



HAL
open science

ERROR ESTIMATES BETWEEN SGD WITH MOMENTUM AND UNDERDAMPED LANGEVIN DIFFUSION

Arnaud Guillin, Yu Wang, Lihu Xu, Haoran Yang

► **To cite this version:**

Arnaud Guillin, Yu Wang, Lihu Xu, Haoran Yang. ERROR ESTIMATES BETWEEN SGD WITH MOMENTUM AND UNDERDAMPED LANGEVIN DIFFUSION. 2024. hal-04746813

HAL Id: hal-04746813

<https://hal.science/hal-04746813v1>

Preprint submitted on 21 Oct 2024

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

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

ERROR ESTIMATES BETWEEN SGD WITH MOMENTUM AND UNDERDAMPED LANGEVIN DIFFUSION

ARNAUD GUILLIN, YU WANG, LIHU XU, AND HAORAN YANG

ABSTRACT. Stochastic gradient descent with momentum is a popular variant of stochastic gradient descent, which has recently been reported to have a close relationship with the underdamped Langevin diffusion. In this paper, we establish a quantitative error estimate between them in the 1-Wasserstein and total variation distances.

Keywords: Stochastic Gradient Descent with momentum (SGDm), Underdamped Langevin Diffusion, 1-Wasserstein Distance, Total Variation Distance, Variant Time Step, Malliavin Calculus

MSC2020 subject classification:60J20, 60H35, 60H30, 60F99.

CONTENTS

1. Introduction	2
1.1. Literature Review	3
1.2. Contributions and Methods	4
1.3. Structure of The Paper	5
2. Preliminary and Main Results	5
2.1. Notations and Assumptions	5
2.2. Main Results	7
3. Auxiliary Lemmas	9
3.1. Moments Estimates	10
3.2. Auxiliary Lemmas for 1-Wasserstein Distance	10
3.3. Auxiliary Lemmas for Total Variation Distance	11
4. Proofs of the Main Results	14
4.1. Proof of Theorem 1	14
4.2. Proof of Theorem 2	15
4.3. Proofs of Corollaries 3 and 4	18
4.4. Proof of Corollary 5	19
References	21
Appendix	23
Appendix A. Supporting Lemmas	23
Appendix B. Proofs of Auxiliary Lemmas	25
B.1. Proofs of Lemmas for Moments Estimations	25
B.2. Proofs of Auxiliary Lemmas for 1-Wasserstein Distance	29
B.3. Proofs of Auxiliary Lemmas in the Total Variation Distance	32

1. INTRODUCTION

Many tasks in machine learning and statistics can be formulated as an optimization problem as follows

$$(1.1) \quad \text{minimize} \quad \mathbb{E}_\xi F(\mathbf{x}, \xi),$$

where ξ is random and its distribution is not known, and $\mathbf{x} \in \mathbb{R}^d$. In practice, one needs to replace the mean $\mathbb{E}_\xi F(\mathbf{x}, \xi)$ with its sample mean as

$$(1.2) \quad \text{minimize} \quad \frac{1}{N} \sum_{i=1}^N F(\mathbf{x}, \xi^i),$$

where ξ^1, \dots, ξ^N are i.i.d. and have the same distribution as ξ . For the further use, we denote

$$(1.3) \quad f(\mathbf{x}) = \mathbb{E}_\xi F(\mathbf{x}, \xi).$$

To solve the minimization problem (1.2), one often uses online stochastic gradient descent (SGD) or online SGD with momentum (SGDm), and the latter usually converges faster or works more efficient than the former. The relationship between SGD and stochastic differential equation (SDE) has been intensively studied recently, see e.g. [L^TW17, L^TW19, CSX23, CLTZ20, FDBD21]. Although there have been papers reporting the connections between stochastic gradient descent with momentum (SGDm) and underdamped Langevin diffusion, see for instance [GGZ22], their relationship has not been well understood, particularly when the time tends to infinity. The primary goal of this paper is to quantify their error bound in the 1-Wasserstein and total variation distances.

