



Stochastic Systems

Publication details, including instructions for authors and subscription information:
<http://pubsonline.informs.org>

Stochastic Inertial Dynamics via Time Scaling and Averaging

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

To cite this article:

Rodrigo Maulen-Soto, Jalal Fadili, Hédya Attouch, Peter Ochs (2026) Stochastic Inertial Dynamics via Time Scaling and Averaging. *Stochastic Systems* 16(1):61-89. <https://doi.org/10.1287/stsy.2024.0068>

This work is licensed under a Creative Commons Attribution 4.0 International License. You are free to copy, distribute, transmit and adapt this work, but you must attribute this work as “*Stochastic Systems*. Copyright © 2026 The Author(s). <https://doi.org/10.1287/stsy.2024.0068>, used under a Creative Commons Attribution License: <https://creativecommons.org/licenses/by/4.0/>.”

Copyright © 2026 The Author(s)

Please scroll down for article—it is on subsequent pages







With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes. For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

Stochastic Inertial Dynamics via Time Scaling and Averaging

 Rodrigo Maulen-Soto,^{a,*} Jalal Fadili,^a Hedy Attouch,^b Peter Ochs^c
^aGREYC: Groupe de recherche en informatique, image et instrumentation de Caen, Ensicaen, Université de Caen, Normandie Université, 14000 Caen, France; ^bIMAG: Institut Montpellierain Alexander Grothendieck, Université de Montpellier, 34090 Montpellier, France;

^cDepartment of Mathematics and Computer Science, Saarland University, 66123 Saarbrücken, Germany

*Corresponding author

Contact: maulen@lpsm.paris,  <https://orcid.org/0009-0007-1383-3297> (RM-S); jalal.fadili@ensicaen.fr,  <https://orcid.org/0000-0002-8165-7578> (JF); hedy.attouch@umontpellier.fr,  <https://orcid.org/0000-0003-2676-0887> (HA); ochs@cs.uni-saarland.de,  <https://orcid.org/0000-0002-4880-7511> (PO)

Received: March 27, 2024

Revised: February 18, 2025; October 3, 2025

Accepted: October 24, 2025

Published Online in Articles in Advance: January 19, 2026

<https://doi.org/10.1287/stsy.2024.0068>
Copyright: © 2026 The Author(s)

Abstract. Our work is part of the close link between continuous-time dissipative dynamical systems and optimization algorithms, and more precisely here, in the stochastic setting. We aim to study stochastic convex minimization problems through the lens of stochastic inertial differential inclusions that are driven by the subgradient of a convex objective function. This will provide a general mathematical framework for analyzing the convergence properties of stochastic second-order inertial continuous-time dynamics involving vanishing viscous damping and measurable stochastic subgradient selections. Our chief goal in this paper is to develop a systematic and unified way that transfers the properties recently studied for first-order stochastic differential equations to second-order ones involving even subgradients in lieu of gradients. This program will rely on two tenets: time scaling and averaging, following an approach recently developed in the literature by one of the coauthors in the deterministic case. Under a mild integrability assumption involving the diffusion term and the viscous damping, our first main result shows that almost surely, there is weak convergence of the trajectory toward a minimizer of the objective function and fast convergence of the values and gradients. We also provide a comprehensive complexity analysis by establishing several new pointwise and ergodic convergence rates in expectation for the convex, strongly convex, and (local) Polyak-Łojasiewicz case. Finally, using Tikhonov regularization with a properly tuned vanishing parameter, we can obtain almost sure strong convergence of the trajectory toward the minimum norm solution.



Open Access Statement: This work is licensed under a Creative Commons Attribution 4.0 International License. You are free to copy, distribute, transmit and adapt this work, but you must attribute this work as “*Stochastic Systems*. Copyright © 2026 The Author(s). <https://doi.org/10.1287/stsy.2024.0068>, used under a Creative Commons Attribution License: <https://creativecommons.org/licenses/by/4.0/>.”

Funding: This research was supported by Agence Nationale de la Recherche (ANR) [Projet-ANR-20-CE92-0037].

Keywords: stochastic convex optimization • inertial gradient systems • stochastic differential equation • Tikhonov regularization • time-dependent viscosity

1. Introduction

1.1. Problem Statement

We consider the minimization problem

$$\min_{x \in \mathbb{H}} F(x) \stackrel{\text{def}}{=} f(x) + g(x), \quad (\text{P})$$

 where \mathbb{H} is a separable real Hilbert space, and the objective F satisfies the following standing assumptions:

$$\left\{ \begin{array}{l} f : \mathbb{H} \rightarrow \mathbb{R} \text{ is continuously differentiable and convex with } L - \text{Lipschitz} \\ \text{continuous gradient;} \\ g : \mathbb{H} \rightarrow \mathbb{R} \cup \{\pm\infty\} \text{ is proper, lower semi-continuous (lsc) and convex;} \\ S_F \stackrel{\text{def}}{=} \arg \min(F) \neq \emptyset. \end{array} \right. \quad (\text{H}_0)$$

1.1.1. First-Order in Time Systems. To solve (P) when $g \equiv 0$, a fundamental dynamic is the gradient flow system:

$$\left\{ \begin{array}{l} \dot{x}(t) + \nabla f(x(t)) = 0, \quad t > t_0; \\ x(t_0) = x_0. \end{array} \right. \quad (\text{GF})$$

The gradient system (GF) is a dissipative dynamical system, whose study dates back to Cauchy Cauchy (1847) in finite dimension. It plays a fundamental role in optimization: It transforms the problem of minimizing f into the study of the asymptotic behavior of the trajectories of (GF). This example was the precursor to the rich connection between continuous dissipative dynamical systems and optimization. It is well known since the founding papers of Brezis, Baillon, and Bruck in the 1970s that, if the solution set $\arg \min(f)$ of (P) is nonempty, then each solution trajectory of (GF) converges weakly, and its (weak) limit belongs to $\arg \min(f)$. Moreover, this dynamic is known to yield a convergence rate of $\mathcal{O}(t^{-1})$ (in fact even $o(t^{-1})$) on the values.

The Euler forward (a.k.a. Euler-Maruyama) discretization (GF), with stepsize sequence $h_k > 0$, is the celebrated gradient descent scheme

$$x_{k+1} = x_k - h_k \nabla f(x_k). \quad (\text{GD})$$

Under (H_0) , and for $(h_k)_{k \in \mathbb{N}} \subset]0, 2/L[$, then we have both the convergence of the values $f(x_k) - \min f = \mathcal{O}(1/k)$ (in fact even $o(1/k)$), and the weak convergence of iterates $(x_k)_{k \in \mathbb{N}}$ to a point in $\arg \min(f)$. This convergence rate can be refined under various additional geometrical properties on the objective f such as error bounds (and the closely related Kurdyka-Łojasiewicz property in the convex case; Bolte et al. 2016).

1.1.2. Second-Order In-Time Systems: Key Role of Inertia. Second-order in-time inertial dynamical systems have been introduced to provably accelerate the convergence behavior compared with (GF). An abundant literature has been devoted to the study of inertial dynamics:

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

The importance of working with an asymptotically vanishing viscosity coefficient $\gamma(t)$ to obtain acceleration was stressed by several authors (Cabot et al. 2009, Attouch and Cabot 2017). Most of the literature focuses on the case $\gamma(t) = \alpha/t$, originating from the seminal work of Su et al. (2016) who showed the rate of convergence $\mathcal{O}(1/t^2)$ of the values for $\alpha = 3$, thus making the link with the Nesterov accelerated gradient method (Nesterov 1983). Since then, an important body of literature has been devoted to this important case and the subtle tuning of parameter α . Taking $\alpha \geq 3$ is mandatory for provable accelerated rate $\mathcal{O}(1/t^2)$ of the values (Attouch et al. 2018b), and $\alpha > 3$ provides an even better rate $o(1/t^2)$ and weak convergence of the trajectory (Attouch and Peypouquet 2016, May 2017). On the other hand, $\alpha < 3$ necessarily leads to a slower rate $\mathcal{O}(1/t^{2\alpha/3})$ (Apidopoulos et al. 2018, Attouch et al. 2019). Another remarkable instance of (IGS_γ) corresponds to the well-known heavy ball with friction (HBF) method, where $\gamma(t)$ is a constant, first introduced by Polyak (1964). When the objective is strongly convex, it was shown that the trajectory converges exponentially with an optimal convergence rate for a properly chosen constant γ (Attouch et al. 2011).

1.1.3. Geometric Hessian-Driven Damping. Because of the inertial aspects, and the asymptotic vanishing viscous damping coefficient, (IGS_γ) may exhibit many small oscillations that are not desirable from an optimization point of view. To remedy this, a powerful tool consists in introducing a geometric damping driven by the Hessian of f into the dynamic. This yields the inertial system with explicit Hessian-driven damping (Attouch et al. 2016, 2022b):

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

and the inertial system with implicit Hessian-driven damping (Alecsa et al. 2021, Muehlebach and Jordan 2021):

$$\ddot{x}(t) + \gamma(t)\dot{x}(t) + \nabla f(x(t) + \beta(t)\dot{x}(t)) = 0. \quad (\text{ISIHD})$$

The rationale behind the use of the term “implicit” in (ISIHD) comes from a Taylor expansion of the gradient term (as $t \rightarrow +\infty$ we expect $\dot{x}(t) \rightarrow 0$). Following the physical interpretation of these Ordinary Differential Equations, we call the nonnegative parameters γ and β as the viscous and geometric damping parameters, respectively.

At first glance, the presence of the Hessian in (ISEHD) may seem to entail numerical difficulties. However, this is not the case because the term $\nabla^2 f(x(t))\dot{x}(t)$ is nothing but the time derivative of $t \mapsto \nabla f(x(t))$. This explains why the time discretization of this dynamic provides efficient first-order algorithms (Attouch et al. 2022b). On the other hand, (ISEHD) can be argued to be truly of second-order nature, that is, close to Newton’s and Levenberg-Marquardt’s dynamics (Castera et al. 2024). This understanding suggests that (ISIHD) may reflect the nature of first-order algorithms more faithfully than (ISEHD). Moreover, when it will come to our stochastic setting, which is the focus of this paper, the approach we propose will only make sense for the implicit form of the Hessian-driven damping. Indeed, in our

stochastic setting, we do not have direct access to the evaluation of the gradient of f . Instead, we model the associated errors with a continuous Itô martingale (denoted as $M(t)$). But for the explicit form of the Hessian-driven damping, this would entail a term involving the time derivative of $\nabla f(X(t)) + M(t)$. This is meaningless because (nonconstant) martingales are not differentiable almost sure (a.s.) This is why from now on, we will solely consider (ISIH).

1.2. Motivations

In the following, \mathbb{H}, \mathbb{K} are real separable Hilbert spaces. In many practical situations, the (sub)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 (sub)gradient are necessary. These errors can also be due to some other exogenous factor. The continuous-time approach through stochastic differential equations (SDEs) 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 SDE perspective to model Stochastic Gradient Descent (SGD)-type algorithms motivated by various reasons; (Li et al. 2017, 2021; Mertikopoulos and Staudigl 2018; Soatto and Chaudhari 2018; Hu et al. 2019b; Orvieto and Lucchi 2019; Shi et al. 2023; Dambrine et al. 2024; Maulen-Soto et al. 2024, 2025). 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 context, the first SDE that comes to mind is

$$\begin{cases} dX(t) = -\nabla f(X(t)) + \sigma(t, X(t))dW(t), & t > t_0; \\ X(t_0) = X_0, \end{cases} \quad (\text{SGF})$$

defined over a complete filtered probability space $(\Omega, \mathcal{F}, (\mathcal{F}_t)_{t \geq t_0}, \mathbb{P})$, where the diffusion (volatility) term $\sigma : \mathbb{R}^+ \times \mathbb{H} \rightarrow \mathcal{L}_2(\mathbb{K}; \mathbb{H})$ is a measurable function, W is a \mathcal{F}_t -adapted \mathbb{K} -valued cylindrical Brownian motion, and the initial data X_0 is an \mathcal{F}_{t_0} -measurable \mathbb{H} -valued random variable. All these notations and concepts are introduced in Section 2. (SGF) is a stochastic counterpart of (GF) where the error is modeled as a stochastic integral with respect to the measure defined by a continuous Itô martingale.

1.2.1. SDE Modeling of SGD. To simplify the discussion, let us focus in this section on the finite-dimensional case ($\mathbb{H} = \mathbb{R}^d, \mathbb{K} = \mathbb{R}^m$). In various areas of science and engineering, in particular, in machine learning, the SGD is a powerful alternative to gradient descent and consists in replacing the full gradient computation by a cheap random version, serving as an unbiased estimator. The SGD updates the iterates according to

$$x_{k+1} = x_k - h(\nabla f(x_k) + e_k), \quad (\text{SGD})$$

where $h \in \mathbb{R}_+$ is the stepsize and e_k is the random noise term on the gradient at the k th iteration. As such, (SGD) can be viewed as instance of the Robbins-Monro stochastic approximation algorithm (Robins and Monro 1951). When the objective takes the form $f(x) = \mathbb{E}[\hat{f}(x, \xi)]$, where the expectation is with respect to the random variable ξ , the single-batch version of SGD reads

$$x_{k+1} = x_k - h\nabla \hat{f}(x_k, \xi_k), \quad (\text{SGD}_{\text{SB}})$$

where $(\xi_k)_{k \in \mathbb{N}}$ are independent and identically distributed random variables with the same distribution as ξ . Of course, (SGD_{SB}) is an instance of (SGD) with $e_k = \nabla \hat{f}(x_k, \xi_k) - \nabla f(x_k)$.

The SDE continuous-time approach is motivated by its close relation to (SGD) or (SGD_{SB}). We first note that when the noise e_k in (SGD) is $\mathcal{N}(0, \sigma_k I_d)$, (SGF) is a better continuous-time model for (SGD) than (GF), as has been shown recently in Dambrine et al. (2024, proposition 2.1). There, the sequence $(x_k)_{k \in \mathbb{N}}$ provided by (SGD), with $e_k \sim \mathcal{N}(0, \sigma_k I_d)$, was proved to be accurately approximated by (SGF) with $\sigma(t, X(t)) = \sqrt{h}\sigma(t)$ and $\sigma(kh) = \sigma_k I_d$.

For the standard single-batch SGD (SGD_{SB}), the argument is more involved. Actually, many recent works (Mandt et al. 2016; Li et al. 2017, 2021, 2024; Soatto and Chaudhari 2018; Hu et al. 2019b; Orvieto and Lucchi 2019; Latz 2021; Xie et al. 2021; Shi et al. 2023) have linked algorithm (SGD_{SB}) with continuous-time first-order stochastic diffusion dynamics such as (SGF). These works show either empirically or theoretically under which conditions (appropriate drift and diffusion terms, regularity of f , etc.) (SGF) can be seen as a good approximation model of (SGD_{SB}) for fixed stepsize. By a good model, we mean that (SGF) is a continuum limit of (SGD_{SB}) as the stepsize h goes to zero, or equivalently that the approximation of (SGD_{SB}) via the diffusion process (SGF) is precise in some weak sense (Li et al. 2017, 2021; Hu et al. 2019b).

As a consequence, using (SGF) as a proxy of (SGD) or (SGD_{SB}) allows to capitalize on the wealth of results in the field of SDEs, Itô calculus, and measure theory, and this in turn opens the door to new insights in the

behavior of (SGD) or (SGD_{SB}) and to transfer all the convergence results that one can prove for (SGF) to (SGD). Actually, this is one of the main messages we want to convey in this work. Our motivation and results are also complementary to those in the literature. Indeed, most, if not all, of works cited in the previous paragraph are primarily motivated by the fact that continuous-time SDE approximation of (SGD) is a crucial tool to study its escape behavior of bad saddle points (a.k.a. traps) in the nonconvex smooth case. Our standpoint, which is line with Maulen-Soto et al. (2024), Mertikopoulos and Staudigl (2018), Dambrine et al. (2024), and Maulen-Soto et al. (2025), is complementary, and we argue that the continuous-time perspective offers a deep insight and unveils the key properties of the dynamic, without being tied to a specific discretization. This enlightens the behavior of the sequence generated by some specific algorithm. In turn, studying the continuous-time SDE will allow to predict the convergence behavior of stochastic algorithms seen as a discretization of the corresponding continuous-time dynamics.

1.2.2. Second-Order SDE Modeling of Inertial SGD. Using a lifting argument to get an equivalent first-order reformulation, a natural generalization of (ISIH_D) to the nonsmooth case yields the differential inclusion

$$\begin{cases} \dot{x}(t) = v(t), & t > t_0; \\ \dot{v}(t) \in -[\gamma(t)v(t) + \partial F(x(t) + \beta(t)v(t))], & t > t_0; \\ x(t_0) = x_0, \quad \dot{x}(t_0) = v_0, \end{cases} \quad (\text{ISIH}_{\text{NS}})$$

where ∂F is the convex subdifferential of F . In this setting, keeping in mind that we want to give a rigorous meaning to (ISIH_{NS}), 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 inclusion (SDI), which is the stochastic counterpart of (ISIH_{NS}):

$$\begin{cases} dX(t) = V(t)dt, \\ dV(t) \in -\gamma(t)V(t)dt - \partial F(X(t) + \beta(t)V(t))dt + \sigma(t, X(t) + \beta(t)V(t))dW(t), \\ X(t_0) = X_0, \quad V(t_0) = V_0. \end{cases} \quad (\text{S-ISIH}_{\text{NS}})$$

This SDI is defined on a filtered probability space $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$, where $\sigma : [t_0, +\infty[\times \mathbb{H} \rightarrow \mathcal{L}_2(\mathbb{K}; \mathbb{H})$ is a measurable function; and W is a \mathbb{K} -valued cylindrical Brownian motion. When $g \equiv 0$, this reduces to the stochastic counterpart of (ISIH_D), given by the following SDE:

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

In the smooth but nonconvex case, and motivated again by strict saddle points avoidance, Hu et al. (2019a, b) study the convergence behavior of a stochastic discrete heavy ball method from its approximating SDE. The latter is a randomly perturbed nonlinear oscillator in Hu et al. (2019c) and a coupled system of nonlinear oscillators in Hu et al. (2019a). These SDEs are very different from the ones considered here. By contrast, other works do not employ SDE techniques and instead analyze the stochastic iteration directly; see, for instance, Barakat et al. (2021), Gadat et al. (2018), and Le (2024), with the last reference also addressing the nonsmooth case. There are several advantages of adopting the SDE perspective in our setting, among which are a more accurate approximation of the discrete stochastic dynamics than the ODE method and the ability to quantify stochastic fluctuations and describe the long-run distribution of the iterates; the approximation aspect is detailed in Appendix A.

