + \beta(t)V(t)) - \min f)}\sqrt{\mathbb{E}(\|V(t)\|^2)} \\ &\quad + \frac{L}{2}\beta^2(t)\mathbb{E}(\|V(t)\|^2) \\ &\leq 2E_0\sqrt{2L}\frac{\beta(t)}{\sqrt{a(t)}}\frac{\sqrt{1 + \frac{b^2(t)}{D}}}{c(t)} + 2LE_0\frac{\beta^2(t)\left(1 + \frac{b^2(t)}{D}\right)}{c^2(t)} \\ &= \mathcal{O}\left(\max\left\{\frac{\beta(t)}{\sqrt{a(t)}}\frac{\sqrt{1 + b^2(t)}}{c(t)}, \frac{\beta^2(t)(1 + b^2(t))}{c^2(t)}\right\}\right)\end{aligned}$$

Then, we notice that

$$\begin{aligned}\mathbb{E}(f(X(t)) - \min f) &= \mathbb{E}[f(X(t)) - f(X(t) + \beta(t)V(t))] + \mathbb{E}[f(X(t) + \beta(t)V(t)) - \min f] \\ &= \mathcal{O}\left(\max\left\{\frac{\beta(t)}{\sqrt{a(t)}}\frac{\sqrt{1 + b^2(t)}}{c(t)}, \frac{\beta^2(t)(1 + b^2(t))}{c^2(t)}, \frac{1}{a(t)}\right\}\right).\end{aligned}$$

□

References

- [1] C. Alecsa, S. László, and T. Pinta. An extension of the second order dynamical system that models Nesterov's convex gradient method. *Appl. Math. Optim.*, 84:1687–1716, 2021.
- [2] M. Muehlebach and M.I. Jordan. A dynamical systems perspective on Nesterov acceleration. In *Proceedings of the 36 th International Conference on Machine Learning*, volume 97. PMLR, 2019.
- [3] Hedy Attouch, Jalal Fadili, and Vyacheslav Kungurtsev. On the effect of perturbations in first-order optimization methods with inertia and Hessian driven damping. *Evolution equations and Control*, 12(1):71–117, 2023.
- [4] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic gradient algorithms. *arXiv:1511.06251*, 2017.
- [5] Antonio Orvieto and Aurelien Lucchi. Continuous-time models for stochastic optimization algorithms. *33rd Conference on Neural Information Processing Systems (NeurIPS 2019)*, 2019.
- [6] Wenqing Hu, Chris Junchi Li, Lei Li, and Jian-Guo Lui. On the diffusion approximation of nonconvex stochastic gradient descent. *arXiv:1705.07562v2*, 2018.
- [7] Bin Shi, Weijie J. Su, and Michael I. Jordan. On learning rates and Schrödinger operators. *Journal of Machine Learning Research*, 24:1–53, 2023.
- [8] Zhiyuan Li, Sadhika Malladi, and Sanjeev Arora. On the validity of modeling sgd with stochastic differential equations. *arXiv:2102.12470*, 2021.
- [9] 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.
- [10] 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.