The noised online SGDm has two variables $(\mathbf{m}_k, \mathbf{x}_k)$, called the moment and position respectively, satisfying

$$(1.4) \quad \begin{cases} \mathbf{m}_{k+1} = \mathbf{m}_k - \gamma \eta_{k+1} \mathbf{m}_k - \frac{\eta_{k+1}}{N} \sum_{i=1}^N \nabla F(\mathbf{x}_k, \xi_{k+1}^i) + \beta \sqrt{\eta_{k+1}} \boldsymbol{\zeta}_{k+1}, & k \geq 0, \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \eta_{k+1} \mathbf{m}_k, & k \geq 0, \end{cases}$$

with initial value $(\mathbf{m}_0, \mathbf{x}_0) \in \mathbb{R}^{2d}$, where the constant $\gamma > 0$ is the friction-coefficient, $\beta > 0$ can be regarded as the temperature, $\{\eta_k, k \geq 1\}$ denote the step size, $\{\xi_k^i, k \geq 1, 1 \leq i \leq N\}$ are i.i.d. copies of ξ , and $\{\boldsymbol{\zeta}_k, k \geq 1\}$ are i.i.d. with standard d -dimensional normal distribution. Note that for each $k \geq 1$, $\nabla^k F(\mathbf{x}, \xi)$ is the k -th order derivative of F with respect to \mathbf{x} throughout this paper.

We will compare the algorithm (1.4) with the underdamped Langevin diffusion $(\mathbf{M}_t, \mathbf{X}_t)_{t \geq 0}$ with state space \mathbb{R}^{2d}

$$(1.5) \quad \begin{cases} d\mathbf{M}_t = -\gamma \mathbf{M}_t dt - \nabla f(\mathbf{X}_t) dt + \beta d\mathbf{B}_t, \\ d\mathbf{X}_t = \mathbf{M}_t dt, \end{cases}$$

where $(\mathbf{B}_t)_{t \geq 0}$ is a d -dimensional standard Brownian motion. The two processes possess the same initial point $(\mathbf{M}_0, \mathbf{X}_0) = (\mathbf{m}_0, \mathbf{x}_0)$. In order to compare the algorithm (1.4) with SDE (1.5), we introduce the following intermediate stochastic system: for

$t \in [t_k, t_{k+1})$ with $k \geq 0$,

$$(1.6) \quad \begin{cases} d\widetilde{\mathbf{M}}_t = -\gamma\widetilde{\mathbf{M}}_{t_k} dt - \frac{1}{N} \sum_{i=1}^N \nabla F(\widetilde{\mathbf{X}}_{t_k}, \xi_{k+1}^i) dt + \beta d\mathbf{B}_t, \\ d\widetilde{\mathbf{X}}_t = \widetilde{\mathbf{M}}_{t_k} dt, \end{cases}$$

where $t_k = \sum_{j=1}^k \eta_j$ and $t_0 = 0$. It is obvious that $\{(\mathbf{m}_k, \mathbf{x}_k), k \geq 0\}$ and $\{(\widetilde{\mathbf{M}}_{t_k}, \widetilde{\mathbf{X}}_{t_k}), k \geq 0\}$ have the same distribution as long as $(\widetilde{\mathbf{M}}_{t_0}, \widetilde{\mathbf{X}}_{t_0}) = (\mathbf{m}_0, \mathbf{x}_0)$.

1.1. Literature Review. Comparisons between stochastic algorithms and stochastic continuous dynamics have been intensively studied recently, e.g. [GM91, RRT17, CSX23]. In particular, the online stochastic gradient Langevin descent (SGLD) related to the minimization problem (1.2) reads as

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla F(\mathbf{x}_k, \xi) + \sqrt{2\beta^{-1}\eta} \boldsymbol{\zeta}_{k+1}, \quad k \geq 0,$$

with a constant learning rate $\eta > 0$. [RRT17] studied the connection between this SGLD and the overdamped Langevin diffusion process $\mathbf{X} = (\mathbf{X}_t)_{t \geq 0}$

$$d\mathbf{X}_t = -\nabla f(\mathbf{X}_t) dt + \sqrt{2\beta^{-1}} d\mathbf{B}_t,$$

see also [Ebe16, Ebe11, CHS87, Pav14] for related distribution sampling problems. Besides, [RRT17] showed that the 2-Wasserstein distance between the SGLD with constant step size η and the continuous time Langevin diffusion can be bounded by $\mathcal{O}(k\eta(\delta^{1/4} + \eta^{1/4}))$ at the k -th iteration. It is obvious that as k has an order higher than $[\eta(\delta^{1/4} + \eta^{1/4})]^{-1}$, this bound is not useful. By the Lindeberg technique, [CSX23] obtained a uniform bound $\mathcal{O}((1 + \log |\eta|)\eta)$ in 1-Wasserstein distance with respect to the time. Furthermore, [PP23] studied the convergence of the Euler-Maruyama (EM) scheme for the SDEs driven by multiplicative noises in the 1-Wasserstein and total variation distances, the steps of the EM scheme are decreasing.