Extending what we have discussed above for (SGF) as a good model of (SGD), we will show in Proposition A.1 (see Appendix A for details) that (S-ISIH_D) is a good model of the natural stochastic inertial algorithm (A.1) obtained by simple Euler forward (a.k.a. Euler-Maruyama) discretization of (S-ISIH_D). The corresponding convergence rate as a function of the time stepsize h is $\mathcal{O}(h)$ and shows that this is much better than the approximation rate of (ISIH_D), which is only $\mathcal{O}(\sqrt{h})$. As a consequence, this justifies and motivates the continuous-time dynamics (S-ISIH_D) (and (S-ISIH_{NS})) as a good proxy of stochastic inertial algorithms and opens the door to new insights in the behavior of such algorithms, and that their convergence properties can be easily derived from those of (S-ISIH_D) with minimal effort. For instance, when f is smooth with Lipschitz-continuous gradient, it is easy to see from the descent lemma that

$$\mathbb{E}[f(X_k) - \min f] = \mathbb{E}[f(X(kh)) - \min f] + \mathcal{O}(h),$$

where X_k are the iterates of (A.1) and $X(t)$ the solution to (S-ISIHD). This means that any rate proved on $\mathbb{E}[f(X(t)) - \min f]$ can be directly transferred to $\mathbb{E}[f(X_k) - \min f]$.

1.3. Objectives and Contributions

In this work, our goal is to provide a general mathematical framework for analyzing the convergence properties of (S-ISIHD) and (S-ISIHD_{NS}). In this context, considering inertial dynamics with a time-dependent vanishing viscosity coefficient γ is a key ingredient to obtain fast convergent methods. We will develop a systematic and unified way that transfers the properties of stochastic first-order dynamics recently studied by Maulen-Soto et al. (2024, 2025) to second-order ones. Our program will then rely on two pillars: *time scaling* and *averaging*, following the methodology recently developed in Attouch et al. (2024) for the deterministic case.

More precisely, we study the stochastic dynamics (S-ISIHD_{NS}) and its long-time behavior in order to solve (P). We conduct a new analysis using specific and careful arguments that are much more involved than in the deterministic case. To get some intuition, keeping the discussion informal at this stage, let us first identify the assumptions needed to expect that the position state of (S-ISIHD) “approaches” $\arg \min(f)$ in the long run. In the case where $\mathbb{H} = \mathbb{K}$, $\gamma(\cdot) \equiv \gamma > 0, \beta \equiv 0$, and $\sigma = \tilde{\sigma}I_{\mathbb{H}}$, where $\tilde{\sigma}$ is a positive real constant; under mild assumptions one can show that (S-ISIHD) has a unique invariant distribution $\pi_{\tilde{\sigma}}$ in (x, v) with density proportional to

$$\exp\left(-\frac{2\gamma}{\tilde{\sigma}^2}\left(f(x) + \frac{v^2}{2}\right)\right)$$

(Pavliotis 2014, proposition 6.1). Clearly, as $\tilde{\sigma} \rightarrow 0^+$, $\pi_{\tilde{\sigma}}$ gets concentrated around $\arg \min(f) \times \{0_{\mathbb{H}}\}$, with $\lim_{\tilde{\sigma} \rightarrow 0^+} \pi_{\tilde{\sigma}}(\arg \min(f) \times \{0_{\mathbb{H}}\}) = 1$; see also Section 1.4 for further discussion. 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 main contributions are summarized as follows:

- We show almost sure weak convergence of the trajectory (see Theorem 5) and convergence rates (see Theorem 6) in expectation in the case of time-dependent coefficients $\gamma(t)$ and a proper choice of $\beta(t)$ (depending on γ). For this analysis, we transfer the results from the Lyapunov analysis of the first-order in-time stochastic (sub)gradient system studied in Maulen-Soto et al. (2024, 2025) from which our inertial system is built through time scaling and averaging.
- We obtain almost sure and ergodic convergence results which correspond precisely to the best-known results in the deterministic case. In particular, if we let $\alpha > 3$, $\gamma(t) = \alpha/t$, $\beta(t) = t/(\alpha - 1)$, then under appropriate assumptions on the diffusion (volatility) term σ , we obtain the rate of convergence $o(1/t^2)$ of the values as well as fast convergence of the gradient in almost sure sense (see Corollary 1). This corresponds to the best known results of second-order in-time Hessian-driven damping dynamics in the deterministic case.
- As far as the interplay between viscous damping, geometric damping, and noise level, our results reveal that viscous damping, which allows for acceleration, enters in the control of the noise and may impose a more stringent summability condition on the diffusion coefficient (see, e.g., Theorem 5, Theorem 6 and Corollary 1). The geometric damping, which is designed to reduce oscillations, does not enter directly in this control and can be chosen in a more flexible way.
- We then turn to providing a local analysis with a local linear convergence rate under the Polyak-Łojasiewicz inequality (see Theorem 8). This is much more challenging in the stochastic case, and even more for second-order systems, as localizing the process in this case is very delicate. Leveraging time scaling and averaging offers an elegant framework to achieve such a local convergence analysis while it is barely possible otherwise. To the best of our knowledge, this is the first result of this kind for these systems in the literature.
- We also show almost sure strong convergence of the trajectory to the minimal norm solution when adding a Tikhonov regularization to our systems (see Theorem 9). Moreover, we show convergence rates in expectation for the objective and the trajectory for a particular Tikhonov regularizer (see Theorem 11).

It is worth observing that, because our approach is based on an averaging technique, it will involve Jensen’s inequality at some point to get fast convergence rates. In this respect, the convexity assumption on the objective function appears unavoidable, at least in this proof strategy. It is also worth stressing the fact that the approach only makes sense for the implicit form of the Hessian-driven damping. Indeed, as explained above, the explicit form of the Hessian-driven damping has a term involving the time derivative of the (sub)gradient at the trajectory. As the noise, modeled here as an Itô martingale, in practice stems from the (sub)gradient evaluation, this time derivative is meaningless with explicit Hessian-driven damping, as (nonconstant) martingales are a.s. nondifferentiable.

1.4. Relation to Prior Work

1.4.1. Kinetic Diffusion Dynamics for Sampling. Let us consider (S-ISIHD) in the case where $\mathbb{H} = \mathbb{K} = \mathbb{R}^n$, $\gamma(t) = \gamma > 0, \beta \equiv 0$, and $\sigma = \sqrt{2\gamma}I$. 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 that has the density $\propto \exp(-f(x) - \|v\|^2/2)$ with respect to the Lebesgue measure on \mathbb{R}^{2n} . 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 nonasymptotic convergence guarantees in various topologies and under various conditions; see Cheng et al. (2018), Ma et al. (2021), and Dalalyan et al. (2019) and references therein. By rescaling the problem, relation between sampling and optimization with quantitative estimates has been investigated in, for example, Dalalyan and Karagulyan (2019) for the strongly convex case.

1.4.2. First-Order Stochastic (Sub)Gradient Systems. In Maulen-Soto et al. (2024) (respectively, Maulen-Soto et al. (2025)), the authors studied (SGF) (respectively, its nonsmooth version as an SDI) as a proxy for (SGD). One of their main results is the almost sure weak convergence of the trajectory to the set of minimizers under integrability conditions on σ , as well as convergence rates in expectation and in almost sure sense of the dynamic under different geometries of f . Our goal here is to take these results to second-order inertial dynamics featuring both viscous and geometric dampings. This turns out to be a challenging task. We will show that the convergence rate on the objective can be achieved provided that the noise vanishes sufficiently fast.

1.4.3. Inexact Inertial Gradient Systems. There is an abundant literature regarding the dynamics of (ISIHD) and (ISEHD), either in the exact case or with errors but only deterministic ones (Schmidt et al. 2011; Haraux and Jendoubi 2012; Villa et al. 2013; Dossal and Aujol 2015; Attouch et al. 2018a, b, 2022a, b, 2023a, b; Alecsa et al. 2021; Castera et al. 2021; Shi et al. 2022). Only a few papers have been devoted to studying stochastic versions of (IGS $_{\gamma}$), either with vanishing damping $\gamma(t) = \alpha/t$ or constant damping $\gamma(t)$ (stochastic HBF) (Gadat and Panloup 2014, Gadat et al. 2018, Orvieto et al. 2020, Dambrine et al. 2024). One of the advantages of considering the stochastic version of (ISIHD) is the possible reduction of oscillations thanks to the introduction of the geometric damping. We show that this effect could be preserved in the stochastic setting, hence obtaining more stable trajectory solutions than other stochastic second-order dynamics. Additionally, having the flexibility to choose the viscous damping would eventually allow us to achieve faster convergence results. For a stochastic version of (IGS $_{\gamma}$), Dambrine et al. (2024) 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. These geometrical assumptions come in the form of *global* growth and flatness of the objective that is very restrictive. Rather, here, our geometrical assumptions on f will be only local. The corresponding stochastic algorithms for these two choices of γ , whose mathematical formulation and analysis are simpler, have been the subject of active research work (Frostig et al. 2015, Allen-Zhu 2017, Jain et al. 2017, Lin et al. 2017, Gadat et al. 2018, Yan 2018, Assran and Rabbat 2020, Laborde and Oberman 2020, Lan 2020, Loizou and Richtárik 2020, Defazio and Jelassi 2022, Driggs et al. 2022, Attouch et al. 2024, Hamadouche et al. 2024).

1.4.4. Time Scaling and Averaging. Attouch et al. (2024) 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). Our aim is to extend the results of Attouch et al. (2024) to the stochastic case. Leveraging these techniques with a general function γ and an appropriate β , we will be able to transfer all the results we obtained in Maulen-Soto et al. (2025) for a first-order SDI to the second-order one (S-ISIHD $_{\text{NS}}$). This avoids in particular to go through an intricate and a dedicated Lyapunov analysis for (S-ISIHD $_{\text{NS}}$). A local convergence analysis becomes also easily accessible through this lens, whereas it is barely possible otherwise. We also specialize our results to the standard case where $\gamma(t) = \alpha/t$ and $\beta(t) = t/(\alpha - 1)$.

The idea of passing from a first-order system to a second-order one via time scaling is not new. In the smooth case ($g \equiv 0$), Cabot (2009) proposes time scaling and a tricky change of variables to show that (IGS $_{\gamma}$) is equivalent to an averaged gradient system, that is, the steepest gradient system (GF) where the instantaneous value of $\nabla f(x(t))$ is replaced by some average of the gradients $\nabla f(x(s))$ over all past positions $s \leq t$. See also Goudou and Munier (2005) for more general gradient systems with memory terms involving kernels. This gives rise to an integro-differential equation whose asymptotic behavior and the equivalent second-order dynamic have been investigated in Cabot (2009). A stochastic version of this integro-differential equation has been studied in Gadat and Panloup (2014), where the long time behavior of the resulting process, in particular its invariant distribution and occupation measure, was investigated under ellipticity assumptions on f and σ and proper behavior of the averaging gradient function. Clearly, the motivation of that work is not on the minimizing properties of the process, whereas it is our focus here. We remark that after simple integration on

(S-ISIHD), we can get that

$$\begin{cases} dX(t) = V(t)dt, \\ V(t) = \frac{V_0}{p(t)} - \frac{1}{p(t)} \int_{t_0}^t p(s) \nabla f(X(s) + \beta(s)V(s)) ds \\ \quad + \frac{1}{p(t)} \int_{t_0}^t p(s) \sigma(s, X(s) + \beta(s)V(s)) dW(s), \end{cases}$$

where $p(t) = \exp(\int_{t_0}^t \gamma(s) ds)$. This representation highlights the link with integro-differential equations.

1.5. Organization of the Paper

Section 2 introduces notations, definitions, and preliminaries that are essential to our exposition. Section 3 is the main part of our study. We develop the passage from the first-order system to the second-order inertial system by using the time scaling and averaging in a stochastic framework. Almost sure and ergodic convergence rates are provided under different geometric properties of the objective function, such as convexity and Polyak-Łojasiewicz geometry. Finally, we show a strong convergence result when adding a Tikhonov regularization. Additional technical results that are needed throughout the paper are gathered in Appendix B.

2. Notation and Preliminaries

2.1. Notation

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 will omit the subscripts \mathbb{H} and \mathbb{K} whenever it is clear from the context); $I_{\mathbb{H}}$ is the identity operator from \mathbb{H} to \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 $M \in \mathcal{L}_1(\mathbb{K})$, the trace is defined by

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

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 $w\text{-lim}$ (respectively, $s\text{-lim}$) the limit for the weak (respectively, strong) 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} . For $f : \mathbb{H} \rightarrow \mathbb{R}$, 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 $p = +\infty$. 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 Rockafellar (1997).

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, that is, functions f such that $f - \mu \|\cdot\|^2/2$ is convex. 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 *subdifferential* of a function $f \in \Gamma_0(\mathbb{H})$ is the set-valued operator $\partial f : \mathbb{H} \rightrightarrows \mathbb{H}$ such that, for every x in \mathbb{H} ,

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

which is nonempty for every point in the relative interior of the domain of f . When f is finite-valued, then f is continuous, and $\partial f(x)$ is a nonempty convex and compact set for every $x \in \mathbb{H}$. If f is differentiable, then $\partial f(x) = \{\nabla f(x)\}$. For every $x \in \mathbb{H}$ such that $\partial f(x) \neq \emptyset$, the minimum norm selection of $\partial f(x)$ is the unique element $\{\partial^0 f(x)\} \stackrel{\text{def}}{=} \arg \min_{u \in \partial f(x)} \|u\|$. The projection of a point $x \in \mathbb{H}$ onto a nonempty closed convex set $C \subseteq \mathbb{H}$ is denoted by $P_C(x)$.

2.2. Assumptions on Volatility and Damping Parameters

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 our SDEs and SDIs (in particular (S-ISIHD) and (S-ISIHD)_{NS}) 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 classical and will be required to ensure the well posedness of (S-ISIHD) and (S-ISIHD)_{NS}. Let us also define $\sigma_\infty : [t_0, +\infty[\rightarrow \mathbb{R}_+$ as

$$\sigma_\infty(t) \stackrel{\text{def}}{=} \sup_{x \in \mathbb{H}} \|\sigma(t, x)\|_{\text{HS}}.$$

Remark 1. (H_σ) implies the existence of $\sigma_* > 0$ such that

$$\|\sigma(t, x)\|_{\text{HS}}^2 = \text{tr}[\Sigma(t, x)] \leq \sigma_*^2,$$

for all $t \geq t_0, x \in \mathbb{H}$, where $\Sigma \stackrel{\text{def}}{=} \sigma \sigma^*$.

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

$$\begin{aligned} p(t) &\stackrel{\text{def}}{=} \exp\left(\int_{t_0}^t \gamma(s) ds\right), \quad \Gamma(t) \stackrel{\text{def}}{=} p(t) \int_t^{+\infty} \frac{ds}{p(s)}, \\ I[h](t) &\stackrel{\text{def}}{=} \exp\left(-\int_{t_0}^t \frac{du}{\Gamma(u)}\right) \int_{t_0}^t h(u) \frac{\exp\left(\int_{t_0}^u \frac{ds}{\Gamma(s)}\right)}{\Gamma(u)} du. \end{aligned}$$

We assume that

$$\begin{cases} \gamma \text{ is upper bounded by a non-increasing function for every } t \geq t_0; \\ \int_{t_0}^{+\infty} \frac{ds}{p(s)} < +\infty. \end{cases} \quad (\text{H}_\gamma)$$

Remark 2. Let us notice that Γ satisfies the relation $\Gamma' = \gamma\Gamma - 1$.

2.3. Reminder of Main Previous Results

Before we delve into our core contributions, it is important to note that we will require some specific results gleaned from Maulen-Soto et al. (2024, 2025). These are the subject of the following paragraphs.

2.3.1. Results on (SGF). Because we are going to show results in the smooth case, we rewrite (H₀) when $g \equiv 0$:

$$\begin{cases} f : \mathbb{H} \rightarrow \mathbb{R} \text{ is continuously differentiable and convex with } L\text{-Lipschitz} \\ \text{continuous gradient;} \\ \mathcal{S} \stackrel{\text{def}}{=} \arg \min(f) \neq \emptyset. \end{cases} \quad (\text{H}_0)$$

Theorem 1 (Maulen-Soto et al. 2024, Theorem 3.1). *Consider f and σ that satisfy Assumptions (H₀) and (H_σ). Let $\nu \geq 2$ and consider the SDE*

$$\begin{cases} dX(t) = -\nabla f(X(t))dt + \sigma(t, X(t))dW(t), \\ X(t_0) = X_0, \end{cases} \quad (1)$$

where $X_0 \in L^\nu(\Omega; \mathbb{H})$. Then, there exists a unique solution $X \in S_{\mathbb{H}}^\nu[t_0]$ (see Section 7.2 for the notation) of (1). Additionally, if $\sigma_\infty \in L^2([t_0, +\infty[)$, then

- i. $\sup_{t \geq t_0} \mathbb{E}[\|X(t)\|^2] < +\infty$.
- ii. $\forall x^* \in \mathcal{S}$, $\lim_{t \rightarrow +\infty} \|X(t) - x^*\|$ exists a.s. and $\sup_{t \geq t_0} \|X(t)\| < +\infty$ a.s.
- iii. $\lim_{t \rightarrow +\infty} \|\nabla f(X(t))\| = 0$ a.s. As a result, $\lim_{t \rightarrow +\infty} f(X(t)) = \min f$ a.s.
- iv. There exists an \mathcal{S} -valued random variable X^* such that $w - \lim_{t \rightarrow +\infty} X(t) = X^*$ a.s.