- [11] Rodrigo Maulen S., Jalal Fadili, and Hedy Attouch. An SDE perspective on stochastic convex optimization. *arXiv:2207.02750*, 2022.
- [12] M. Dambrine, C. Dossal, B. Puig, and A. Rondepierre. Stochastic differential equations for modeling first order optimization methods. *Hal*, 2022.
- [13] Rodrigo Maulen-Soto, Jalal Fadili, and Hedy Attouch. Tikhonov regularization for stochastic non-smooth convex optimization in Hilbert spaces. *arXiv:2403.06708*, 2024.
- [14] Hedy Attouch and Alexandre Cabot. Asymptotic stabilization of inertial gradient dynamics with time-dependent viscosity. *J. Differential Equations*, 263(9):5412–5458, 2017.
- [15] Grigorios A. Pavliotis. Stochastic processes and applications. *Springer*, 2014.
- [16] 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.
- [17] Weijie Su, Stephen Boyd, and Emmanuel J. Candès. A differential equation for modeling Nesterov’s accelerated gradient method: Theory and insights. *Journal of Machine Learning Research*, 17:1–43, 2016.
- [18] Y.E. Nesterov. A method of solving a convex programming problem with convergence rate $\mathcal{O}(1/k^2)$. *Doklady Akademii Nauk SSSR*, 269(3):543–547, 1983.
- [19] H. Attouch and A. Cabot. Asymptotic stabilization of inertial gradient dynamics with time-dependent viscosity. *Journal of Differential Equations*, 263-9:5412–5458, 2017.
- [20] Hedy Attouch and Juan Peypouquet. The rate of convergence of Nesterov’s accelerated forward-backward method is actually faster than $\frac{1}{k^2}$. *SIAM Journal on Optimization*, 26(3):1824–1834, 2016.
- [21] Hedy Attouch, Zaki Chbani, Jalal Fadili, and Hassan Riahi. First-order optimization algorithms via inertial systems with Hessian driven damping. *Mathematical Programming*, 193(4), 2020.
- [22] Camille Castera, Hedy Attouch, Jalal Fadili, and Peter Ochs. Continuous Newton-like methods featuring inertia and variable mass. *arXiv:2301.08726*, 2023.
- [23] Xiang Cheng, Niladri S. Chatterji, Peter L. Bartlett, and Michael I. Jordan. Underdamped Langevin MCMC: A non-asymptotic analysis. *arXiv:1707.03663*, 2017.
- [24] Yi-An Ma, Niladri Chatterji, Xiang Cheng, Nicolas Flammarion, Peter Bartlett, and Michael I. Jordan. Is there an analog of Nesterov acceleration for MCMC? *Bernoulli*, 27 (3):1942–1992, 2021.
- [25] Arnak S. Dalalyan, Lionel Riou-Durand, and Avetik G. Karagulyan. Bounding the error of discretized langevin algorithms for non-strongly log-concave targets. *J. Mach. Learn. Res.*, 23:235:1–235:38, 2019.
- [26] A. Haraux and Jendoubi M.A. On a second order dissipative ODE in Hilbert space with an integrable source term. *Acta Math. Sci.*, 32:155–163, 2012.
- [27] B. Shi, S.S. Du, M.I. Jordan, and Su W.J. Understanding the acceleration phenomenon via high resolution differential equations. *Math. Program.*, 2021.
- [28] H. Attouch, A. Cabot, Chbani Z., and H. Riahi. Accelerated forward-backward algorithms with perturbations: Application to Tikhonov regularization. *J. Optim. Theory Appl.*, 179:1–36, 2018.
- [29] H. Attouch, Z. Chbani, J. Peypouquet, and P. Redont. Fast convergence of inertial dynamics and algorithms with asymptotic vanishing viscosity. *Math. Program. Ser. B*, 168:123–175, 2018.
- [30] C. Dossal and J.F. Aujol. Stability of over-relaxations for the forward-backward algorithm, application to fista. *SIAM J. Optim.*, 25:2408–2433, 2015.
- [31] M. Schmidt, N. Le Roux, and F. Bach. Convergence rates of inexact proximal-gradient methods for convex optimization. *NIPS’11*, 25th Annual Conference, 2011.
- [32] S. Villa, S. Salzo, and Baldassarres L. Accelerated and inexact forward-backward. *SIAM J. Optim.*, 23:1607–1633, 2013.
- [33] H. Attouch, J. Peypouquet, and P. Redont. Fast convex minimization via inertial dynamics with Hessian driven damping. *J. Differential Equations*, 261:5734–5783, 2016.
- [34] H. Attouch, Z. Chbani, J. Fadili, and H. Riahi. First order optimization algorithms via inertial systems with Hessian driven damping. *Math. Program.*, 2020.
- [35] H. Attouch, Z. Chbani, J. Fadili, and H. Riahi. Convergence of iterates for first-order optimization algorithms with inertia and Hessian driven damping. *Optimization*, 2021.
- [36] A. Orvieto, J. Kohler, and A. Lucchi. The role of memory in stochastic optimization. *Proceeding of Machine Learning Research*, 115:356–366, 2020.
- [37] Sébastien Gadat and Fabien Panloup. Long time behaviour and stationary regime of memory gradient diffusions. *Annales de l’Institut Henri Poincaré - Probabilités et Statistiques*, 2014.
- [38] Hongzhou Lin, Julien Mairal, and Zaid Harchaoui. Catalyst acceleration for first-order convex optimization: from theory to practice. *J. Mach. Learn. Res.*, 18:Paper No. 212, 54, 2017.