The underdamped Langevin diffusion (1.5) has been extensively studied in recent years, see for instance [EGZ19, Pav14, BGM10, GM16, Wu01, Sch24, Nea11, CCBJ18]. An important advantage of the underdamped diffusion is that it often converges to the stationary faster than the overdamped diffusion, see for instance [EGZ19, CLW23, AAMN24] either in Wasserstein distance or in L^2 when the spectral gap is small by appropriately choosing friction coefficient. [Wu01] showed that the system converges to its equilibrium measure with exponential rate by constructing appropriate Lyapunov test functions. [EGZ19, Sch24] established the contraction property in the Wasserstein type distances for two solutions of (1.5) with different initial data, and [SW24] studied the EM scheme of (1.5) and its convergence in the 1-Wasserstein distance. [CCAY+18] considered the problem of sampling from distribution $p^*(\mathbf{x}) \propto e^{-f(\mathbf{x})}$ where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth everywhere and m -strongly convex outside a ball of finite radius R . Given a tolerance error ε , they proved that the complexity of the underdamped Langevin Markov chain Monte Carlo (MCMC) is $\mathcal{O}(\sqrt{d}/\varepsilon)$ if the step size is $\eta = \mathcal{O}(\varepsilon/\sqrt{d})$. What's more, [CR22, GGZ22] considered the stochastic gradient Hamilton Monte Carlo (SGHMC) for sampling a similar distribution and found a complexity bound $\mathcal{O}((k\eta)^{3/2} \sqrt{\log k\eta} (\delta^{1/4} + \eta^{1/4}) + k\eta\sqrt{\eta})$ in the 2-Wasserstein distance when the step

size is a constant η . For other applications related to underdamped Langevin diffusion, we refer the reader to [CDC15, CR22, GGZ22] for stochastic gradient Hamilton MC and to [CCBJ18, CCAY⁺18] for underdamped Langevin MCMC.

For other SGD algorithms such as Mini-batch SGD and Nesterov SGD, we refer the reader to [KLRT15, LZCS14, CL20, LR20, BLB17, MJ19, KF16] and the references therein. These algorithms differ in their approaches to updating the parameters during training neural networks and can be useful in improving convergence and reducing overfitting.

1.2. Contributions and Methods. Our contributions and methods are summarised as the below:

(1) Although there have been several papers studying the connections between the accelerated algorithms with continuous dynamics, see for instance [CR22, GGZ22], but most of them only provide qualitative limit without a quantitative error bound, particularly when the time is large. To the best of our knowledge, this paper first provides error bounds between the SGDM and the underdamped Langevin diffusion uniformly with respect to the time, these bounds depend on the dimension d polynomially and reveal convergence rates. Our results can be formally formulated as the following:

$$\begin{aligned} d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{M}_{t_n}, \mathbf{X}_{t_n}), \mathcal{L}(\mathbf{m}_n, \mathbf{x}_n)) &\leq C\sqrt{d}(1 + 1/N)\sqrt{\eta_n}, \\ d_{\text{TV}}(\mathcal{L}(\mathbf{M}_{t_n}, \mathbf{X}_{t_n}), \mathcal{L}(\mathbf{m}_n, \mathbf{x}_n)) &\leq Cd^4(\sqrt{\eta_n} + 1/\sqrt{N}). \end{aligned}$$

We clearly see a significant difference between the above two bounds, one being of an order $O(\sqrt{\eta_n} + \sqrt{\eta_n}/N)$ in the 1-Wasserstein distance and the other $O(\sqrt{\eta_n} + 1/\sqrt{N})$ in the total variation distance. We think that the rate $O(1/\sqrt{N})$ in the total variation distance is due to the singularity of this distance and that it is hard to improve the $O(1/\sqrt{N})$ without paying a price on the rate $\sqrt{\eta_n}$.

To the best of our knowledge, most of the known results related to underdamped Langevin diffusion samplings and discretization schemes have a constant step size η , while our SGDM has the non-increasing step sizes η_k which include the constant step size case.