Theorem 2 (Maulen-Soto et al. 2024, Theorem 3.4). *Let $v \geq 2$ and consider the SDE (1) with initial data $X_0 \in L^v(\Omega; \mathbb{H})$, where f and σ satisfy Assumptions (H_0) and (H_σ) . Moreover, we assume that σ satisfies $t \mapsto t\sigma_\infty^2(t) \in L^1([t_0, +\infty[)$ and that either \mathbb{H} is finite-dimensional or $f \in C^2(\mathbb{H})$. Then, the solution trajectory $X \in S_{\mathbb{H}}^v[t_0]$ is unique and we have that*

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

Moreover, if $f \in C^2(\mathbb{H})$, then the following hold:

ii. $t \mapsto t\|\nabla f(X(t))\|^2 \in L^1([t_0, +\infty[)$ a.s.

iii. $f(X(t)) - \min f = o(t^{-1})$ a.s.

2.3.2. Results on (SGI). The far more intricate nonsmooth version of (SGF) has been recently studied in Maulen-Soto et al. (2025). Let F, σ satisfy (H_0) and (H_σ) . We consider the stochastic differential inclusion

$$\begin{cases} dX(t) \in -\partial F(X(t))dt + \sigma(t, X(t))dW(t), & t > t_0, \\ X(t_0) = X_0. \end{cases} \quad (\text{SGI})$$

The following definition makes precise the notion of solution that we are interested in.

Definition 1. A solution of (SGI) is a couple (X, η) of \mathcal{F}_t -adapted processes such that almost surely

i. X is continuous with sample paths in the domain of ∂g .

ii. $\eta : [t_0, +\infty[\rightarrow \mathbb{H}$ is absolutely continuous, such that $\eta(t_0) = 0$, and $\forall T > t_0, \eta' \in L^2([t_0, T]; \mathbb{H}), \eta'(t) \in \partial g(X(t))$ for almost all $t \geq t_0$;

iii. For $t > t_0$,

$$\begin{cases} X(t) = X_0 - \int_{t_0}^t \nabla f(X(s))ds - \eta(t) + \int_{t_0}^t \sigma(s, X(s))dW(s), \\ X(t_0) = X_0. \end{cases} \quad (2)$$

For brevity, we sometimes omit the process η and say that X is a solution of (SGI), meaning that, there exists a process η such that (X, η) satisfies the previous definition. The definition of uniqueness for the process X is given in Section 7.2.

In order to show the main results for (SGI), we consider the sequence of solutions $(X_\lambda)_{\lambda>0}$ of the SDEs

$$\begin{cases} dX_\lambda(t) = -\nabla(f + g_\lambda)(X_\lambda(t))dt + \sigma(t, X_\lambda(t))dW(t), & t > t_0, \\ X_\lambda(t_0) = X_0, \end{cases} \quad (\text{SDE}_\lambda)$$

where g_λ is the Moreau envelope of g with parameter $\lambda > 0$. Observe that system (SDE λ) was also directly studied in Maulen-Soto et al. (2024) to solve (P), where a complete discussion is provided. It is worth emphasizing that the Moreau envelope of g , and the corresponding system (SDE λ) is only needed as a means to establish existence and uniqueness of solution of (SGI). This follows the same reasoning as in the deterministic case where the Moreau-Yosida regularization and nonlinear semigroup theory have proven very useful to show existence and uniqueness for differential inclusions Brézis (1973). However, in the convergence analysis, we never replace g by its Moreau envelope.

We define the integrability condition that for every $T > t_0$,

$$\limsup_{\lambda \downarrow 0} \int_{t_0}^T \mathbb{E}(\|\nabla g_\lambda(X_\lambda(t))\|^2)dt < +\infty, \quad (\text{H}_\lambda)$$

Remark 3. Some conditions on g for (H_λ) to be satisfied are discussed in Pettersson (1995). Of particular interest, which covers a broad class of functions, is when ∂g has full domain and there exists $C_0 > 0$ such that

$$\|\partial^0 g(x)\| \leq C_0(1 + \|x\|), \quad \forall x \in \mathbb{H}.$$

This is, for instance, the case when g is Lipschitz continuous.

The following results on (SGI) were proved in Maulen-Soto et al. (2025).

Theorem 3 (Maulen-Soto et al. 2025). *Consider $F = f + g$ and σ satisfying (H_0) and (H_σ) . Suppose further that g verifies (H_λ) . Let $v \geq 2, t_0 \geq 0$, and consider the dynamic (SGI) with initial data $X_0 \in L^v(\Omega; \mathbb{H})$. Then, (SGI) has a unique solution in the sense of Definition 1 $(X, \eta) \in S_{\mathbb{H}}^v[t_0] \times C^1([t_0, +\infty[; \mathbb{H})$.*

Moreover, if $\sigma_\infty \in L^2([t_0, +\infty[)$, then the following holds:

i. $\mathbb{E}[\sup_{t \geq t_0} \|X(t)\|^v] < +\infty$.

ii. $\forall x^* \in S_F, \lim_{t \rightarrow +\infty} \|X(t) - x^*\|$ exists a.s. and $\sup_{t \geq t_0} \|X(t)\| < +\infty$ a.s.

iii. If g is continuous, then ∇f is constant on \mathcal{S}_F , there exists $x^* \in \mathcal{S}_F$ such that $s\text{-}\lim_{t \rightarrow +\infty} \nabla f(X(t)) = \nabla f(x^*)$ a.s., and

$$\int_{t_0}^{+\infty} (F(X(t)) - \min F) dt < +\infty \quad \text{a.s.}$$

iv. There exists an \mathcal{S}_F -valued random variable X^* such that $w\text{-}\lim_{t \rightarrow +\infty} X(t) = X^*$.

2.3.3. Tikhonov Regularization. Let us now turn to a Tikhonov regularization of (SGI), that is,

$$\begin{cases} dX(t) \in -\partial F(X(t)) - \varepsilon(t)X(t) + \sigma(t, X(t))dW(t), & t \geq t_0, \\ X(t_0) = X_0. \end{cases} \quad (\text{SGI-TA})$$

Solution existence and uniqueness for (SGI-TA) are established in Maulen-Soto et al. (2025, theorem 3.3). We also have the following.

Theorem 4 (Maulen-Soto et al. 2025, Theorem 4.1). *Let $\nu \geq 2$ and consider the dynamic (SGI-TA) with initial data $X_0 \in L^\nu(\Omega; \mathbb{H})$, where $F = f + g$ and σ satisfy Assumptions (\mathbf{H}_0) and (\mathbf{H}_σ) . Furthermore, assume that g satisfies (\mathbf{H}_λ) . Then, there exists a unique solution $X \in S_{\mathbb{H}}^\nu[t_0]$ of (SGI-TA). Let $x^* = P_{\mathcal{S}_F}(0)$ be the minimum norm solution, and for $\varepsilon > 0$ let x_ε be the unique minimizer of $F_\varepsilon(x) \stackrel{\text{def}}{=} F(x) + \varepsilon\|x\|^2/2$. Suppose that $\sigma_\infty \in L^2([t_0, +\infty[)$, and that $\varepsilon : [t_0, +\infty[\rightarrow \mathbb{R}_+$ satisfies the conditions:*

(T₁) $\varepsilon(t) \rightarrow 0$ as $t \rightarrow +\infty$;

(T₂) $\int_{t_0}^{+\infty} \varepsilon(t) dt = +\infty$;

(T₃) $\int_{t_0}^{+\infty} \varepsilon(t)(\|x^*\|^2 - \|x_{\varepsilon(t)}\|^2) dt < +\infty$.

Then, we have

i. $\sup_{t \geq t_0} \|X(t)\| < +\infty$ a.s., and

ii. $s\text{-}\lim_{t \rightarrow +\infty} X(t) = x^*$ a.s.

This means that we can obtain almost sure strong convergence of the trajectory to a particular (nonrandom) solution: the one of minimal norm.

3. From First-Order to Second-Order Systems

3.1. Time Scaling and Averaging

We apply a time scaling and then an averaging technique to the system (SGI) to derive an insightful reparametrization of a particular case of our second-order system ($\mathbf{S}\text{-ISIHD}_{\text{NS}}$), specifically, the case when $\beta \equiv \Gamma$. The main advantage of this method is that the results of (SGI) directly carry over to obtain results on the convergence behavior of ($\mathbf{S}\text{-ISIHD}_{\text{NS}}$) without passing through a dedicated Lyapunov analysis. However, as discussed in the introduction, the averaging technique is restricted to convex objectives, as it heavily relies on Jensen's inequality.

Let $\nu \geq 2$, $s_0 > 0$. We consider the potential $F = f + g$ where g satisfies (\mathbf{H}_λ) . Let σ_1 be a diffusion term in the time parametrization by s . We will study the dynamic (SGI) in s , starting at s_0 , with diffusion term σ_1 under Hypotheses (\mathbf{H}_0) and (\mathbf{H}_σ) . Let $\sigma_{1*} > 0$ be such that

$$\|\sigma_1(s, x)\|_{\text{HS}} \leq \sigma_{1*}^2, \quad \forall s \geq s_0, \forall x \in \mathbb{H},$$

and $\sigma_{1\infty}(s) \stackrel{\text{def}}{=} \sup_{x \in \mathbb{H}} \|\sigma_1(s, x)\|_{\text{HS}}$. We rewrite (SGI) in a format adapted to our case,

$$\begin{cases} dZ(s) \in -\partial F(Z(s))ds + \sigma_1(s, Z(s))dW(s), & s > s_0, \\ Z(s_0) = Z_0, \end{cases} \quad (3)$$

where $Z_0 \in L^\nu([s_0, +\infty[; \mathbb{H})$.

Let us make the change of time $s = \tau(t)$ in the dynamic (3), where τ is an increasing function from $[t_0, +\infty[$ to $[s_0, +\infty[$, which is twice differentiable and which satisfies $\lim_{t \rightarrow +\infty} \tau(t) = +\infty$. Denote $Y(t) \stackrel{\text{def}}{=} Z(s)$ and t_0 be such that $s_0 = \tau(t_0)$. By the chain rule and Öksendal 2003 (theorem 8.5.7), we have

$$\begin{cases} dY(t) \in -\tau'(t)\partial F(Y(t))dt + \sqrt{\tau'(t)}\sigma_1(\tau(t), Y(t))dW(t), & t > t_0, \\ Y(t_0) = Z_0. \end{cases} \quad (4)$$

Consider the smooth case, that is, when $g \equiv 0$ and the hypotheses of Theorem 2 ($f \in C_L^2(\mathbb{H})$ and $\sigma_1 \in L^2([s_0, +\infty[))$), then we can conclude that the convergence rate of (4) (when $g \equiv 0$) is the following:

$$f(Y(t)) - \min f = o\left(\frac{1}{\tau(t)}\right) \text{ a.s.} \quad (5)$$

By introducing a function τ that grows faster than the identity ($\tau(t) \geq t$), we have accelerated the dynamic, passing from the asymptotic convergence rate $1/s$ for (3) to $1/\tau(t)$ for (4). The price to pay is that the drift term in (4) is nonautonomous; furthermore, when the coefficient in front of the gradient tends to infinity as $t \rightarrow +\infty$, it will preclude the use of an explicit discretization in time. To overcome this, we adapt from Attouch et al. (2024) the following approach, which is called averaging.

Consider (4) and define the two stochastic processes $X, V : \Omega \times [t_0, +\infty[\rightarrow \mathbb{H}$ as

$$\begin{cases} dX(t) = V(t)dt, & t > t_0, \\ Y(t) = X(t) + \tau'(t)V(t), & t > t_0, \\ X(t_0) = X_0, & V(t_0) = V_0, \end{cases} \quad (6)$$

where $Y(t)$ is the process in (4), and $X_0, V_0 \in L^r(\Omega; \mathbb{H})$ are initial data. This leads us to set $Z_0 \stackrel{\text{def}}{=} X_0 + \tau'(t_0)V_0$ in order for the equations to fit. According to the averaging, the differential form of $Y(t)$ is

$$dY(t) = dX(t) + \tau''(t)V(t)dt + \tau'(t)dV(t).$$

Combining the previous equation with (4), we have that

$$-\tau'(t)\partial F(Y(t))dt + \sqrt{\tau'(t)}\sigma_1(\tau(t), Y(t))dW(t) \ni dX(t) + \tau''(t)V(t)dt + \tau'(t)dV(t).$$

Using that $dX(t) = V(t)dt$ and dividing by τ' , we then have

$$-\partial F(X(t) + \tau'(t)V(t))dt + \frac{1}{\sqrt{\tau'(t)}}\sigma_1(\tau(t), X(t) + \tau'(t)V(t))dW(t) \ni \frac{1 + \tau''(t)}{\tau'(t)}V(t)dt + dV(t).$$

Therefore, after the time scaling and averaging, we obtain the following dynamic:

$$\begin{cases} dX(t) = V(t)dt, & t > t_0, \\ dV(t) \in -\frac{1 + \tau''(t)}{\tau'(t)}V(t)dt - \partial F(X(t) + \tau'(t)V(t))dt \\ \quad + \frac{1}{\sqrt{\tau'(t)}}\sigma_1(\tau(t), X(t) + \tau'(t)V(t))dW(t), & t > t_0, \\ X(t_0) = X_0, & V(t_0) = V_0. \end{cases} \quad (\text{SIHD-S.1})$$

Let $\gamma : [t_0, +\infty[\rightarrow \mathbb{R}_+$ satisfy (H γ). We are going to determine τ in order to obtain a viscous damping coefficient equal to γ , that is,

$$\frac{1 + \tau''(t)}{\tau'(t)} = \gamma(t).$$

Clearly, τ' solves the following ODE in ζ :

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

As observed in Remark 2, the function Γ also satisfies the same ODE, and thus we can adjust the initial condition of τ' to obtain

$$\tau'(t) = \Gamma(t) = p(t) \int_t^{+\infty} \frac{du}{p(u)} \quad \forall t \geq t_0.$$

We then integrate and take $\tau(t) = s_0 + \int_{t_0}^t \Gamma(u)du$ to get $\tau(t_0) = s_0$ as required. This is a valid selection of τ because $t \mapsto s_0 + \int_{t_0}^t \Gamma(u)du$ is an increasing function from $[t_0, +\infty[$ to $[s_0, +\infty[$, twice differentiable and $\Gamma \notin L^1([t_0, +\infty[)$ because Γ is lower bounded by a nondecreasing function because γ is upper bounded by a nonincreasing function (Attouch and Cabot 2017, proposition 2.2) by (H γ). For this particular selection of τ , and defining

$\tilde{\sigma}_1(t, \cdot) \stackrel{\text{def}}{=} \sigma_1(\tau(t), \cdot) / \sqrt{\Gamma(t)}$, we have that (ISIHD-S.1) is equivalent to

$$\begin{cases} dX(t) = V(t)dt, & t > t_0, \\ dV(t) \in -\gamma(t)V(t)dt - \partial F(X(t) + \Gamma(t)V(t))dt \\ \quad + \tilde{\sigma}_1(t, X(t) + \Gamma(t)V(t))dW(t), & t > t_0, \\ X(t_0) = X_0, \quad V(t_0) = V_0. \end{cases} \quad (\text{ISIHD-S.2})$$

Clearly, (ISIHD-S.2) is nothing but (S-ISIHD_{NS}) when $\beta \equiv \Gamma$ and $\sigma \equiv \tilde{\sigma}_1$.

In order to be able to transfer the convergence results on Z in (3) (via (4)) to X in (ISIHD-S.2), it remains to express X in terms of Y only. For this, let

$$a(t) \stackrel{\text{def}}{=} \frac{1}{\tau'(t)}, \quad A(t) \stackrel{\text{def}}{=} \int_{t_0}^t a(u)du.$$

Recalling the averaging in (6), we need to integrate the following equation:

$$V(t) + a(t)X(t) = a(t)Y(t). \quad (7)$$

Multiplying both sides by $e^{A(t)}$ and using (6), we get equivalently

$$d(e^{A(t)}X(t)) = a(t)e^{A(t)}Y(t)dt. \quad (8)$$

Integrating and again using (6), we obtain

$$\begin{aligned} X(t) &= e^{-A(t)}X(t_0) + e^{-A(t)} \int_{t_0}^t a(u)e^{A(u)}Y(u)du \\ &= e^{-A(t)}Y(t_0) + e^{-A(t)} \int_{t_0}^t a(u)e^{A(u)}Y(u)du - e^{-A(t)}\tau'(t_0)V(t_0). \end{aligned}$$

Then we can write

$$X(t) = \int_{t_0}^t Y(u)d\mu_t(u) + \xi(t), \quad (9)$$

where μ_t is the probability measure on $[t_0, t]$ defined by

$$\mu_t \stackrel{\text{def}}{=} e^{-A(t)}\delta_{t_0} + a(u)e^{A(u)-A(t)}du, \quad (10)$$

where δ_{t_0} is the Dirac measure at t_0 , $a(u)e^{A(u)-A(t)}du$ is the measure with density $a(\cdot)e^{A(\cdot)-A(t)}$ with respect to the Lebesgue measure on $[t_0, t]$, and $\xi(t)$ is a random process since V_0 is a random variable, that is,

$$\xi(t) \stackrel{\text{def}}{=} \xi(\omega, t) = -e^{-A(t)}\tau'(t_0)V_0(\omega) \quad \forall \omega \in \Omega. \quad (11)$$

3.2. Convergence of the Trajectory and Convergence Rates Under General γ , and $\beta \equiv \Gamma$

We state the main results of this section. We first show almost sure convergence of the trajectory of (S-ISIHD_{NS}) to a random variable taking values in the set of minimizers of F . When $g \equiv 0$, we also provide convergence rates.

The following result imposes an integrability condition on the diffusion term, which is essential for ensuring the almost sure weak convergence of the trajectory. More precisely, this integrability condition allows to show that the trajectory asymptotically converges almost surely, in the weak topology, to a random variable that takes values in the set of minimizers of F .