- [39] R. Frostig, S. Kakade R. Ge, and A. Sidford. Un-regularizing: approximate proximal point and faster stochastic algorithms for empirical risk minimization. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37, pages 2540–2548, 2015.
- [40] Prateek Jain, Praneeth Netrapalli, Sham M. Kakade, Rahul Kidambi, and Aaron Sidford. Parallelizing stochastic gradient descent for least squares regression: mini-batching, averaging, and model misspecification. *J. Mach. Learn. Res.*, 18:Paper No. 223, 42, 2017.
- [41] M. Assran and M. Rabbat. On the convergence of nesterov’s accelerated gradient method in stochastic settings. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 410–420, 2020.
- [42] Zeyuan Allen-Zhu. Katyusha: the first direct acceleration of stochastic gradient methods. *J. Mach. Learn. Res.*, 18:Paper No. 221, 51, 2017.
- [43] Bowei Yan. *Theoretical Analysis for Convex and Non-Convex Clustering Algorithms*. ProQuest LLC, Ann Arbor, MI, 2018. Thesis (Ph.D.)—The University of Texas at Austin.
- [44] Sébastien Gadat, Fabien Panloup, and Sofiane Saadane. Stochastic heavy ball. *Electron. J. Stat.*, 12(1):461–529, 2018.
- [45] Nicolas Loizou and Peter Richtárik. Momentum and stochastic momentum for stochastic gradient, Newton, proximal point and subspace descent methods. *Comput. Optim. Appl.*, 77(3):653–710, 2020.
- [46] Maxime Laborde and Adam Oberman. A lyapunov analysis for accelerated gradient methods: from deterministic to stochastic case. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 602–612. PMLR, 26–28 Aug 2020.
- [47] Guanghui Lan. *First-order and stochastic optimization methods for machine learning*. Springer Series in the Data Sciences. Springer, Cham, [2020] ©2020.
- [48] Aaron Defazio and Samy Jelassi. Adaptivity without compromise: a momentumized, adaptive, dual averaged gradient method for stochastic optimization. *J. Mach. Learn. Res.*, 23:Paper No. [144], 34, 2022.
- [49] Derek Driggs, Matthias J. Ehrhardt, and Carola-Bibiane Schönlieb. Accelerating variance-reduced stochastic gradient methods. *Math. Program.*, 191(2, Ser. A):671–715, 2022.
- [50] Anis Hamadouche, Yun Wu, Andrew M. Wallace, and João F. C. Mota. Sharper bounds for proximal gradient algorithms with errors. *SIAM Journal on Optimization*, 34(1):278–305, 2024.
- [51] Hedy Attouch, Jalal Fadili, and Vyacheslav Kungurtsev. The stochastic ravine accelerated gradient method with general extrapolation coefficients, 2024.
- [52] Hedy Attouch, Radu Ioan Bot, and Dang-Khoa Nguyen. Fast convex optimization via time scale and averaging of the steepest descent. *arXiv:2208.08260*, 2022.
- [53] Rodrigo Maulen-Soto, Jalal Fadili, Hedy Attouch, and Peter Ochs. Stochastic inertial dynamics via time scaling and averaging. *arXiv:2403.16775*, 2024.
- [54] R.T. Rockafellar. *Convex analysis*. Princeton university press, 28, 1997.
- [55] Leszek Gawarecki and Vidyadhar Mandrekar. *Stochastic differential equations in infinite dimensions*. Springer, 2011.
- [56] Z. Opial. Weak convergence of the sequence of successive approximations for nonexpansive mappings. *Bull. Amer. Math. Soc.*, 73:591–597, 1967.
- [57] F.W.J. Olver, D.w. Lozier, R.F Boisvert, and Clark C.W. *Nist handbook of mathematical functions*. Cambridge University Press, 2010.
- [58] G.N. Watson. The harmonic functions associated with the parabolic cylinder. *Proceedings of the London Mathematical Society*, 2, no. 17:116–148, 1918.
- [59] Xuerong Mao. *Stochastic differential equations and applications*. Elsevier, 2007.