(2) There are two systems of noises, ξ_k^i and ζ_{k+1} , in SGDM (1.4), whose interplays make the SGDM much more complex than the underdamped Langevin sampling by discretising SDE (1.5). These complexities can be clearly seen in bounding the total variation distance, where we need to estimate a Malliavin matrix related to ζ_{k+1} by splitting the regimes of ξ_k^i . What's more, ξ_k^i can be heavy tailed. More precisely, we only need to assume that ξ_k^i has second and fourth moments to bound $d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{M}_{t_n}, \mathbf{X}_{t_n}), \mathcal{L}(\mathbf{m}_n, \mathbf{x}_n))$ and $d_{\text{TV}}(\mathcal{L}(\mathbf{M}_{t_n}, \mathbf{X}_{t_n}), \mathcal{L}(\mathbf{m}_n, \mathbf{x}_n))$ respectively.

(3) The term $N^{-1} \sum_{i=1}^N \nabla F(x, \xi_{k+1}^i)$ is an unbiased estimator of $\nabla f(x)$. According to the law of large number, it converges to $\nabla f(x)$ almost surely as $N \rightarrow \infty$. This, together with the exponential ergodicity of $(\mathbf{M}_t, \mathbf{X}_t)$, immediately implies an error bound in the order $O(\sqrt{\eta_n})$ for the underdamped Langevin sampling as n is large. The rate $O(\sqrt{\eta_n})$ may be not optimal due to the heavy tail effect of ξ_k^i , we conjecture that the rate $O(\sqrt{\eta_n})$ can be improved by assuming that ξ_k^i has a high order moment. We leave this possible improvement to the future research.

(4) Although the approach of our proofs is still via the classical Lindeberg principle, in contrast to the overdamped Langevin diffusion, the underdamped Langevin diffusion is degenerate and the following difficulties naturally arise: (i) the regularity

problems are very involved and we need Malliavin calculus to handle them; (ii) there are interplays between the randomnesses of ξ_k^i and ζ_{k+1} , it is much more difficult for us to estimate the Malliavin matrix and figure out the polynomial dependence on the dimension d ; (iii) the related Lyapunov function is much more subtle and not intuitive.

1.3. Structure of The Paper. In the next section, we will introduce our assumptions and main theorems. Section 3 will provide all auxiliary lemmas, which include estimates of p -th moments, the contraction property of the 1-Wasserstein distance, and additional estimates related to the total variation distance. The proofs for the main results and their corollaries will be provided in Section 4. In Appendix A, we give the supporting lemmas related to the Lyapunov function. At last, the lemmas in Section 3 will be proved in Appendix B.

Acknowledgements: We would like to thank Professor Feng-Yu Wang for very helpful discussion and pointing out the references [Sch24] and [SW24] to us. A. Guillin is benefited from a government grant managed by the Agence Nationale de la Recherche under the France 2030 investment plan ANR-23-EXMA-0001”n and under the grant ANR-23-CE40-0003. L. Xu is supported by the National Natural Science Foundation of China No. 12071499, the Science and Technology Development Fund (FDCT) of Macau S.A.R. FDCT 0074/2023/RIA2, and the University of Macau grants MYRG2020-00039-FST, MYRG-GRG2023-00088-FST.

2. PRELIMINARY AND MAIN RESULTS

2.1. Notations and Assumptions. We use normal font for scalars (e.g. a, A, \dots) and boldface for vectors and matrices (e.g. $\mathbf{x}, \mathbf{y}, \mathbf{A}, \mathbf{B}, \dots$). For any vectors $\mathbf{u} = (u_1, \dots, u_d)$, $\mathbf{v} = (v_1, \dots, v_d)$ in \mathbb{R}^d , their standard inner product is denoted by $\langle \mathbf{u}, \mathbf{v} \rangle = \sum_{i=1}^d u_i v_i$, and denote the corresponding Euclidean norm as $|\mathbf{u}| = (\sum_{i=1}^d u_i^2)^{1/2}$. The HilbertSchmidt inner product for matrices $\mathbf{A} = (A_{ij})_{d \times d}$, $\mathbf{B} = (B_{ij})_{d \times d} \in \mathbb{R}^{d \times d}$ is denoted by $\langle \mathbf{A}, \mathbf{B} \rangle_{\text{HS}} = \sum_{i,j=1}^d A_{ij} B_{ij}$, the HilbertSchmidt norm is defined as $\|\mathbf{A}\|_{\text{HS}} = (\sum_{i,j=1}^d A_{ij}^2)^{1/2}$. Besides, the operator norm of \mathbf{A} is denoted by $\|\mathbf{A}\|_{\text{op}} = \sup_{|\mathbf{v}|=1} |\mathbf{A}\mathbf{v}|$, which has the following relationship with $\|\mathbf{A}\|_{\text{HS}}$