Theorem 5. *Let $v \geq 2$ and consider the dynamic (S-ISIHD_{NS}) with initial data $X_0, V_0 \in L^v(\Omega; \mathbb{H})$, where $\gamma: [t_0, +\infty[\rightarrow \mathbb{R}_+$ satisfies (H _{γ}), and $\beta \equiv \Gamma$. Besides, $F = f + g$ and σ satisfy Assumptions (H₀) and (H _{σ}). Moreover, suppose that g satisfies (H _{λ}). Then, there exists a unique solution $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^v[t_0]$ of (S-ISIHD_{NS}). Additionally, if $\Gamma\sigma_\infty \in L^2([t_0, +\infty[)$, then there exists an \mathcal{S}_F -valued random variable X^* such that $w\text{-}\lim_{t \rightarrow +\infty} X(t) = X^*$ a.s. and $w\text{-}\lim_{t \rightarrow +\infty} \Gamma(t)V(t) = 0$ a.s.*

Proof. Let $\theta(t) \stackrel{\text{def}}{=} \int_{t_0}^t \Gamma(u) du$, $\tilde{\sigma}(s, \cdot) \stackrel{\text{def}}{=} \sigma(\theta^{-1}(s), \cdot) \sqrt{\Gamma(\theta^{-1}(s))}$, and $\tilde{\sigma}_\infty(s) \stackrel{\text{def}}{=} \sup_{x \in \mathbb{H}} \|\tilde{\sigma}(s, x)\|_{\text{HS}}$. Then $\Gamma\sigma_\infty \in L^2([t_0, +\infty[)$ is equivalent to $\tilde{\sigma}_\infty \in L^2(\mathbb{R}_+)$. Consider the dynamic:

$$\begin{cases} dZ(s) \in -\partial F(Z(s)) + \tilde{\sigma}(s, Z(s)) dW(s), & s > 0, \\ Z(0) = X_0 + \Gamma(t_0)V_0. \end{cases} \quad (12)$$

By Theorem 3, we have that there exists a unique solution $(Z, \eta) \in S_{\mathbb{H}}^v \times C^1(\mathbb{R}_+; \mathbb{H})$ of (12) and an \mathcal{S}_F -valued random variable X^* such that $w\text{-}\lim_{s \rightarrow +\infty} Z(s) = X^*$ a.s. Moreover, using the time scaling $\tau \equiv \theta$ and the averaging described in this section, we end up with the dynamic (S-ISIHD_{NS}) in the case where $\beta \equiv \Gamma$.

It is direct to check that the time scaling and averaging preserves the uniqueness of a solution $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^0[t_0]$. Now let us validate $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^v[t_0]$. Because

$$\mathbb{E} \left(\sup_{s \in [0, T]} \|Z(s)\|^v \right) < +\infty, \quad \forall T > 0,$$

we directly obtain

$$\mathbb{E} \left(\sup_{t \in [t_0, T]} \|Y(t)\|^v \right) < +\infty, \quad \forall T > t_0.$$

Thanks to Relation (9), the following holds:

$$\begin{aligned} \|X(t)\|^v &\leq v \left(\left\| X(t) - \int_{t_0}^t Y(u) d\mu_t(u) \right\|^v + \left\| \int_{t_0}^t Y(u) d\mu_t(u) \right\|^v \right) \\ &\leq v \left(\|\xi(t)\|^v + (t - t_0)^{v-1} \int_{t_0}^t \|Y(u)\|^v d\mu_t(u) \right). \end{aligned}$$

Let $T > t_0$ be arbitrary. Taking supremum over $[t_0, T]$ and then expectation at both sides, we obtain that

$$\mathbb{E} \left(\sup_{t \in [t_0, T]} \|X(t)\|^v \right) \leq v \left(\mathbb{E}(\|V_0\|^v) \|\Gamma(t_0)\|^v + (T - t_0)^{v-1} \mathbb{E} \left(\sup_{t \in [t_0, T]} \|Y(t)\|^v \right) \right) < +\infty.$$

Because $V(t) = \frac{Y(t) - X(t)}{\Gamma(t)}$, we have

$$\|V(t)\|^v \leq \frac{v}{\Gamma^v(t)} (\|Y(t)\|^v + \|X(t)\|^v).$$

Similar to before, we let $T > t_0$ be arbitrary, and take the supremum over $[t_0, T]$ and then expectation at both sides to obtain

$$\mathbb{E} \left(\sup_{t \in [t_0, T]} \|V(t)\|^v \right) \leq v \sup_{t \in [t_0, T]} \frac{1}{\Gamma^v(t)} \left(\mathbb{E} \left(\sup_{t \in [t_0, T]} (\|Y(t)\|^v + \|X(t)\|^v) \right) \right).$$

Because Γ is a continuous positive function, by the extreme value theorem, we have that there exists $t_T \in [t_0, T]$ such that $\sup_{t \in [t_0, T]} 1/\Gamma^v(t) = 1/\Gamma^v(t_T) < +\infty$, and we conclude that $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^v[t_0]$.

Now we prove that there exists an \mathcal{S}_F -valued random variable X^* such that $w\text{-}\lim_{t \rightarrow +\infty} X(t) = X^*$ a.s. By virtue of Theorem 3, there exists an \mathcal{S}_F -valued random variable X^* such that $w\text{-}\lim_{s \rightarrow +\infty} Z(s) = X^*$ a.s. We also notice that we have directly $w\text{-}\lim_{t \rightarrow +\infty} Y(t) = X^*$ a.s. Let $h \in \mathbb{H}$ be arbitrary and use Relation (9) as follows:

$$\begin{aligned} |\langle X(t) - X^*, h \rangle| &\leq \left| \left\langle X(t) - \int_{t_0}^t Y(u) d\mu_t(u), h \right\rangle \right| + \left| \left\langle \int_{t_0}^t Y(u) d\mu_t(u) - X^*, h \right\rangle \right| \\ &= |\langle \xi(t), h \rangle| + \left| \left\langle \int_{t_0}^t (Y(u) - X^*) d\mu_t(u), h \right\rangle \right| \\ &= |\langle \xi(t), h \rangle| + \left| \int_{t_0}^t \langle Y(u) - X^*, h \rangle d\mu_t(u) \right| \\ &\leq \|\xi(t)\| \|h\| + \int_{t_0}^t |\langle Y(u) - X^*, h \rangle| d\mu_t(u), \end{aligned}$$

where the second equality comes from the dominated convergence theorem, because $\sup_{s>t_0} \|Y(s)\| < +\infty$ a.s. (by (ii) of Theorem 3).

Now let $a(t) = 1/\Gamma(t)$ and $A(t) = \int_{t_0}^t du/\Gamma(u)$. By Lemma 3 (defined in Section 7.1) and the fact that $V_0 \in L^v(\Omega; \mathbb{H})$, we have that $\lim_{t \rightarrow +\infty} \|\xi(t)\| = 0$ a.s. On the other hand, it holds that

$$\begin{aligned} \int_{t_0}^t |\langle Y(u) - X^*, h \rangle| d\mu_t(u) &\leq e^{-A(t)} |\langle Y(t_0) - X^*, h \rangle| \\ &+ e^{-A(t)} \int_{t_0}^t a(u) e^{A(u)} |\langle Y(u) - X^*, h \rangle| du. \end{aligned}$$

Now let $b(u) = |\langle Y(u) - X^*, h \rangle|$. Because we already proved that $\lim_{u \rightarrow +\infty} b(u) = 0$ a.s., and we have that $a \notin L^1([t_0, +\infty[)$ by Lemma 3, we utilize Lemma 2 (also defined in Section 7.1) with our respective a, b functions. This let us conclude that

$$\lim_{t \rightarrow +\infty} |\langle X(t) - X^*, h \rangle| = 0 \quad \text{a.s.}$$

Thus, $w\text{-}\lim_{t \rightarrow +\infty} X(t) = X^*$ a.s. Finally, because

$$Y(t) = X(t) + \Gamma(t)V(t),$$

and that X and Y have (a.s.) the same limit, we conclude that

$$w\text{-}\lim_{t \rightarrow +\infty} \Gamma(t)V(t) = 0 \quad \text{a.s.} \quad \square$$

In the smooth case, we also have convergence rates on the objective value and the gradient. In particular, the following two theorems will provide general abstract convergence rates under the same integrability condition on the diffusion term. These results will be specialized to specific choice of the parameters in Section 3.3.

Theorem 6. Let $v \geq 2$ and consider the dynamic (S-ISIHD) with initial data $X_0, V_0 \in L^v(\Omega; \mathbb{H})$, such that f and σ satisfy (H_0) and (H_γ) , and in the case where γ satisfies $(H\gamma)$, $\beta \equiv \Gamma$. Moreover, suppose that either \mathbb{H} is finite dimensional or $f \in C^2(\mathbb{H})$, and

$$t \mapsto \sqrt{\theta(t)}\Gamma(t)\sigma_\infty(t) \in L^2([t_0, +\infty[),$$

where $\theta(t) \stackrel{\text{def}}{=} \int_{t_0}^t \Gamma(u)du$. Then the solution trajectory $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^v[t_0]$ is unique and satisfies

$$\mathbb{E}[f(X(t)) - \min f] = \mathcal{O}\left(\max\left\{e^{-A(t)}, I\left[\frac{1}{\theta}\right](t)\right\}\right), \quad \forall t > t_0,$$

where $A(t) \stackrel{\text{def}}{=} \int_{t_0}^t \frac{du}{\Gamma(u)}$ and we recall that $I\left[\frac{1}{\theta}\right](t) = e^{-A(t)} \int_{t_0}^t \frac{1}{\theta(u)} \frac{e^{A(u)}}{\Gamma(u)} du$.

From Hypothesis (H_A) , we have that $\lim_{t \rightarrow +\infty} e^{-A(t)} = 0$, and because $\Gamma \notin L^1([t_0, +\infty[)$, we can use Lemma 2 to check that $\lim_{t \rightarrow +\infty} I\left[\frac{1}{\theta}\right](t) = 0$.

Remark 4. If \mathbb{H} is finite-dimensional and g is a L_0 -Lipschitz function, we could obtain convergence rates for the nonsmooth setting (P) by considering (S-ISIHD) with potential $F_\rho \stackrel{\text{def}}{=} f + g_\rho$, where g_ρ is the Moreau envelope of g with parameter $\rho > 0$. With the same hypotheses on the initial conditions, σ, γ, β as in Theorem 6, and proceeding as in Maulen-Soto et al. (2024, section 5), we could obtain

$$\mathbb{E}[F(X_\rho(t)) - \min F] = \mathcal{O}\left(\max\left\{e^{-A(t)}, I\left[\frac{1}{\theta}\right](t)\right\}\right) + \frac{L_0^2}{2}\rho, \quad \forall t > t_0,$$

where $(X_\rho, V_\rho) \in S_{\mathbb{H} \times \mathbb{H}}^v[t_0]$ is the solution of (S-ISIHD) with potential F_ρ .

Proof. We will utilize the averaging technique used in Theorem 5 and Jensen's inequality. First, we have

$$\begin{aligned} \mathbb{E}(f(X(t)) - \min f) &= \mathbb{E}\left(f(X(t)) - f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right)\right) \\ &+ \mathbb{E}\left(f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right) - \min f\right). \end{aligned}$$

Let us recall that $\mathbb{E}(\sup_{s \geq 0} \|Z(s)\|) < +\infty$, which implies that $\mathbb{E}(\sup_{t \geq t_0} \|X(t)\|) < +\infty$. We bound the first term using the convexity of f and Cauchy-Schwarz inequality, that is, that for every $x, y \in \mathbb{H}$,

$$f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle \leq \|\nabla f(x)\| \|y - x\|,$$

we get that

$$\begin{aligned} f(X(t)) - f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right) &\leq \|\nabla f(X(t))\| \|\xi(t)\| \\ &\leq \|\xi(t)\| (L\|X(t)\| + \|\nabla f(0)\|) \\ &\leq \|\xi(t)\| \left(L \sup_{t \geq t_0} \|X(t)\| + \|\nabla f(0)\| \right), \end{aligned}$$

and we conclude that

$$\mathbb{E}\left(f(X(t)) - f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right)\right) = \mathcal{O}(e^{-A(t)}).$$

For the second term, we use Jensen's inequality to obtain

$$\begin{aligned} \mathbb{E}\left(f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right) - \min f\right) &\leq \int_{t_0}^t \mathbb{E}[f(Y(u)) - \min f] d\mu_t(u) \\ &\leq e^{-A(t)} \mathbb{E}[f(Y(t_0)) - \min f] \\ &\quad + e^{-A(t)} \int_{t_0}^t \frac{e^{A(u)}}{\Gamma(u)} \mathbb{E}[f(Y(u)) - \min f] du. \end{aligned}$$

Because $\sqrt{\theta}\Gamma\sigma_\infty \in L^2([t_0, +\infty[)$ is equivalent to $s \mapsto s\tilde{\sigma}_\infty^2(s) \in L^1(\mathbb{R}_+)$, by Theorem 2, we have that there exists $C > 0$ such that $\mathbb{E}(f(Z(s)) - \min f) \leq C/s$. Then, we have $\mathbb{E}(f(Y(t)) - \min f) \leq C/\theta(t)$. Hence, there exists $C_0 > 0$ such that

$$\mathbb{E}(f(X(t)) - \min f) \leq C_0 e^{-A(t)} + CI \left[\frac{1}{\theta} \right](t). \quad \square$$

Theorem 7. Let $v \geq 2$ and consider the dynamic (S-ISIHD) with initial data $X_0, V_0 \in L^v(\Omega; \mathbb{H})$, such that f and σ satisfy (H_0) and (H_σ) , and in the case where γ satisfies (H_γ) , $\beta \equiv \Gamma$. Moreover, suppose that $f \in C^2(\mathbb{H})$ and

$$t \mapsto \sqrt{\theta(t)}\Gamma(t)\sigma_\infty(t) \in L^2([t_0, +\infty[),$$

where $\theta(t) \stackrel{\text{def}}{=} \int_{t_0}^t \Gamma(u) du$. Then the solution trajectory $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^v[t_0]$ is unique and satisfies

$$\int_{t_0}^{+\infty} \theta(u)\Gamma(u)\|\nabla f(X(u) + \Gamma(u)V(u))\|^2 du < +\infty \quad \text{a.s.} \quad (13)$$

Proof. Consider (12) and the technique used in Theorem 5. We have that $t \mapsto \theta(t)\Gamma^2(t)\sigma_\infty^2(t) \in L^1([t_0, +\infty[)$ is equivalent to $s \mapsto s\tilde{\sigma}_\infty^2(s) \in L^1(\mathbb{R}_+)$. Therefore, we can use Theorem 2 to state that

$$\int_0^{+\infty} s\|\nabla f(Z(s))\|^2 ds < +\infty \quad \text{a.s.}$$

Using the time scaling $\tau \equiv \theta$ and making the change of variable $\theta(t) = s$ in the previous integral, we obtain

$$\int_{t_0}^{+\infty} \theta(t)\Gamma(t)\|\nabla f(Y(t))\|^2 dt < +\infty \quad \text{a.s.}$$

Recalling that in the averaging, we impose that $Y = X + \Gamma V$, we conclude. \square

3.3. Fast Convergence Under $\alpha > 3$, $\gamma(t) = \alpha/t$ and $\beta(t) = t/(\alpha - 1)$

In the following, we show fast convergence results in expectation and in almost sure sense. This result imposes an integrability condition on the diffusion term to ensure a convergence rate of $1/t^2$ for the function values. This is highly desirable, as it represents the fastest convergence rate we can expect to achieve for gradient-based

second-order dynamics featuring inertia when applied to a general convex function. Besides, these results match those in Attouch et al. (2024) when there is no noise.

Corollary 1 (Case α/t). Let $v \geq 2, \alpha > 3$ and consider the dynamic (S-ISIHD) with initial data $X_0, V_0 \in L^v(\Omega; \mathbb{H})$, in the case where $\gamma(t) = \alpha/t$ and $\beta(t) = t/(\alpha - 1)$. Besides, consider that f and σ satisfy (H_0) and (H_σ) . Moreover, let $f \in C^2(\mathbb{H})$ and $t \mapsto t^2 \sigma_\infty(t) \in L^2([t_0, +\infty[)$. Then the solution trajectory $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^v[t_0]$ is unique and satisfies

- i. $f(X(t)) - \min f = o(t^{-2})$ a.s.
- ii. $\mathbb{E}[f(X(t)) - \min f] = \mathcal{O}(t^{-2})$.

$$\int_{t_0}^{+\infty} t^3 \|\nabla f\left(X(t) + \frac{t}{\alpha-1} V(t)\right)\|^2 dt < +\infty \quad \text{a.s.}$$

Proof. Consider (12) with $\Gamma(t) = t/(\alpha - 1)$ and $\theta(t) = (t^2 - t_0^2)/(2(\alpha - 1))$. Let $\tilde{\sigma}(s, \cdot) = \sigma(\theta^{-1}(s), \cdot) \sqrt{\Gamma(\theta^{-1}(s))}$. Notice that $t \mapsto t^2 \sigma_\infty(t) \in L^2([t_0, +\infty[)$ is equivalent to $s \mapsto s \tilde{\sigma}_\infty^2(s) \in L^1(\mathbb{R}_+)$. We apply Theorem 2 to deduce that

$$f(Z(s)) - \min f = o(s^{-1}) \quad \text{a.s.}$$

Using the time scaling $\tau \equiv \theta$ and then the averaging technique as in the proof of Theorem 5, we have that

$$f(Y(t)) - \min f = o(t^{-2}) \quad \text{a.s.}$$

Moreover, it holds that

$$X(t) = \int_{t_0}^t Y(u) d\mu_t(u) + \xi(t).$$

- i. Now we prove the first point in the following way:

$$\begin{aligned} t^2(f(X(t)) - \min f) &= t^2\left(f(X(t)) - f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right)\right) \\ &\quad + t^2\left(f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right) - \min f\right). \end{aligned}$$

Let us bound the first term using the convexity of f :

$$\begin{aligned} t^2\left(f(X(t)) - f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right)\right) &\leq t^2 \|\nabla f(X(t))\| \|\xi(t)\| \\ &\leq t^2 \|\xi(t)\| (L \|X(t)\| + \|\nabla f(0)\|) \\ &\leq t^2 \|\xi(t)\| \left(L \sup_{t \geq t_0} \|X(t)\| + \|\nabla f(0)\|\right). \end{aligned}$$

Let us recall that $\sup_{s \geq 0} \|Z(s)\| < +\infty$ a.s. Because of the time scaling and averaging, it is direct to check that $\sup_{t \geq t_0} \|X(t)\| < +\infty$ a.s. On the other hand, $\|\xi(t)\| = \mathcal{O}(t^{1-\alpha})$ a.s. Therefore, we have

$$t^2\left(f(X(t)) - f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right)\right) = \mathcal{O}(t^{3-\alpha}) \quad \text{a.s.} \quad (14)$$

Now let us bound the second term using Jensen's inequality:

$$\begin{aligned} t^2\left(f\left(\int_{t_0}^t Y(u) d\mu_t(u)\right) - \min f\right) &\leq t^2\left(\int_{t_0}^t [f(Y(u)) - \min f] d\mu_t(u)\right) \\ &= \frac{t_0^{\alpha-1}}{t^{\alpha-3}} [f(Y(t_0)) - \min f] + \frac{\alpha-1}{t^{\alpha-3}} \int_{t_0}^t u^{\alpha-4} (u^2(f(Y(u)) - \min f)) du. \end{aligned}$$

In order to calculate the limit of this second term, let $a(t) = (\alpha - 1)/t$, $b(u) = u^2(f(Y(u)) - \min f)$; by Lemma 2, we

have that

$$\lim_{t \rightarrow +\infty} \frac{\alpha - 1}{t^{\alpha-1}} \int_{t_0}^t u^{\alpha-2} b(u) du = 0 \quad \text{a.s.}$$

Because $\alpha > 3$, we also have that

$$\lim_{t \rightarrow +\infty} \frac{\alpha - 3}{t^{\alpha-3}} \int_{t_0}^t u^{\alpha-4} b(u) du = 0 \quad \text{a.s.} \quad (15)$$

Therefore, we conclude that

$$\lim_{t \rightarrow +\infty} t^2 (f(X(t)) - \min f) = 0 \quad \text{a.s.}$$

ii. By Theorem 6, in the case $\gamma(t) = \alpha/t$, we have that $e^{-A(t)} = t_0^{\alpha-1} t^{1-\alpha}$ and $\theta(t) = (t^2 - t_0^2)/(2(\alpha - 1))$. On the other hand,

$$I \left[\frac{1}{\theta} \right] (t) = 2(\alpha - 1)^2 t^{1-\alpha} \int_{t_0}^t \frac{u^{\alpha-2}}{u^2 - t_0^2} = \mathcal{O}(t^{1-\alpha} + t^{-2}).$$

Because $\alpha > 3$, we have that $\mathcal{O}(t^{1-\alpha})$ is also $\mathcal{O}(t^{-2})$, and we conclude that

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

iii. This point follows directly from Theorem 7 in the case $\gamma(t) = \frac{\alpha}{t}$. \square

3.4. Convergence Rate Under Polyak-Łojasiewicz Inequality

In this section, we show a local convergence rate under Polyak-Łojasiewicz inequality. The Polyak-Łojasiewicz property is a special case of the Łojasiewicz property (Łojasiewicz 1963, 1965, 1984) and is commonly used to prove linear convergence of gradient descent algorithms.

Definition 2 (Polyak-Łojasiewicz Inequality). Let $f : \mathbb{H} \rightarrow \mathbb{R}$ be a differentiable function with $S \neq \emptyset$. Then, f satisfies the Polyak-Łojasiewicz (PL) inequality on S , if there exists $r > \min f$ and $\mu > 0$ such that

$$2\mu(f(x) - \min f) \leq \|\nabla f(x)\|^2, \quad \forall x \in [\min f < f < r], \quad (16)$$

and we will write $f \in \text{PL}_\mu(S)$.

Theorem 8. Let $\nu \geq 2$ and consider the dynamic (S-ISIHD) with initial data $X_0, V_0 \in L^\nu(\Omega; \mathbb{H})$, where f satisfies (H'_0) , and σ satisfies (H_σ) . Besides, $f \in \text{PL}_\mu(S)$ and suppose that either \mathbb{H} is finite dimensional or $f \in C^2(\mathbb{H})$. Let also, $\gamma \equiv \sqrt{2\mu}$, $\beta \equiv \Gamma \equiv 1/\sqrt{2\mu}$, and such that $\sigma_\infty \in L^2([t_0, +\infty])$.

Then the solution trajectory $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$ is unique. Moreover, letting $\delta > 0$, then there exists $\hat{t}_\delta > t_0, \hat{s}_\delta > 0, K_{\mu, \delta}, C_l, C_f > 0$ such that

$$\mathbb{E}(f(X(t)) - \min f) \leq K_{\mu, \delta} e^{-\frac{\mu}{2}(t-\hat{t}_\delta)} + \frac{1}{\mu} l_\delta \left(\frac{t + 3\hat{t}_\delta - 4t_0}{4\mu} \right) + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta, \quad (17)$$

where

$$l_\delta(s) = \frac{L}{2} \sigma_\infty^2(s) + C_l \sqrt{\delta} \frac{\sigma_\infty^2(s)}{2 \sqrt{\int_{\hat{s}_\delta}^s \sigma_\infty^2(u) du}}.$$

Besides, if $f \in \text{PL}_\mu(S)$ holds on the entire space (i.e., $r = +\infty$), then we have that there exists $K_\mu > 0$ such that

$$\mathbb{E}(f(X(t)) - \min f) \leq K_\mu e^{-\frac{\mu}{2}(t-t_0)} + \frac{L}{2\mu} \sigma_\infty^2 \left(\frac{t-t_0}{4\mu} \right), \quad \forall t > t_0. \quad (18)$$

Remark 5. If f is μ -strongly convex, then $f \in \text{PL}_\mu(S)$ holds on the entire space (i.e., $r = +\infty$).

Remark 6. For a proper discussion on the arbitrarily small (but nonzero) term $C_f\sqrt{\delta}$, we refer to (Maulen-Soto et al. 2024, section 4).

Proof. Consider the dynamic (S-ISIHD) with $\gamma \equiv c$, $\beta \equiv \Gamma \equiv 1/c$, where $c > 0$ is a constant that will be fixed later.

Let us also define $\theta(t) \stackrel{\text{def}}{=} \int_{t_0}^t \Gamma(u)du = (t - t_0)/c$ and $\tilde{\sigma}(s, \cdot) \stackrel{\text{def}}{=} \sigma(\theta^{-1}(s), \cdot)\sqrt{\Gamma(\theta^{-1}(s))}$. Then $\sigma_\infty \in L^2([t_0, +\infty[)$ is equivalent to $\tilde{\sigma}_\infty \in L^2(\mathbb{R}_+)$. Now consider the dynamic

$$\begin{cases} dZ(s) = -\nabla f(Z(s)) + \tilde{\sigma}(s, Z(s))dW(s), & s > 0, \\ Z(0) = X_0 + \Gamma(t_0)V_0. \end{cases} \quad (19)$$

Let $\delta > 0$ and apply the result of Maulen-Soto et al. (2024, theorem 4.5(i-b)) (with coefficient $\sqrt{2\mu}$); that is, there exists $\hat{s}_\delta > 0$ such that for every $\lambda \in]0, 1[$,

$$\begin{aligned} \mathbb{E}(f(Z(s)) - \min f) &\leq e^{-2\mu(s-\hat{s}_\delta)} \mathbb{E}(f(Z(\hat{s}_\delta)) - \min f) \\ &\quad + e^{-2\mu(1-\lambda)(s-\hat{s}_\delta)} \left(\frac{LC_\infty^2}{2} + C_l C_\infty \sqrt{\delta} \right) \\ &\quad + \frac{l_\delta(\hat{s}_\delta + \lambda(s - \hat{s}_\delta))}{2\mu} + C_f \sqrt{\delta}, \quad \forall s > \hat{s}_\delta, \end{aligned} \quad (20)$$

where $C_\infty, C_l, C_f > 0$ and the establishment of l_δ are detailed in Maulen-Soto et al. (2024, section 4.2), in particular

$$C_\infty = \sqrt{\int_{t_0}^{+\infty} \sigma_\infty^2(s)ds}.$$

Considering the time scaling $\tau \equiv \theta$, $Y(t) = Z(\theta(t))$ and $\hat{t}_\delta > t_0$ such that $\theta(\hat{t}_\delta) = \hat{s}_\delta$ (i.e., $\hat{t}_\delta = c\hat{s}_\delta + t_0$), we have that

$$\begin{aligned} \mathbb{E}(f(Y(t)) - \min f) &\leq e^{-2\mu(\theta(t)-\hat{s}_\delta)} \mathbb{E}(f(Y(\hat{t}_\delta)) - \min f) \\ &\quad + e^{-2\mu(1-\lambda)(\theta(t)-\hat{s}_\delta)} \left(\frac{LC_\infty^2}{2} + C_l C_\infty \sqrt{\delta} \right) \\ &\quad + \frac{l_\delta(\hat{s}_\delta + \lambda(\theta(t) - \hat{s}_\delta))}{2\mu} + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta. \end{aligned} \quad (21)$$

Let $a(t) = c$ and $A(t) = c(t - \hat{t}_\delta)$. Now, we consider the averaging as in (8) but change the initial condition to \hat{t}_δ . Thus, we have

$$X(t) = \int_{\hat{t}_\delta}^t Y(u) d\tilde{\mu}_t(u) + \tilde{\xi}(t), \quad (22)$$

where $\tilde{\mu}_t$ is the probability measure on $[\hat{t}_\delta, t]$ defined by

$$\tilde{\mu}_t = e^{-c(t-\hat{t}_\delta)} \delta_{\hat{t}_\delta} + ce^{c(u-t)} du, \quad (23)$$

where $\delta_{\hat{t}_\delta}$ is the Dirac measure at \hat{t}_δ and

$$\tilde{\xi}(t) \stackrel{\text{def}}{=} -\frac{1}{c} e^{-c(t-\hat{t}_\delta)} V(\hat{t}_\delta). \quad (24)$$

Then

$$\mathbb{E}(f(X(t)) - \min f) = \mathbb{E} \left(f(X(t)) - f \left(\int_{\hat{t}_\delta}^t Y(u) d\mu_t(u) \right) \right) + \mathbb{E} \left(f \left(\int_{\hat{t}_\delta}^t Y(u) d\mu_t(u) \right) - \min f \right).$$

We can bound the first term using convexity and Cauchy-Schwarz inequality in the following way:

$$\begin{aligned} &\mathbb{E} \left(f(X(t)) - f \left(\int_{\hat{t}_\delta}^t Y(u) d\tilde{\mu}_t(u) \right) \right) \\ &\leq \sqrt{\mathbb{E}(\|\nabla f(X(t))\|^2)} \sqrt{\mathbb{E}(\|\tilde{\xi}(t)\|^2)} \\ &\leq \frac{\sqrt{\mathbb{E}(\|V(\hat{t}_\delta)\|^2)}}{c} \sqrt{2\|\nabla f(0)\|^2 + 2L^2 \mathbb{E} \left(\sup_{t \geq t_0} \|X(t)\|^2 \right)} e^{-c(t-\hat{t}_\delta)}, \end{aligned}$$

where $\mathbb{E}(\sup_{t \geq t_0} \|X(t)\|^2) < +\infty$ as mentioned in Corollary 1.

On the other hand, we can bound the second term using Jensen’s inequality and then (21):

$$\begin{aligned}
 & \mathbb{E} \left(f \left(\int_{\hat{t}_\delta}^t Y(u) d\tilde{\mu}_t(u) \right) - \min f \right) \\
 & \leq \int_{\hat{t}_\delta}^t \mathbb{E}(f(Y(u)) - \min f) d\tilde{\mu}_t(u) \\
 & \leq \int_{\hat{t}_\delta}^t e^{-2\mu(\theta(u) - \hat{s}_\delta)} \mathbb{E}(f(Y(\hat{t}_\delta)) - \min f) d\tilde{\mu}_t(u) \\
 & + \int_{\hat{t}_\delta}^t e^{-2\mu(1-\lambda)(\theta(u) - \hat{s}_\delta)} \left(\frac{LC_\infty^2}{2} + C_l C_\infty \sqrt{\delta} \right) d\tilde{\mu}_t(u) \\
 & + \int_{\hat{t}_\delta}^t \frac{l_\delta(\hat{s}_\delta + \lambda(\theta(u) - \hat{s}_\delta))}{2\mu} d\tilde{\mu}_t(u) + C_f \sqrt{\delta} \\
 & = \left(\mathbb{E}(f(Y(\hat{t}_\delta)) - \min f) + \left(\frac{LC_\infty^2}{2} + C_l C_\infty \sqrt{\delta} \right) + \frac{l_\delta(\hat{s}_\delta)}{2\mu} \right) e^{-c(t - \hat{t}_\delta)} \\
 & + c \mathbb{E}(f(Y(\hat{t}_\delta)) - \min f) e^{-ct} \int_{\hat{t}_\delta}^t e^{-2\mu(\theta(u) - \hat{s}_\delta)} e^{cu} du \\
 & + c \left(\frac{LC_\infty^2}{2} + C_l C_\infty \sqrt{\delta} \right) e^{-ct} \int_{\hat{t}_\delta}^t e^{-2\mu(1-\lambda)(\theta(u) - \hat{s}_\delta)} e^{cu} du \\
 & + \frac{c}{2\mu} e^{-ct} \int_{\hat{t}_\delta}^t l_\delta(\hat{s}_\delta + \lambda(\theta(u) - \hat{s}_\delta)) e^{cu} du + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta.
 \end{aligned}$$

We bound the first integral as follows:

$$e^{-ct} \int_{\hat{t}_\delta}^t e^{-2\mu(\theta(u) - \hat{s}_\delta)} e^{cu} du \leq e^{2\mu \left(\frac{t_0}{c} + \hat{s}_\delta \right)} e^{-\frac{2\mu}{c}t}.$$

The second integral is bounded in the same way:

$$e^{-ct} \int_{\hat{t}_\delta}^t e^{-2\mu(1-\lambda)(\theta(u) - \hat{s}_\delta)} e^{cu} du \leq e^{2\mu(1-\lambda) \left(\frac{t_0}{c} + \hat{s}_\delta \right)} e^{-\frac{2\mu(1-\lambda)}{c}t}.$$

To treat the third integral, we are going to split the integral in two in order to find a useful convergence rate. Let

us recall that $l_\delta \in L^1([\hat{s}_\delta, +\infty[)$ and that l_δ is decreasing. Let us define $\varphi_{\lambda,c,\delta}(t) \stackrel{\text{def}}{=} \hat{s}_\delta + \lambda \left(\frac{t + \hat{t}_\delta - 2t_0}{2c} - \hat{s}_\delta \right)$, then

$$\begin{aligned}
 e^{-ct} \int_{\hat{t}_\delta}^t l_\delta(\hat{s}_\delta + \lambda(\theta(u) - \hat{s}_\delta)) e^{cu} du & = e^{-ct} \int_{\hat{t}_\delta}^{\frac{\hat{t}_\delta+t}{2}} l_\delta(\hat{s}_\delta + \lambda(\theta(u) - \hat{s}_\delta)) e^{cu} du \\
 & + e^{-ct} \int_{\frac{\hat{t}_\delta+t}{2}}^t l_\delta(\hat{s}_\delta + \lambda(\theta(u) - \hat{s}_\delta)) e^{cu} du \\
 & \leq \frac{c}{\lambda} e^{\frac{c\hat{t}_\delta}{2}} C_\infty e^{-\frac{ct}{2}} + l_\delta(\varphi_{\lambda,c,\delta}(t)).
 \end{aligned}$$

Now that we have bounded all the terms, we have the following bound:

$$\begin{aligned}
 \mathbb{E}(f(X(t)) - \min f) & \leq \frac{\sqrt{\mathbb{E}(\|V(\hat{t}_\delta)\|^2)}}{c} \sqrt{2\|\nabla f(0)\|^2 + 2L^2 \mathbb{E} \left(\sup_{t \geq t_0} \|X(t)\|^2 \right)} e^{-c(t - \hat{t}_\delta)} \\
 & + \left(\mathbb{E}(f(Y(\hat{t}_\delta)) - \min f) + \left(\frac{LC_\infty^2}{2} + C_l C_\infty \sqrt{\delta} \right) + \frac{l_\delta(\hat{s}_\delta)}{2\mu} \right) e^{-c(t - \hat{t}_\delta)} \\
 & + c \mathbb{E}(f(Y(\hat{t}_\delta)) - \min f) e^{2\mu \left(\frac{t_0 + \hat{t}_\delta}{c} + \hat{s}_\delta \right)} e^{-\frac{2\mu}{c}(t - \hat{t}_\delta)} \\
 & + c \left(\frac{LC_\infty^2}{2} + C_l C_\infty \sqrt{\delta} \right) e^{2\mu(1-\lambda) \left(\frac{t_0 + \hat{t}_\delta}{c} + \hat{s}_\delta \right)} e^{-\frac{2\mu(1-\lambda)}{c}(t - \hat{t}_\delta)} \\
 & + \frac{c}{2\mu} \left(\frac{c}{\lambda} C_\infty e^{-\frac{c(t - \hat{t}_\delta)}{2}} + l_\delta(\varphi_{\lambda,c,\delta}(t)) \right) + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta.
 \end{aligned}$$

Letting $\lambda = 1/2$ and $c = \sqrt{2\mu}$, we obtain

$$\begin{aligned} \mathbb{E}(f(X(t)) - \min f) &\leq \frac{\sqrt{\mathbb{E}(\|V(\hat{t}_\delta)\|^2)}}{\sqrt{2\mu}} \sqrt{2\|\nabla f(0)\|^2 + 2L^2\mathbb{E}\left(\sup_{t \geq t_0} \|X(t)\|^2\right)} e^{-\sqrt{2\mu}(t-\hat{t}_\delta)} \\ &\quad + \left(\mathbb{E}(f(Y(\hat{t}_\delta)) - \min f) + \left(\frac{LC_\infty^2}{2} + C_I C_\infty \sqrt{\delta}\right) + \frac{l_\delta(\hat{s}_\delta)}{2\mu}\right) e^{-\sqrt{2\mu}(t-\hat{t}_\delta)} \\ &\quad + \sqrt{2\mu}\mathbb{E}(f(Y(\hat{t}_\delta)) - \min f) e^{2\mu\left(\frac{t_0+\hat{t}_\delta}{\sqrt{2\mu}} + \hat{s}_\delta\right)} e^{-\sqrt{2\mu}(t-\hat{t}_\delta)} \\ &\quad + \sqrt{2\mu}\left(\frac{LC_\infty^2}{2} + C_I C_\infty \sqrt{\delta}\right) e^{\mu\left(\frac{t_0+\hat{t}_\delta}{\sqrt{2\mu}} + \hat{s}_\delta\right)} e^{-\frac{\sqrt{2\mu}}{2}(t-\hat{t}_\delta)} \\ &\quad + 2C_\infty e^{-\frac{\sqrt{2\mu}(t-\hat{t}_\delta)}{2}} + \frac{1}{\sqrt{2\mu}} l_\delta(\varphi_{\frac{1}{2}, \sqrt{2\mu}, \delta}(t)) + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta. \end{aligned}$$

Letting $K_{\mu, \delta} \stackrel{\text{def}}{=} \sqrt{2\mu}\left(\frac{LC_\infty^2}{2} + C_I C_\infty \sqrt{\delta}\right) e^{\mu\left(\frac{t_0+\hat{t}_\delta}{\sqrt{2\mu}} + \hat{s}_\delta\right)} + 2C_\infty$, we conclude that

$$\mathbb{E}(f(X(t)) - \min f) \leq K_{\mu, \delta} e^{-\frac{\sqrt{2\mu}}{2}(t-\hat{t}_\delta)} + \frac{1}{\sqrt{2\mu}} l_\delta(\varphi_{\frac{1}{2}, \sqrt{2\mu}, \delta}(t)) + C_f \sqrt{\delta}, \quad \forall t > \hat{t}_\delta. \quad \square \quad (25)$$

4. From Weak to Strong Convergence Under General γ and $\beta \equiv \Gamma$

4.1. General Result

We consider the Tikhonov regularization of the dynamic (S-ISIHD_{NS}), that is, for $t > 0$,

$$\begin{cases} dX(t) = V(t)dt, & t > t_0, \\ dV(t) \in -\gamma(t)V(t)dt - \partial F(X(t) + \Gamma(t)V(t))dt \\ \quad + \tilde{\sigma}_1(t, X(t) + \Gamma(t)V(t))dW(t), & t > t_0, \\ X(t_0) = X_0, \quad V(t_0) = V_0. \end{cases} \quad (\text{S-ISIHD-S.2})$$

We show some conditions (on $\gamma, \beta, \varepsilon$) under which we can obtain strong convergence of the trajectory.