$$\|\mathbf{A}\|_{\text{op}} \leq \|\mathbf{A}\|_{\text{HS}} \leq \sqrt{d} \|\mathbf{A}\|_{\text{op}}.$$

If a matrix \mathbf{A} is positive semi-definite, we define $\lambda_{\max}(\mathbf{A})$ as its maximal eigenvalue and $\lambda_{\min}(\mathbf{A})$ as its minimal eigenvalue.

Let $\mathcal{C}(\mathbb{R}^d, \mathbb{R})$ denote the set of all continuous functions from \mathbb{R}^d to \mathbb{R} and $\mathcal{C}_b(\mathbb{R}^d, \mathbb{R})$ denote the set of all bounded continuous functions from \mathbb{R}^d to \mathbb{R} . For $k \geq 0$, denote by $\mathcal{C}^k(\mathbb{R}^d, \mathbb{R})$ the set of functions from \mathbb{R}^d to \mathbb{R} which has continuous 0-th, ..., k -th order derivatives, further denote by $\mathcal{C}_b^k(\mathbb{R}^d, \mathbb{R})$ the set of functions from \mathbb{R}^d to \mathbb{R} which has bounded continuous 0-th, ..., k -th order derivatives and by $\mathcal{C}_p^k(\mathbb{R}^d, \mathbb{R})$ the set of functions from \mathbb{R}^d to \mathbb{R} whose 0-th, ..., k -th order derivatives have polynomial growth. For $g \in \mathcal{C}^3(\mathbb{R}^d, \mathbb{R})$ and $\mathbf{v}, \mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{x} \in \mathbb{R}^d$, we denote that

$$\begin{aligned} \nabla_{\mathbf{v}} g(\mathbf{x}) &= \lim_{\varepsilon \rightarrow 0} \frac{g(\mathbf{x} + \varepsilon \mathbf{v}) - g(\mathbf{x})}{\varepsilon}, \\ \nabla_{\mathbf{v}_2} \nabla_{\mathbf{v}_1} g(\mathbf{x}) &= \lim_{\varepsilon \rightarrow 0} \frac{\nabla_{\mathbf{v}_1} g(\mathbf{x} + \varepsilon \mathbf{v}_2) - \nabla_{\mathbf{v}_1} g(\mathbf{x})}{\varepsilon}, \end{aligned}$$

$$\nabla_{\mathbf{v}_3} \nabla_{\mathbf{v}_2} \nabla_{\mathbf{v}_1} g(\mathbf{x}) = \lim_{\varepsilon \rightarrow 0} \frac{\nabla_{\mathbf{v}_2} \nabla_{\mathbf{v}_1} g(\mathbf{x} + \varepsilon \mathbf{v}_3) - \nabla_{\mathbf{v}_2} \nabla_{\mathbf{v}_1} g(\mathbf{x})}{\varepsilon},$$

as the directional derivatives of g . We know $\nabla g(\mathbf{x}) \in \mathbb{R}^d$, $\nabla^2 g(\mathbf{x}) \in \mathbb{R}^{d \times d}$, $\nabla^3 g(\mathbf{x}) \in \mathbb{R}^{d \times d \times d}$. Moreover, we define the operator norm of $\nabla^k g(\mathbf{x})$, $k = 2, 3$ with

$$\|\nabla^k g(\mathbf{x})\|_{\text{op}} = \sup_{|\mathbf{v}_i|=1, i=1, \dots, k} |\nabla_{\mathbf{v}_k} \dots \nabla_{\mathbf{v}_1} g(\mathbf{x})|.$$