Theorem 9. Consider that $\gamma : [t_0, +\infty[\rightarrow \mathbb{R}_+$ satisfies (H γ). Besides, $F = f + g$ and σ satisfy Assumptions (H $_0$) and (H $_\sigma$). Moreover, suppose that g satisfies (H $_\lambda$) and let $\nu \geq 2$. Consider (S-ISIHD_{NS-TA}) with $\beta \equiv \Gamma$ and 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 (S-ISIHD_{NS-TA}). Additionally, let $x^* \stackrel{\text{def}}{=} P_{S_F}(0)$ be the minimum norm solution, and for $\varepsilon > 0$, let x_ε be the unique minimizer of $F_\varepsilon(x) \stackrel{\text{def}}{=} F(x) + \varepsilon\|x\|^2/2$. If we suppose that $\Gamma\sigma_\infty \in L^2([t_0, +\infty[)$, and that $\varepsilon : [t_0, +\infty[\rightarrow \mathbb{R}_+$ satisfies the conditions:

$$(T'_1) \quad \varepsilon(t) \rightarrow 0 \text{ as } t \rightarrow +\infty;$$

$$(T'_2) \quad \int_{t_0}^{+\infty} \varepsilon(t)\Gamma(t)dt = +\infty; \text{ and}$$

$$(T'_3) \quad \int_{t_0}^{+\infty} \varepsilon(t)\Gamma(t)(\|x^*\|^2 - \|x_{\varepsilon(t)}\|^2)dt < +\infty.$$

Then $s\text{-}\lim_{t \rightarrow +\infty} X(t) = x^*$ a.s., and $V(t) = o\left(\frac{1}{\Gamma(t)}\right)$ a.s.

Proof. Let $s_0 > 0$, $\theta(t) \stackrel{\text{def}}{=} s_0 + \int_{t_0}^t \Gamma(u)du$; $\tilde{\varepsilon}(t) = \varepsilon(\theta^{-1}(t))$; and $\tilde{\sigma}(s, \cdot) \stackrel{\text{def}}{=} \sigma(\theta^{-1}(s), \cdot)\sqrt{\Gamma(\theta^{-1}(s))}$. Then ε satisfying (T $'_1$), (T $'_2$), and (T $'_3$) is equivalent to $\tilde{\varepsilon}$ satisfying (T $_1$), (T $_2$), and (T $_3$) defined in Theorem 4. Besides, $\Gamma\sigma_\infty \in L^2([t_0, +\infty[)$ is equivalent to $\tilde{\sigma}_\infty \in L^2(\mathbb{R}_+)$. Consider the dynamic:

$$\begin{cases} dZ(s) \in -\partial F(Z(s)) - \tilde{\varepsilon}(s)Z(s) + \tilde{\sigma}(s, Z(s))dW(s), & s > s_0, \\ Z(s_0) = X_0 + \Gamma(t_0)V_0. \end{cases} \quad (26)$$

By Theorem 4, we have that there exists a unique solution $Z \in S_{\mathbb{H}}^\nu[s_0]$, and that $\lim_{s \rightarrow +\infty} Z(s) = x^*$ a.s. (recall that $x^* \stackrel{\text{def}}{=} P_{S_F}(0)$). Using the time scaling $\tau \equiv \theta$ and the averaging described at the beginning of this section, we end up with the dynamic (S-ISIHD_{NS-TA}) in the case where $\beta \equiv \Gamma$. The existence and uniqueness of solution, and the fact that $(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^\nu[t_0]$ goes analogously as in the proof of Theorem 5.

Now we prove the claim: Because $\lim_{s \rightarrow +\infty} Z(s) = x^*$ a.s., this implies directly that $\lim_{t \rightarrow +\infty} Y(t) = x^*$ a.s. Besides, we have Relation (9), that is,

$$X(t) = \int_{t_0}^t Y(u) d\mu_t(u) + \xi(t),$$

where μ_t and ξ are defined in (10) and (11), respectively. Consequently, we have

$$\begin{aligned} \|X(t) - x^*\| &\leq \left\| X(t) - \int_{t_0}^t Y(u) d\mu_t(u) \right\| + \left\| \int_{t_0}^t Y(u) d\mu_t(u) - x^* \right\| \\ &\leq \|\xi(t)\| + \left\| \int_{t_0}^t Y(u) d\mu_t(u) - x^* \right\|. \end{aligned}$$

Let $a(t) = 1/\Gamma(t)$ and $A(t) = \int_{t_0}^t du/\Gamma(u)$. By Lemma 3, we have that $\lim_{t \rightarrow +\infty} \|\xi(t)\| = 0$. On the other hand,

$$\begin{aligned} \left\| \int_{t_0}^t Y(u) d\mu_t(u) - x^* \right\| &= \left\| \int_{t_0}^t (Y(u) - x^*) d\mu_t(u) \right\| \\ &\leq \int_{t_0}^t \|Y(u) - x^*\| d\mu_t(u) \\ &= e^{-A(t)} \|Y(t_0) - x^*\| + e^{-A(t)} \int_{t_0}^t a(u) e^{A(u)} \|Y(u) - x^*\| du. \end{aligned}$$

Let $b(u) = \|Y(u) - x^*\|$. Because we already proved that $\lim_{u \rightarrow +\infty} b(u) = 0$ a.s., and we have that $a \notin L^1([t_0, +\infty[)$ by Lemma 3, we utilize Lemma 2 with our respective a, b functions. This let us conclude that

$$\lim_{t \rightarrow +\infty} \left\| \int_{t_0}^t Y(u) d\mu_t(u) - x^* \right\| = 0 \quad \text{a.s.}$$

Thus, $\lim_{t \rightarrow +\infty} X(t) = x^*$ a.s. Finally, because

$$Y(t) = X(t) + \Gamma(t)V(t),$$

and the fact that X and Y have (a.s.) the same limit, we conclude that

$$\lim_{t \rightarrow +\infty} \Gamma(t)V(t) = 0 \quad \text{a.s.} \quad \square$$

4.2. Concrete Cases

In order to give some conditions when (T'_1) , (T'_2) , and (T'_3) of Theorem 9 are satisfied, we need to introduce the following definition.

Definition 3 (Hölderian Error Bound). Let $f : \mathbb{H} \rightarrow \mathbb{R}$ be a proper function such that $\mathcal{S} \neq \emptyset$; f satisfies a Hölderian (or power-type) error bound inequality on \mathcal{S} with exponent $p \geq 1$, if there exists $\kappa > 0$ and $r > \min f$ such that

$$f(x) - \min f \geq \kappa \text{dist}(x, \mathcal{S})^p, \quad \forall x \in [\min f \leq f \leq r], \quad (27)$$

and we will write $f \in \text{EB}^p(\mathcal{S})$.

Remark 7. Let $f : \mathbb{H} \rightarrow \mathbb{R}$ be a differentiable function such that $\mathcal{S} \neq \emptyset$. If f satisfies the P-L inequality on \mathcal{S} , then f satisfies a Hölderian error bound inequality with exponent $p = 2$.

Theorem 10. Consider the setting of Theorem 9 and suppose that $F = f + g \in \text{EB}^p(\mathcal{S}_F)$. Let $s_0 > 0$ and denote $\theta(t) \stackrel{\text{def}}{=} s_0 + \int_{t_0}^t \Gamma(s) ds$, then taking the Tikhonov parameter $\varepsilon(t) = 1/\theta^r(t)$ with

$$1 \geq r > \frac{2p}{2p+1},$$

the three conditions (T'_1) , (T'_2) , and (T'_3) of Theorem 9 are satisfied simultaneously. In particular, for any solution $(X, V) \in \mathcal{S}_{\mathbb{H} \times \mathbb{H}}^v[t_0]$ of (9), we get almost sure (strong) convergence of $X(t)$ to the minimal norm solution named $x^* = P_{\mathcal{S}_F}(0)$ and that $V(t) = o(1/\Gamma(t))$.

Proof. We proceed as in the proof of Theorem 9 and arrive to the dynamic (26); because $\tilde{\varepsilon}(t) = \varepsilon(\theta^{-1}(t)) = 1/t^r$, the proof goes as in Maulen-Soto et al. (2025, theorem 4.8). \square

Theorem 11. Let $v \geq 2, f \in \Gamma_0(\mathbb{H}) \cap C_L^2(\mathbb{H})$ such that S is nonempty, and also $f \in \text{EB}^p(S), \sigma$ satisfying (H_σ) , and $\Gamma\sigma_\infty \in L^2([t_0, +\infty[)$ and is nonincreasing. Let us consider $\varepsilon(t) = 1/t^r$, where $0 < r < 1$, then we evaluate $(S\text{-ISIH}\mathbb{D}_{\text{NS}}\text{-TA})$ in the case where γ satisfies $(H_\gamma), g \equiv 0, \beta \equiv \Gamma$, and with initial data $X_0, V_0 \in L^v(\Omega; \mathbb{H})$, that is, for $t > t_0$,

$$\begin{cases} dX(t) = V(t)dt, \\ dV(t) = -\gamma(t)V(t)dt - \nabla f(X(t) + \Gamma(t)V(t))dt - \varepsilon(t)(X(t) + \Gamma(t)V(t)) \\ \quad + \sigma(t, X(t) + \Gamma(t)V(t))dW(t), \\ X(t_0) = X_0, \quad V(t_0) = V_0. \end{cases} \tag{28}$$

For $\varepsilon > 0$, let us define $f_\varepsilon(x) \stackrel{\text{def}}{=} f(x) + \varepsilon\|x\|^2/2$, and let x_ε be the unique minimizer of f_ε . Moreover, let $s_0 > 0$ and for $s_1 > s_0$ consider

$$R(s) \stackrel{\text{def}}{=} e^{-\frac{s^{1-r}}{1-r}} \int_{s_1}^s e^{\frac{u^{1-r}}{1-r}} \sigma_\infty^2(\theta^{-1}(u)) \Gamma(\theta^{-1}(u)) du, \tag{29}$$

where $\theta(t) \stackrel{\text{def}}{=} s_0 + \int_{t_0}^t \Gamma(u)du$. Let also $x^* \stackrel{\text{def}}{=} P_S(0), A(s) \stackrel{\text{def}}{=} \int_{s_1}^s du/\Gamma(u)$, and $t_1 \stackrel{\text{def}}{=} \theta^{-1}(s_1)$. Then, the solution trajectory

$(X, V) \in S_{\mathbb{H} \times \mathbb{H}}^v[t_0]$ is unique, and we have that

- i. $R(\theta(t)) \rightarrow 0$ as $t \rightarrow +\infty$.
- ii. Let $\bar{\sigma}(t) = \Gamma(t)\sigma_\infty^2(t)$, then

$$R(\theta(t)) = \mathcal{O}\left(\exp(-\theta^r(t)(1 - 2^{-r})) + \theta^r(t)\bar{\sigma}\left(\frac{s_1 + \theta(t)}{2}\right)\right).$$

Moreover, if $\bar{\sigma}(t) = \mathcal{O}(\theta^{-\Delta}(t))$ for $\Delta > 1$, then $R(\theta(t)) = \mathcal{O}(\theta^{r-\Delta}(t))$.

Besides, we have the following convergence rate in expectation:

- iii. For the values, we have

$$\mathbb{E}[f(X(t)) - \min(f)] = \mathcal{O}(\max\{e^{-A(t)}, I[h_1](t)\}),$$

where $h_1(t) = 1/\theta^r(t) + R(\theta(t))$.

- iv. Also, for the trajectory, we obtain

$$\mathbb{E}[\|X(t) - x^*\|^2] = \mathcal{O}(\max\{e^{-A(t)}, I[h_2](t)\}),$$

where $h_2(t) = \theta^{r-1}(t) + \theta^{-\frac{r}{p}}(t) + \theta^r(t)R(\theta(t))$.

Proof. We proceed as in the proof of Theorem 5 and define analogously $\tilde{\sigma}, \tilde{\varepsilon}$; we also consider the dynamic (12). By Maulen-Soto et al. (2025, theorem 4.11) we obtain that

$$R(s) = e^{-\frac{s^{1-r}}{1-r}} \int_{s_1}^s e^{\frac{u^{1-r}}{1-r}} \tilde{\sigma}_\infty^2(u) du,$$

where $\tilde{\sigma}_\infty^2 \in L^2([s_0, +\infty[)$, satisfies the following:

- $R(s) \rightarrow +\infty$ as $t \rightarrow +\infty$.
- $R(s) = \mathcal{O}(\exp(-s^r(1 - 2^{-r})) + s^r \tilde{\sigma}_\infty^2(\frac{s_1+s}{2}))$. Moreover if $\tilde{\sigma}_\infty^2(s) = \mathcal{O}(s^{-\Delta})$ for $\Delta > 1$, then $R(s) = \mathcal{O}(s^{r-\Delta})$.

Also, evaluating at $s = \theta(t)$, we obtain the first two items of the theorem. For the third and fourth items, we used that

- $\mathbb{E}[f(Z(s)) - \min(f)] = \mathcal{O}(\frac{1}{s^r} + R(s))$.
- $\mathbb{E}[\|Z(s) - x^*\|^2] = \mathcal{O}(\frac{1}{s^{1-r}} + \frac{1}{s^p} + s^r R(s))$.

Then, proceeding as in the proof of Theorem 6, we obtain the desired results. \square

Corollary 2. Consider Theorem 10 in the case where $\gamma(t) = \alpha/t$ for $\alpha > 1, \beta(t) = t/(\alpha - 1)$; then, we have that

- 1. If $\sigma_\infty^2(t) = \mathcal{O}(t^{-2(\Delta+1)})$ for $\Delta > 1$, and $\alpha \neq \{1 + 2r, 1 + 2(\Delta - r)\}$, then

$$\mathbb{E}[f(X(t)) - \min(f)] = \mathcal{O}(\max\{t^{-(\alpha-1)}, t^{-2r}, t^{-2(\Delta-r)}\}).$$

In particular, if $\alpha > 3$,

- 2. If $\sigma_\infty^2(t) = \mathcal{O}(t^{-2(\Delta+1)})$ for $\Delta > \max\{1, 2r\}$, and $\alpha \neq \{3 - 2r, 1 + 2r/p, 1 + 2(2r - \Delta)\}$, then

$$\mathbb{E}[\|X(t) - x^*\|^2] = \mathcal{O}(\max\{t^{-(\alpha-1)}, t^{-2(1-r)}, t^{-\frac{2r}{p}}, t^{-2(2r-\Delta)}\}).$$

5. Conclusion

This work uncovers the global and local convergence properties of trajectories of the inertial system with implicit Hessian-driven damping under stochastic errors both in the smooth and nonsmooth setting. The aim is to solve convex optimization problems with noisy gradient input with vanishing variance. We have shed light on these properties and provided a comprehensive local and global complexity analysis both in the case where the Hessian damping parameter β was dependent on the geometric damping γ and when it was zero. We believe that this work, along with the technique of time scaling and averaging, paves the way for important extensions and research avenues. Among them, we mention extension to the situation where the drift term is a nonpotential cocoercive operator.

Appendix A. SDEs Are Better Approximations to Stochastic Inertial Algorithms

In this section, we will argue that (S-ISIHD) is an accurate continuous-time model of a natural stochastic inertial algorithm with an accuracy error that scales as $\mathcal{O}(h)$, and this is much better than (ISIHD) whose accuracy is only $\mathcal{O}(\sqrt{h})$. This generalizes (Dambrine et al. 2024, proposition 2.1) and holds for a more general model of the diffusion coefficient and under weaker regularity assumptions. The proof is reminiscent of strong consistency bounds of Euler-type methods for nonlinear SDEs (Kloeden and Platen 1992) that we will refine by exploiting the structure of our SDE (S-ISIHD). We also restrict our discussion to the finite-dimensional case where $\mathbb{H} = \mathbb{K} = \mathbb{R}^d$. Similar approximation results can be found, for instance, in Hu et al. (2019b), Li et al. (2021), and Fontaine et al. (2021), whereas for momentum-based algorithms, the tools developed in Li et al. (2017) may provide useful insights for deriving alternative approximation results, once adapted to our Hessian-driven damping setting.