If g is bounded, we denote its infinity norm by $\|g\|_{\infty} = \sup_{\mathbf{x}} |g(\mathbf{x})|$. If g is a Lipschitz function, its Lipschitz norm is denoted by $\|g\|_{\text{Lip}} = \sup_{\mathbf{x} \neq \mathbf{y}} |g(\mathbf{x}) - g(\mathbf{y})| / |\mathbf{x} - \mathbf{y}|$. Additionally, for a function $h: (\mathbf{m}, \mathbf{x}) \in \mathbb{R}^{2d} \mapsto h(\mathbf{m}, \mathbf{x}) \in \mathbb{R}$, we denote its derivatives with respect to \mathbf{m} and \mathbf{x} as $\nabla_{\mathbf{m}} h$ and $\nabla_{\mathbf{x}} h$ respectively. Whenever we need to emphasize the initial value $(\mathbf{M}_0, \mathbf{X}_0) = (\mathbf{m}, \mathbf{x})$ for given $(\mathbf{m}, \mathbf{x}) \in \mathbb{R}^{2d}$, we will write $(\mathbf{M}_t^{\mathbf{m}}, \mathbf{X}_t^{\mathbf{x}})$. Similarly, we use the notation $(\widetilde{\mathbf{M}}_t^{\mathbf{m}}, \widetilde{\mathbf{X}}_t^{\mathbf{x}})$. The operator semigroup induced by $(\mathbf{M}_t^{\mathbf{m}}, \mathbf{X}_t^{\mathbf{x}})_{t \geq 0}$ from (1.5) is given by

$$P_t h(\mathbf{m}, \mathbf{x}) = \mathbb{E} h(\mathbf{M}_t^{\mathbf{m}}, \mathbf{X}_t^{\mathbf{x}}), \quad h \in \mathcal{C}_b(\mathbb{R}^{2d}, \mathbb{R}), \quad t > 0.$$

Its infinitesimal generator \mathcal{A} is defined as

$$\mathcal{A} h(\mathbf{m}, \mathbf{x}) := -\langle \nabla_{\mathbf{m}} h, \gamma \mathbf{m} + \nabla f(\mathbf{x}) \rangle + \langle \nabla_{\mathbf{x}} h, \mathbf{m} \rangle + \frac{1}{2} \beta^2 \Delta_{\mathbf{m}} h,$$

for any $h \in \mathcal{C}_p^2(\mathbb{R}^{2d}, \mathbb{R})$. The exact form for the domain of \mathcal{A} is not necessarily figured out in this paper.

Let $\mathcal{P}(\mathbb{R}^{2d})$ denote the space of probability distributions on \mathbb{R}^{2d} . For any $\mu, \nu \in \mathcal{P}(\mathbb{R}^{2d})$, denote their 1-Wasserstein distance by

$$d_{\mathcal{W}_1}(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^{2d} \times \mathbb{R}^{2d}} |\mathbf{z}^1 - \mathbf{z}^2| d\pi(\mathbf{z}^1, \mathbf{z}^2), \quad \mu, \nu \in \mathcal{P}(\mathbb{R}^{2d}),$$

where $\Pi(\mu, \nu)$ denotes the set of probability distributions on $\mathbb{R}^{2d} \times \mathbb{R}^{2d}$ with marginal distributions μ and ν . Kantorovich-Rubinstein Theorem tells us that

$$d_{\mathcal{W}_1}(\mu, \nu) = \sup_{\|h\|_{\text{Lip}} \leq 1} \int_{\mathbb{R}^{2d}} h(\mathbf{z}) d(\mu - \nu)(\mathbf{z}), \quad \mu, \nu \in \mathcal{P}(\mathbb{R}^{2d}).$$

Besides, the total variation distance between μ and ν can be defined by

$$d_{\text{TV}}(\mu, \nu) = \sup_{\|g\|_{\infty} \leq 1} \int_{\mathbb{R}^{2d}} g(\mathbf{z}) d(\mu - \nu)(\mathbf{z}), \quad \mu, \nu \in \mathcal{P}(\mathbb{R}^{2d}).$$

Throughout this paper, we propose the following assumptions.