Given a step-size $h > 0$, the Euler-Maruyama discretization applied to (S-ISIHD) computes approximations X_k and V_k of $X(t_k)$ and $V(t_k)$, where $t_k = t_0 + kh$, $k \in \mathbb{N}$, by initializing $(X_0, V_0) = (X(t_0), V(t_0))$ and forming the following iterative scheme:

$$\begin{cases} X_{k+1} = X_k + hV_k, \\ V_{k+1} = (1 - \gamma_k h)V_k - h\nabla f(X_k + \beta_k V_k) + \sqrt{h}\sigma_k G_k, \end{cases} \quad (\text{A.1})$$

where $G_k \sim \mathcal{N}(0, I_d)$, $\gamma_k = \gamma(t_k)$, $\sigma_k = \sigma(t_k, X_k + \beta_k V_k)$ and $\beta_k = \beta(t_k)$. This is a stochastic version of the algorithm proposed in Alecsa et al. (2021). Setting $Y_k = (X_k, V_k)$, it will be convenient to rewrite (A.1) in the product space \mathbb{R}^{2d} as

$$\begin{cases} Y_{k+1} = Y_k + h\Phi(t_k, Y_k) + \sqrt{h}\zeta(t_k, Y_k)(W(t_{k+1}) - W(t_k)), \\ Y_0 = (X_0, V_0), \end{cases} \quad (\text{A.2})$$

where W is a \mathbb{R}^{2d} -valued Brownian motion and

$$\Phi(t, (x, v)) = (v, -\gamma(t)v - \nabla f(x + \beta(t)v)), \quad \text{and} \quad \zeta(t, (x, v)) = \begin{pmatrix} 0_{d \times d} & 0_{d \times d} \\ 0_{d \times d} & \sigma(t, x + \beta(t)v) \end{pmatrix}.$$

This then motivates the SDE

$$\begin{cases} dY(t) = \Phi(t, Y(t))dt + \sqrt{h}\zeta(t, Y(t))dW(t), \\ Y(t_0) = (X_0, V_0), \end{cases} \quad (\text{A.3})$$

where we set $Y(t) = (X(t), V(t))$ which is (S-ISIHD) with an extra \sqrt{h} in front the diffusion.

As classical in numerical analysis of evolution equations, we define the continuous-time piece-wise linear extension of the sequence $(Y_k)_{k \in \mathbb{N}}$:

$$\begin{aligned} \bar{Y}(t) &\stackrel{\text{def}}{=} Y_k + (t - t_k)\Phi(t_k, Y_k) + \sqrt{h}\zeta(t_k, Y_k)(W(t) - W(t_k)), \quad t \in [t_k, t_{k+1}[\\ &= Y_0 + \int_{t_0}^t \hat{\Phi}(s, \hat{Y}(s))ds + \sqrt{h} \int_{t_0}^t \hat{\zeta}(s, \hat{Y}(s))dW(s), \end{aligned}$$

where for $t \in [t_k, t_{k+1}[$, $\hat{Y}(t) \stackrel{\text{def}}{=} Y_k$, $\hat{\Phi}(t, \hat{Y}(t)) \stackrel{\text{def}}{=} \Phi(t_k, Y_k)$, $\hat{\zeta}(t, \hat{Y}(t)) \stackrel{\text{def}}{=} \zeta(t_k, Y_k)$. We also define $\hat{\gamma}(t)$ and $\hat{\beta}(t)$ similarly. We thus have $\hat{\sigma}(t, \hat{X}(t) + \hat{\beta}(t)\hat{V}(t)) = \sigma_k$ for $t \in [t_k, t_{k+1}[$.

Our result requires the following assumption.

Assumption A.1. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuously differentiable function with L -Lipschitz continuous gradient. For $T > t_0$, γ and β are respectively L_γ - and L_β -Lipschitz continuous on $[t_0, T]$. σ verifies, $\forall t, s \in [t_0, T]$ and $\forall x, x' \in \mathbb{R}^d$,

$$\begin{aligned} \|\sigma(t, x) - \sigma(t, x')\| &\leq L_\sigma \|x - x'\|, \\ \|\sigma(t, x)\| &\leq K_\sigma (1 + \|x\|), \\ \|\sigma(t, x) - \sigma(s, x)\| &\leq L_\sigma (1 + \|x\|) |t - s|^{1/2}. \end{aligned}$$

where all constants $L_\gamma, L_\beta, L_\sigma, K_\sigma$ do not depend on h .

All these assumptions are verified in the instances studied in Section 3 when f is smooth. Observe that these assumptions also ensure existence and uniqueness of the solution $(X(t), V(t))$ to (A.3) (hence, (S-ISIHD)).

We have strong consistency results that show that (A.3) is a better continuous-time approximation to (A.1) than (ISIHD).

Proposition A.1. *Suppose that Assumption 1 holds and that $X_0, V_0 \in L^2(\Omega; \mathbb{R}^d)$. Let $(X(t), V(t))$ be the solution to (A.3) and $(x(t), v(t) = \dot{x}(t))$ the solution to (ISIHD). Then the iterates of algorithm (A.1) satisfy, as $h \rightarrow 0$,*

$$\mathbb{E} \left[\sup_{0 \leq k \leq N-1} \|X_k - X(t_k)\| + \|V_k - V(t_k)\| \right] = \mathcal{O}(h), \quad (\text{A.4})$$

$$\mathbb{E} \left[\sup_{0 \leq k \leq N-1} \|X_k - x(t_k)\| + \|V_k - v(t_k)\| \right] = \mathcal{O}(\sqrt{h}). \quad (\text{A.5})$$

Proof. By Assumption 1, it is not difficult to see that Φ and ζ verify the following the Lipschitz continuity and linear growth properties:

$$\begin{aligned} \|\Phi(t, y) - \Phi(t, z)\|^2 &\leq \max(4L^2, 1 + 2\gamma_{\max^2} + 4L^2\beta_{\max^2})\|y - z\|^2, & \forall t \in [t_0, T], y, z \in \mathbb{R}^d, \\ \|\Phi(t, y) - \Phi(s, y)\|^2 &\leq (2L_\gamma^2 + 4L^2L_\beta^2)\|y\|^2|t - s|^2, & \forall s, t \in [t_0, T], y \in \mathbb{R}^d, \\ \|\Phi(t, y)\|^2 &\leq \max(8L^2, 1 + 2\gamma_{\max^2} + 8L^2\beta_{\max^2})\|y\|^2 + 4\|\nabla f(0)\|^2, & \forall t \in [t_0, T], y \in \mathbb{R}^d. \end{aligned} \quad (\text{A.6})$$

In the rest of the proof, C is any positive constant that does not depend on h (and may depend on $d, T - t_0$, the different Lipschitz and growth constants in Assumption 1), and that may change from one line to another. Using Jensen's inequality, Doob's martingale inequality, and Mao (2007, theorem 1.5.21), we have for any $\tau \in [t_0, T]$,

$$\begin{aligned} &\mathbb{E} \left[\sup_{t_0 \leq t \leq \tau} \|\bar{Y}(t) - Y(t)\|^2 \right] \\ &= \mathbb{E} \left[\sup_{t_0 \leq t \leq \tau} \left\| \int_{t_0}^t (\hat{\Phi}(s, \hat{Y}(s)) - \Phi(s, Y(s))) ds + \sqrt{h} \int_{t_0}^t (\hat{\zeta}(s, \hat{Y}(s)) - \zeta(s, Y(s))) dW(s) \right\|^2 \right] \\ &\leq C \mathbb{E} \left[\int_{t_0}^{\tau} \|\hat{\Phi}(s, \hat{Y}(s)) - \Phi(s, Y(s))\|^2 ds \right] \\ &\quad + Ch \mathbb{E} \left[\int_{t_0}^{\tau} \|\hat{\sigma}(s, \hat{X}(s) + \hat{\beta}(s)\hat{V}(s)) - \sigma(s, X(s) + \beta(s)V(s))\|^2 ds \right]. \end{aligned} \quad (\text{A.7})$$

Using Jensen's inequality, (A.6), and Assumption 1 on σ , we get

$$\begin{aligned} &\mathbb{E} \left[\sup_{t_0 \leq t \leq \tau} \|\bar{Y}(t) - Y(t)\|^2 \right] \\ &\leq C \mathbb{E} \left[\int_{t_0}^{\tau} \|\hat{\Phi}(s, \hat{Y}(s)) - \Phi(s, \hat{Y}(s))\|^2 ds \right] \\ &\quad + Ch \mathbb{E} \left[\int_{t_0}^{\tau} \|\hat{\sigma}(s, \hat{X}(s) + \hat{\beta}(s)\hat{V}(s)) - \sigma(s, \hat{X}(s) + \hat{\beta}(s)\hat{V}(s))\|^2 ds \right] \\ &\quad + C(1+h) \int_{t_0}^{\tau} \mathbb{E} \left[\sup_{t_0 \leq r \leq s} \|\bar{Y}(r) - Y(r)\|^2 \right] ds + C(1+h) \mathbb{E} \left[\int_{t_0}^{\tau} \|\bar{Y}(t) - \hat{Y}(t)\|^2 ds \right] \\ &\quad + Ch \mathbb{E} \left[\sup_{t_0 \leq t \leq T} \|Y(t)\|^2 \right] \int_{t_0}^{\tau} |\beta(s) - \hat{\beta}(s)|^2 ds, \end{aligned}$$

where the expectation in the last display is bounded by Kloeden and Platen (1992, theorem 4.5.4). Let $N = \lfloor (T - t_0)/h \rfloor$. Lipschitz continuity of β gives

$$\begin{aligned} \int_{t_0}^{\tau} |\beta(s) - \hat{\beta}(s)|^2 ds &\leq \sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} |\beta(s) - \beta(t_k)|^2 ds \\ &\leq L_\beta^2 \sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} (s - t_k)^2 ds \\ &= \frac{L_\beta^2 (T - t_0)}{3} h^2. \end{aligned}$$

Now, we have for any $s \in [t_k, t_{k+1}[$,

$$\bar{Y}(s) - \hat{Y}(s) = -(s - t_k)\Phi(t_k, Y_k) - \sqrt{h}\zeta(t_k, Y_k)(W(s) - W(t_k)).$$

It then follows from (A.6), Assumption 1 on σ , and Kloeden and Platen (1992, theorem 4.5.4) that

$$\begin{aligned} & \mathbb{E} \left[\int_{t_0}^{\tau} \|\bar{Y}(t) - \hat{Y}(t)\|^2 ds \right] \\ & \leq \sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} \mathbb{E} [\|(s - t_k)\Phi(t_k, Y_k) + \sqrt{h}\zeta(t_k, Y_k)(W(s) - W(t_k))\|^2] ds \\ & \leq Ch^2. \end{aligned}$$

With similar arguments, we have

$$\mathbb{E} \left[\int_{t_0}^{\tau} \|\hat{\Phi}(s, \hat{Y}(s)) - \Phi(s, Y_k)\|^2 ds \right] = \sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} \mathbb{E} [\|\Phi(t_k, Y_k) - \Phi(s, Y_k)\|^2] ds \leq Ch^2,$$

and using Assumption 1 on σ

$$\begin{aligned} & \mathbb{E} \left[\int_{t_0}^{\tau} \|\hat{\sigma}(s, \hat{X}(s) + \hat{\beta}(s)\hat{V}(s)) - \sigma(s, X_k + \beta_k V_k)\|^2 ds \right] \\ & = \sum_{k=0}^{N-1} \int_{t_k}^{t_{k+1}} \mathbb{E} [\|\sigma(t_k, X_k + \beta_k V_k) - \sigma(s, X_k + \beta_k V_k)\|^2] ds \leq Ch. \end{aligned}$$

Collecting all bounds, we get

$$\mathbb{E} \left[\sup_{t_0 \leq t \leq \tau} \|\bar{Y}(t) - Y(t)\|^2 \right] \leq C(h^2 + h^3) + C(1 + h) \int_{t_0}^{\tau} \mathbb{E} \left[\sup_{t_0 \leq r \leq s} \|\bar{Y}(r) - Y(r)\|^2 \right] ds.$$

Applying the Jensen and Gronwall inequalities then gives

$$\mathbb{E} \left[\sup_{t_0 \leq t \leq T} \|\bar{Y}(t) - Y(t)\| \right] \leq \mathbb{E} \left[\sup_{t_0 \leq t \leq T} \|\bar{Y}(t) - Y(t)\|^2 \right]^{1/2} \leq C(h + h^{3/2})e^{C(1+h)T} = \mathcal{O}(h).$$

In turn

$$\mathbb{E} \left[\sup_{0 \leq k \leq N-1} \|Y_k - Y(t_k)\| \right] = \mathbb{E} \left[\sup_{0 \leq k \leq N-1} \|\bar{Y}(t_k) - Y(t_k)\| \right] \leq \mathbb{E} \left[\sup_{t_0 \leq t \leq T} \|\bar{Y}(t) - Y(t)\| \right] = \mathcal{O}(h),$$

which gives (A.4). For (A.5), it is sufficient to see that with System (A.3), the term $\zeta(s, y(s))$ disappears in the main inequality (A.7). Starting from there, this becomes the leading term that scales as $\mathcal{O}(h)$ (instead of $\mathcal{O}(h^2)$ previously), hence entailing (A.5). \square

Appendix B. Auxiliary Results

B.1. Deterministic Results

Lemma B.1. 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 B.2. Let $a, b : [t_0, +\infty[\rightarrow \mathbb{R}_+$ be two functions such that $a \notin L^1([t_0, +\infty[)$, $\lim_{u \rightarrow +\infty} b(u) = 0$, and define $A(t) \stackrel{\text{def}}{=} \int_{t_0}^t a(u)du$ and $B(t) \stackrel{\text{def}}{=} e^{-A(t)} \int_{t_0}^t a(u)e^{A(u)}b(u)du$. Then $\lim_{t \rightarrow +\infty} B(t) = 0$.

Proof. Let $\varepsilon > 0$ arbitrary, let us take T_ε such that $t_0 < T_\varepsilon$ and $b(u) \leq \varepsilon$ for $u \geq T_\varepsilon$. For $t > T_\varepsilon$, let us write

$$\begin{aligned} B(t) &= e^{-A(t)} \int_{t_0}^{T_\varepsilon} a(u)e^{A(u)}b(u)du + e^{-A(t)} \int_{T_\varepsilon}^t a(u)e^{A(u)}b(u)du \\ &\leq e^{-A(t)} \int_{t_0}^{T_\varepsilon} a(u)e^{A(u)}b(u)du + \varepsilon. \end{aligned}$$

Because $a \notin L^1([t_0, +\infty[)$, then $\lim_{t \rightarrow +\infty} e^{-A(t)} = 0$, we get

$$\limsup_{t \rightarrow +\infty} B(t) \leq \varepsilon.$$

This being true for any $\varepsilon > 0$, we infer that $\lim_{t \rightarrow +\infty} B(t) = 0$, which gives the claim. \square

Lemma B.3. Under Hypothesis (H_γ) , then

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

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

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

B.2. On Stochastic Processes

Let us recall some elements of stochastic analysis. Throughout the paper, $(\Omega, \mathcal{F}, \mathbb{P})$ is a probability space and $\{\mathcal{F}_t | t \geq 0\}$ is a filtration of the σ -algebra \mathcal{F} . Given $\mathcal{C} \in 2^\Omega$, we will denote $\sigma(\mathcal{C})$ the σ -algebra generated by \mathcal{C} . We denote $\mathcal{F}_\infty \stackrel{\text{def}}{=} \sigma(\cup_{t \geq 0} \mathcal{F}_t) \in \mathcal{F}$.

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

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

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

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

An \mathbb{H} -valued stochastic process starting at $t_0 \geq 0$ is a function $X : \Omega \times [t_0, +\infty[\rightarrow \mathbb{H}$. It is said to be continuous if $X(\omega, \cdot) \in C([t_0, +\infty[; \mathbb{H})$ for almost all $\omega \in \Omega$. We will denote $X(t) \stackrel{\text{def}}{=} X(\cdot, t)$. We are going to study SDEs and SDIs, and in order to ensure the uniqueness of a solution, we introduce a relation over stochastic processes. Two stochastic processes $X, Y : \Omega \times [t_0, T] \rightarrow \mathbb{H}$ are said to be equivalent if $X(t) = Y(t)$, $\forall t \in [t_0, T]$, \mathbb{P} -a.s. This leads us to define the equivalence relation \mathcal{R} , which associates the equivalent stochastic processes in the same class.

Furthermore, we will need some properties about the measurability of these processes. A stochastic process $X : \Omega \times [t_0, +\infty[\rightarrow \mathbb{H}$ is progressively measurable if for every $t \geq t_0$, the map $\Omega \times [t_0, t] \rightarrow \mathbb{H}$ defined by $(\omega, s) \rightarrow X(\omega, s)$ is $\mathcal{F}_t \otimes \mathcal{B}([t_0, t])$ -measurable, where \otimes is the product σ -algebra and \mathcal{B} is the Borel σ -algebra. On the other hand, X is \mathcal{F}_t -adapted if $X(t)$ is \mathcal{F}_t -measurable for every $t \geq t_0$. It is a direct consequence of the definition that if X is progressively measurable, then X is \mathcal{F}_t -adapted.

Let us define the quotient space:

$$S_{\mathbb{H}}^0[t_0, T] \stackrel{\text{def}}{=} \{X : \Omega \times [t_0, T] \rightarrow \mathbb{H}, X \text{ is a prog. measurable cont. stochastic process}\} / \mathcal{R}.$$

Set $S_{\mathbb{H}}^0[t_0] \stackrel{\text{def}}{=} \cap_{T \geq t_0} S_{\mathbb{H}}^0[t_0, T]$. For $\nu > 0$, we define $S_{\mathbb{H}}^\nu[t_0, T]$ as the subset of processes $X(t)$ in $S_{\mathbb{H}}^0[t_0, T]$ such that

$$S_{\mathbb{H}}^\nu[t_0, T] \stackrel{\text{def}}{=} \left\{ X \in S_{\mathbb{H}}^0[t_0, T] : \mathbb{E} \left(\sup_{t \in [t_0, T]} \|X_t\|^\nu \right) < +\infty \right\}.$$

We define $S_{\mathbb{H}}^\nu[t_0] \stackrel{\text{def}}{=} \cap_{T \geq t_0} S_{\mathbb{H}}^\nu[t_0, T]$.

Following the notation of Gawarecki and (Mandrekar 2011, section 2.1.2), we say that W_t is a \mathbb{K} -valued cylindrical Brownian motion defined on the filtered space $(\Omega, \mathcal{F}, \mathcal{F}_t, \mathbb{P})$ if

1. For an arbitrary $t \geq 0$, the mapping $W_t : \mathbb{K} \rightarrow L^2(\Omega; \mathbb{R})$ is linear;
2. For an arbitrary $k \in \mathbb{K}$, $W_t(k)$ is an \mathcal{F}_t Brownian motion; and
3. For arbitrary $k, k' \in \mathbb{K}$ and $t \geq 0$, $\mathbb{E}[W_t(k)W_t(k')] = t\langle k, k' \rangle_{\mathbb{K}}$.

Remark B.1. There is no \mathbb{K} -valued process \tilde{W}_t such that

$$W_t(k)(\omega) = \langle \tilde{W}_t(\omega), k \rangle_{\mathbb{K}}.$$

However, if Q is a nonnegative definite symmetric trace-class operator on \mathbb{K} , then a \mathbb{K} -valued Q -Brownian motion can be defined (Da Prato and Zabszyk 2014, section 4.1; Gawarecki and Mandrekar 2011, definition 2.6). Moreover, if $\mathbb{K} = \mathbb{R}^m$ then $W_t(k) = \langle \tilde{W}_t, k \rangle_{\mathbb{K}}$, where \tilde{W}_t corresponds to the standard m -dimensional Brownian motion.

Besides, let $G : \Omega \times \mathbb{R}_+ \rightarrow \mathcal{L}_2(\mathbb{K}; \mathbb{H})$ be a measurable and \mathcal{F}_t -adapted process, then we can define $\int_0^t G(s) dW(s)$, which is the stochastic integral of G , and we have that $G \rightarrow \int_0^\cdot G(s) dW(s)$ is an isometry between the measurable and \mathcal{F}_t -adapted

$\mathcal{L}_2(\mathbb{K}; \mathbb{H})$ -valued processes and the space of \mathbb{H} -valued continuous square-integrable martingales (Gawarecki and Mandrekar 2011, theorem 2.4).

References

- Alecsa C, László S, Pinta T (2021) An extension of the second order dynamical system that models Nesterov's convex gradient method. *Appl. Math. Optim.* 84(2):1687–1716.
- Allen-Zhu Z (2017) Katyusha: The first direct acceleration of stochastic gradient methods. *J. Machine Learn. Res.* 18(221):1–51.
- Apidopoulos V, Aujol JF, Dossal C (2018) The differential inclusion modeling the FISTA algorithm and optimality of convergence rate in the case $b \leq 3$. *SIAM J. Optim.* 28(1):551–574.
- Assran M, Rabbat M (2020) On the convergence of Nesterov's accelerated gradient method in stochastic settings. Shawe-Taylor J, Zemel RS, Bartlett P, Pereira FCN, Weinberger KQ, eds. *Proc. 37th Internat. Conf. Machine Learn.*, vol. 119 (Curran Associates, Inc., Red Hook, NY), 410–420.
- Attouch H, Cabot A (2017) Asymptotic stabilization of inertial gradient dynamics with time-dependent viscosity. *J. Differential Equations* 263(9):5412–5458.
- Attouch H, Peypouquet J (2016) The rate of convergence of Nesterov's accelerated forward-backward method is actually faster than $\frac{1}{\sqrt{\kappa}}$. *SIAM J. Optim.* 26(3):1824–1834.
- Attouch H, Bot R, Nguyen DK (2024) Fast convex optimization via time scale and averaging of the steepest descent. *Math. Oper. Res.* 50(4):2633–2665.
- Attouch H, Chbani Z, Riahi H (2019) Rate of convergence of the Nesterov accelerated gradient method in the subcritical case $\alpha \leq 3$. *ESAIM: Control Optim. Calculus Variations (ESAIM-COCV)*, 25(2).
- Attouch H, Fadili J, Kungurtsev V (2023a) On the effect of perturbations in first-order optimization methods with inertia and Hessian driven damping. *Evolution Equations Control Theory* 12(1):71–117.
- Attouch H, Fadili J, Kungurtsev V (2024) The stochastic Ravine accelerated gradient method with general extrapolation coefficients. Preprint, Submitted March 7, <https://arxiv.org/abs/2403.04860>.
- Attouch H, Goudou X, Redont P (2011) The heavy ball with friction method. I: The continuous dynamical system. *Comm. Contemporary Math.* 2(1):1–34.
- Attouch H, Peypouquet J, Redont P (2016) Fast convex optimization via inertial dynamics with Hessian driven damping. *J. Differential Equations* 261(10):5734–5783.
- Attouch H, Balhag A, Chbani Z, Riahi H (2022a) Fast convex optimization via inertial dynamics combining viscous and Hessian-driven damping with time rescaling. *Evolution Equations Control Theory* 11(2):487–514.
- Attouch H, Cabot A, Chbani Z, Riahi H (2018a) Accelerated forward-backward algorithms with perturbations: Application to Tikhonov regularization. *J. Optim. Theory Appl.* 179(1):1–36.
- Attouch H, Chbani Z, Fadili J, Riahi H (2023b) Convergence of iterates for first-order optimization algorithms with inertia and Hessian driven damping. *Optimization* 72(5):1199–1238.
- Attouch H, Chbani Z, Fadili J, Riahi H (2022b) First-order optimization algorithms via inertial systems with Hessian driven damping. *Math. Programming* 193(1):113–155.
- Attouch H, Chbani Z, Peypouquet J, Redont P (2018b) Fast convergence of inertial dynamics and algorithms with asymptotic vanishing viscosity. *Math. Programming Ser. B* 168(1–2):123–175.
- Barakat A, Bianchi P, Hachem W, Schechtman S (2021) Stochastic optimization with momentum: Convergence, fluctuations, and traps avoidance. *J. Statist.* 15(2):3892–3947.
- Bolte J, Nguyen T, Peypouquet J, Suter BW (2016) From error bounds to the complexity of first-order descent methods for convex functions. *Math. Programming* 165(2):471–507.
- Brézis H (1973) *Opérateurs Maximaux Monotones et Semi-Groupes de Contractions dans les Espaces de Hilbert*, Mathematics Studies, vol. 5 (North-Holland, New York).
- Cabot A (2009) Asymptotics for a gradient system with memory term. *Proc. Amer. Math. Soc.* 137(9):3013–3024.
- Cabot A, Engler H, Gadat S (2009) On the long time behavior of second order differential equations with asymptotically small dissipation. *Trans. Amer. Math. Soc.* 361(11):5983–6017.
- Castera C, Attouch H, Fadili J, Ochs P (2024) Continuous Newton-like methods featuring inertia and variable mass. *SIAM J. Optim.* 34(1):251–277.
- Castera C, Bolte J, Févotte C, Pauwels E (2021) An inertial Newton algorithm for deep learning. *J. Machine Learn. Res.* 22(134):1–31.
- Cauchy A (1847) Méthode générale pour la résolution des systèmes d'équations simultanées. *Comptes Rendus Hebdomadaires des Séances de l'Académie des Sciences* 25:536–538.
- Cheng X, Chatterji N, Bartlett PL, Jordan MI (2018) Underdamped Langevin MCMC: A non-asymptotic analysis. *Proc. Machine Learn. Res.* 75:1–24.
- Da Prato G, Zabszyk J (2014) *Stochastic Equations in Infinite Dimensions*, 2nd ed. (Cambridge University Press, Cambridge, UK).
- Dalalyan AS, Karagulyan A (2019) User-friendly guarantees for the Langevin Monte Carlo with inaccurate gradient. *Stochastic Processes Appl.* 129(12):5278–5311.
- Dalalyan A, Riou-Durand L, Karagulyan A (2019) Bounding the error of discretized Langevin algorithms for non-strongly log-concave targets. *J. Machine Learn. Res.* 23(235):1–38.
- Dambrine M, Dossal C, Puij B, Rondepierre A (2024) Stochastic differential equations for modeling first order optimization methods. *SIAM J. Optim.* 34(2):1402–1426.
- Defazio A, Jelassi S (2022) Adaptivity without compromise: A momentumized, adaptive, dual averaged gradient method for stochastic optimization. *J. Machine Learn. Res.* 23(144):1–34.
- Dossal C, Aujol J (2015) Stability of over-relaxations for the forward-backward algorithm, application to FISTA. *SIAM J. Optim.* 25(4):2408–2433.

- Driggs D, Ehrhardt MJ, Schönlieb C (2022) Accelerating variance-reduced stochastic gradient methods. *Math. Programming* 191(2):671–715.
- Fontaine X, De Bortoli V, Durmus A (2021) Convergence rates and approximation results for SGD and its continuous-time counterpart. Belkin M, Kpotufe S, eds. *Proc. Thirty Fourth Conf. Learn. Theory*, Proceedings of Machine Learning Research, vol. 134 (PMLR, New York), 1965–2058.
- Frostig R, Ge R, Kakade S, Sidford A (2015) Un-regularizing: Approximate proximal point and faster stochastic algorithms for empirical risk minimization. Bach F, Blei D, eds. *Proc. 32nd Internat. Conf. Machine Learn.*, Proceedings of Machine Learning Research, vol. 37 (PMLR, New York), 2540–2548.
- Gadat S, Panloup F (2014) Long time behaviour and stationary regime of memory gradient diffusions. *Annales de L'Institut Henri Poincaré - Probabilités et Statistiques* 50(2):564–601.
- Gadat S, Panloup F, Saadane S (2018) Stochastic heavy ball. *Electronic J. Statist.* 12(1):461–529.
- Gawarecki L, Mandrekar V (2011) *Stochastic Differential Equations in Infinite Dimensions: With Applications to Stochastic Partial Differential Equations* (Springer, New York).
- Goudou X, Munier J (2005) Asymptotic behavior of solutions of a gradient-like integrodifferential Volterra inclusion. *Adv. Math. Sci. Appl.* 15(2):509–525.
- Hamadouche A, Wu Y, Wallace AM, Mota JC (2024) Sharper bounds for proximal gradient algorithms with errors. *SIAM J. Optim.* 34(1):278–305.
- Haraux A, Jendoubi M (2012) On a second order dissipative ODE in Hilbert space with an integrable source term. *Acta Math. Sci.* 32(1):155–163.
- Hu W, Li C, Su W (2019c) On the global convergence of continuous-Time stochastic heavy-ball method for nonconvex optimization. 2019 *IEEE Internat. Conf. Big Data* (IEEE, Piscataway, NJ), 94–104.
- Hu W, Li C, Zhou X (2019a) On the global convergence of continuous-time stochastic heavy-ball method for nonconvex optimization. Baru CK, Huan J, Khan L, Hu X, Ak R, Tian Y, Barga RS, Zaniolo C, Lee K, Ye Y, eds. *Proc. 2019 IEEE Internat. Conf. Big Data* (Institute of Electrical and Electronics Engineers, Piscataway, NJ), 94–104.
- Hu W, Li C, Li L, Lui JG (2019b) On the diffusion approximation of nonconvex stochastic gradient descent. *Ann. Math. Sci. Appl.* 4(1):3–32.
- Jain P, Netrapalli P, Kakade SM, Kidambi R, Sidford A (2017) Parallelizing stochastic gradient descent for least squares regression: Minibatching, averaging, and model misspecification. *J. Machine Learn. Res.* 18(223):1–42.
- Kloeden PE, Platen E (1992) *Numerical Solution of Stochastic Differential Equations* (Springer-Verlag, Berlin).
- Laborde M, Oberman A (2020) A Lyapunov analysis for accelerated gradient methods: From deterministic to stochastic case. Chiappa S, Calandra R, eds. *Proc. Twenty Third Internat. Conf. Artificial Intelligence Statist.*, Proceedings of Machine Learning Research, vol. 108 (PMLR, New York), 602–612.
- Lan G (2020) *First-order and Stochastic Optimization Methods for Machine Learning* (Springer Nature, Cham, Switzerland).
- Latz J (2021) Analysis of stochastic gradient descent in continuous time. *Statist. Comput.* 31:39.
- Le T (2024) Nonsmooth nonconvex stochastic heavy ball. *J. Optim. Theory Appl.* 201(2):699–719.
- Li Z, Malladi S, Arora S (2021) On the validity of modeling SGD with stochastic differential equations. *Advances in Neural Information Processing Systems*, vol. 21 (Curran Associates Inc., Red Hook, NY).
- Li Q, Tai C, Weinan E (2017) Stochastic modified equations and adaptive stochastic gradient algorithms. *Proc. 34th Internat. Conf. Machine Learn.*, vol. 70 (PMLR, New York), 2101–2110.
- Li X, Shen Z, Zhang L, He N (2024) A Hessian-aware stochastic differential equation for modelling SGD. Preprint, submitted May 28, <https://arxiv.org/abs/2405.18373>.
- Lin H, Mairal J, Harchaoui Z (2017) Catalyst acceleration for first-order convex optimization: From theory to practice. *J. Machine Learn. Res.* 18(212):1–54.
- Loizou N, Richtárik P (2020) Momentum and stochastic momentum for stochastic gradient, Newton, proximal point and subspace descent methods. *Comput. Optim. Appl.* 77(3):653–710.
- Łojasiewicz S (1963) Une propriété topologique des sous-ensembles analytiques réels. *Colloque du C.N.R.S. sur les Équations aux Dérivées Partielles (Paris)*, 87–89.
- Łojasiewicz S (1965) *Ensembles Semi-analytiques* (Prépublication) (Institut des Hautes Études Scientifiques, Bures-sur-Yvette, France).
- Łojasiewicz S (1984) Sur les trajectoires du gradient d'une fonction analytique. *Seminari di Geometria* 1982/1983 (Universita di Bologna, Dipartimento di Matematica, Bologna, Italy), 115–117.
- Ma YA, Chatterji N, Cheng X, Flammarion N, Bartlett P, Jordan MI (2021) Is there an analog of Nesterov acceleration for MCMC? *Bernoulli* 27(3):1942–1992.
- Mandt S, Hoffman M, Blei D (2016) A variational analysis of stochastic gradient algorithms. *Proc. 33rd Internat. Conf. Machine Learn.* (PMLR, New York).
- Mao X (2007) *Stochastic Differential Equations and Applications*, 2nd ed. (Woodhead Publishing, Chichester, UK).
- Maulen-Soto R, Fadili J, Attouch H (2024) An SDE perspective on stochastic convex optimization. *Math. Oper. Res.* 50(4):3190–3221.
- Maulen-Soto R, Fadili J, Attouch H (2025) Stochastic differential inclusions and Tikhonov regularization for stochastic non-smooth convex optimization in Hilbert spaces. *Open J. Math. Optim.* 6, article no.9.
- May R (2017) Asymptotic for a second order evolution equation with convex potential and vanishing damping term. *Turkish J. Math.* 41(3):681–685.
- Mertikopoulos P, Staudigl M (2018) On the convergence of gradient-like flows with noisy gradient input. *SIAM J. Optim.* 28(1):163–197.
- Muehlebach M, Jordan MI (2021) Optimization with momentum: Dynamical, control-theoretic, and symplectic perspectives. *J. Machine Learn. Res.* 22(73):1–50.
- Nesterov Y (1983) A method of solving a convex programming problem with convergence rate $\mathcal{O}(1/k^2)$. *Doklady Akademii Nauk SSSR (Proc. USSR Acad. Sci.)* 269(3):543–547.
- Öksendal B (2003) *Stochastic Differential Equations*, 6th ed. (Springer-Verlag, Berlin).
- Orvieto A, Lucchi A (2019) Continuous-time models for stochastic optimization algorithms. Larochelle H, Ranzato M, Hadsell R, Balcan MF, Lin H, eds. *Advances in Neural Information Processing Systems*, vol. 32 (Curran Associates, Inc., Red Hook, NY).

- Orvieto A, Kohler J, Lucchi A (2020) The role of memory in stochastic optimization. Adams RP, Gogate V, eds. *Proc. 35th Conf. Uncertainty Artificial Intelligence*, Proceedings of Machine Learning Research, vol. 115 (PMLR, New York), 356–366.
- Pavliotis GA (2014) *Stochastic Processes and Applications: Diffusion Processes, the Fokker-Planck Equation and Large Deviations* (Springer, New York).
- Petterson R (1995) Yosida approximations for multivalued stochastic differential equations. *Stochastics Stochastics Rep.* 52(1–2):107–120.
- Polyak B (1964) Some methods of speeding up the convergence of iteration methods. *USSR Comput. Math. Math. Physics* 4(5):1–17.
- Robins H, Monro S (1951) A stochastic approximation method. *Ann. Math. Statist.* 22(3):400–407.
- Rockafellar R (1997) *Convex Analysis* (Princeton University Press, Princeton, NJ).
- Schmidt M, Le Roux N, Bach F (2011) Convergence rates of inexact proximal-gradient methods for convex optimization. Shawe-Taylor J, Zemel RS, Bartlett PL, Pereira FCN, Weinberger KQ, eds. *Advances in Neural Information Processing Systems*, vol. 24 (Curran Associates, Inc., Red Hook, NY).
- Shi B, Su WJ, Jordan MI (2023) On learning rates and Schrödinger operators. *J. Machine Learn. Res.* 24(379):1–53.
- Shi B, Du S, Jordan M, Su WJ (2022) Understanding the acceleration phenomenon via high resolution differential equations. *Math. Programming* 195:79–148.
- Soatto S, Chaudhari P (2018) Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. Choi SP, Yilmaz O, Poor HV, eds. *Proc. 2018 Inform. Theory Appl. Workshop (ITA)* (IEEE, Piscataway, NJ), 1–10.
- Su W, Boyd S, Candès E (2016) A differential equation for modeling Nesterov’s accelerated gradient method: Theory and insights. *J. Machine Learn. Res.* 17(153):1–43.
- Villa S, Salzo S, Baldassarres L (2013) Accelerated and inexact forward-backward. *SIAM J. Optim.* 23(3):1607–1633.
- Xie Z, Sato I, Sugiyama M (2021) A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima. *Proc. Ninth Internat. Conf. Learn. Representations (ICLR)* (OpenReview.net).
- Yan B (2018) Theoretical analysis for convex and non-convex clustering algorithms Doctoral dissertation, The University of Texas at Austin, Austin.