Assumption I. *The function $f(\mathbf{x})$, defined by (1.3), is non-negative and satisfies the following conditions:*

(i) *There exist constants $A, B \geq 0$ such that*

$$(2.1) \quad |f(\mathbf{0})| \leq A, \quad |\nabla f(\mathbf{0})| \leq B.$$

(ii) *f is L -smooth for some constant $L > 0$, that is,*

$$(2.2) \quad |\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})| \leq L |\mathbf{x} - \mathbf{y}|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

(iii) *For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, there exist positive constants a and b such that*

$$(2.3) \quad \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq a |\mathbf{x} - \mathbf{y}|^2 - b.$$

Assumption II. Let \mathcal{U} be a measurable space. ξ is a \mathcal{U} -valued random variable such that $\nabla F(\mathbf{x}, \xi)$ satisfies the following conditions:

(i) For any $\mathbf{x} \in \mathbb{R}^d$,

$$\mathbb{E}_\xi[\nabla F(\mathbf{x}, \xi)] = \nabla f(\mathbf{x}).$$

(ii) There exists a constant $A_0 > 0$ such that for some $q \geq 2$,

$$(2.4) \quad \sup_{\mathbf{x} \in \mathbb{R}^d} \left\{ \mathbb{E}_\xi [|\nabla f(\mathbf{x}) - \nabla F(\mathbf{x}, \xi)|^q] \right\}^{1/q} \leq A_0.$$

Here, $F(\mathbf{x}, \xi)$ and $f(\mathbf{x})$ are in (1.1) and (1.3) respectively.

In practical applications, the choice of step size often varies with each iteration k . It is necessary and important to add an additional assumption as following to restrict the behavior of η_k .

Assumption III. Let $(\eta_k)_{k \geq 1}$ be a non-increasing and positive sequence and let $t_n = \sum_{k=1}^n \eta_k$, which satisfy the following conditions:

(i) $\lim_{n \rightarrow \infty} t_n = +\infty$;

(ii) There exists a constant $\omega \in (0, 2\theta)$ such that

$$(2.5) \quad \eta_k \leq \frac{2\theta - \omega}{2\theta^2} \quad \text{and} \quad \eta_{k-1} - \eta_k \leq \omega \eta_k^2, \quad \forall k \geq 1,$$

where the constant θ will be given in Lemma 3.6.

A typical example is $\eta_k = \eta/k^\alpha$ for some constants $\eta > 0$ and $\alpha \in (0, 1)$, then condition (2.5) holds for sufficiently large k .

2.2. Main Results. We are now in the position to state our main results, whose proofs will be given in Section 4. Define the Lyapunov function $\mathcal{V}(\mathbf{m}, \mathbf{x}) : \mathbb{R}^{2d} \rightarrow \mathbb{R}$ of SDE (1.5) as

$$(2.6) \quad \mathcal{V}(\mathbf{m}, \mathbf{x}) := f(\mathbf{x}) + \frac{\gamma^2}{4} \left(\left| \mathbf{x} + \frac{1}{\gamma} \mathbf{m} \right|^2 + \left| \frac{1}{\gamma} \mathbf{m} \right|^2 - \lambda |\mathbf{x}|^2 \right),$$

where constant $0 < \lambda \leq \frac{1}{4} \wedge \frac{a}{4L + \gamma^2}$. More details about this Lyapunov function can be found in Appendix A.

We denote by $\mathcal{L}(\mathbf{M}_{t_n}, \mathbf{X}_{t_n})$ and $\mathcal{L}(\mathbf{m}_n, \mathbf{x}_n)$ the laws of $(\mathbf{M}_{t_n}, \mathbf{X}_{t_n})$ from (1.5) and $(\mathbf{m}_n, \mathbf{x}_n)$ from (1.4) respectively.

Throughout this paper, we shall use the letters c, C, C_1, C_2, C_3 to denote positive numbers which may depend on the parameters $A_0, A, B, L, a, b, q, \theta, \omega$ in Assumptions I, II, and III above, and A'_0, q', B_1, B_2 in Assumption IV below. Their values may vary from line to line, but do NOT depend on the dimension d and the sample size N . In some cases, they may depend on other parameters such as p and we will stress this in a way like ' C also depends on p ' if necessary.

Theorem 1. Under Assumptions I, II, and III, we assume that $\eta_1 \leq c$ for some positive constant c , and $\gamma > \sqrt{2}(2L + a)/\sqrt{a}$ additionally. Then we have

$$d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{M}_{t_n}, \mathbf{X}_{t_n}), \mathcal{L}(\mathbf{m}_n, \mathbf{x}_n)) \leq C \sqrt{\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0) + d} \left(1 + \frac{1}{N} \right) \sqrt{\eta_n}, \quad \forall n \geq 0.$$

To estimate the total variation distance, we rely on the following assumptions regarding the derivatives of the function $F(\mathbf{x}, \xi)$, which are stronger than Assumption II.

Assumption IV. Functions $F(\mathbf{x}, \xi)$ and $f(\mathbf{x})$ are in (1.1) and (1.3) respectively. We assume that

(i) For any $\mathbf{x} \in \mathbb{R}^d$,

$$\mathbb{E}_\xi [\nabla F(\mathbf{x}, \xi)] = \nabla f(\mathbf{x}) \quad \text{and} \quad \mathbb{E}_\xi [\nabla^2 F(\mathbf{x}, \xi)] = \nabla^2 f(\mathbf{x}).$$

(ii) There exists a constant $A'_0 > 0$ such that for some $q' \geq 4$,

$$(2.7) \quad \sup_{\mathbf{x} \in \mathbb{R}^d} \left\{ \mathbb{E}_\xi [|\nabla f(\mathbf{x}) - \nabla F(\mathbf{x}, \xi)|^{q'}] \right\}^{1/q'} \leq A'_0.$$

(iii) There exist some positive constants B_1 and B_2 such that

$$\left\{ \mathbb{E} \left[\sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla^2 F(\mathbf{x}, \xi)\|_{\text{op}}^8 \right] \right\}^{1/8} \leq B_1,$$

$$\max \left\{ \sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla^3 f(\mathbf{x})\|_{\text{op}}, \sup_{\mathbf{x} \in \mathbb{R}^d} \left(\mathbb{E} \|\nabla^3 F(\mathbf{x}, \xi)\|_{\text{op}}^4 \right)^{1/4} \right\} \leq B_2.$$

Note that (2.7) covers the (2.4). We have the following result about total variation distance.

Theorem 2. Under Assumptions I, II, III and IV, we assume that $\eta_1 \leq cd^{-2}$ for some constant c , and $\gamma > \sqrt{2}(2L + a)/\sqrt{a}$ additionally. Then we have

$$d_{\text{TV}}(\mathcal{L}(\mathbf{M}_{t_n}, \mathbf{X}_{t_n}), \mathcal{L}(\mathbf{m}_n, \mathbf{x}_n)) \leq Cd^{7/2} \sqrt{\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0) + d} \left(\sqrt{\eta_n} + \frac{1}{\sqrt{N}} \right), \quad \forall n \geq 0.$$

Above two results yield the following corollaries for the case of $\eta_k = \eta/k^\alpha$ for some constant $\eta > 0$ and $\alpha \in (0, 1)$.

Corollary 3. Under the same conditions in Theorem 1, let η be a (small) positive constant and $\eta_k = \eta/k^\alpha$ with $\alpha \in (0, 1)$. Then, there exists a constant C such that for all $n \geq 0$

$$d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{M}_{t_n}, \mathbf{X}_{t_n}), \mathcal{L}(\mathbf{m}_n, \mathbf{x}_n)) \leq C \sqrt{\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0) + d} \left(1 + \frac{1}{N} \right) \sqrt{\frac{\eta}{n^\alpha}}.$$

Corollary 4. Under the same conditions in Theorem 2, let η be a (small) positive constant and $\eta_k = \eta/k^\alpha$ with $\alpha \in (0, 1)$. Then, there exists a constant C such that for all $n \geq 0$

$$d_{\text{TV}}(\mathcal{L}(\mathbf{M}_{t_n}, \mathbf{X}_{t_n}), \mathcal{L}(\mathbf{m}_n, \mathbf{x}_n)) \leq Cd^{7/2} \sqrt{\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0) + d} \left(\sqrt{\frac{\eta}{n^\alpha}} + \frac{1}{\sqrt{N}} \right).$$

With regard to the original minimization problem, we have the following generalization error bound.

Corollary 5. Under the same conditions in Theorem 1, let $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$ be the minimizer of $f(\mathbf{x})$ in (1.3). It holds

$$\mathbb{E}F(\mathbf{x}_n, \xi) - f(\mathbf{x}^*) \leq \left\{ C \left[\left(1 + \frac{1}{N} \right) \sqrt{\eta_n} \right]^{\frac{q-2}{q-1}} + B_e e^{-\kappa t_n} \right\} (\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0) + d) + Jd,$$

where B_e and κ are the constants in (A.3), q is in Assumption II, and J has the form of

$$J = \frac{\beta^2}{4} \log \left[\frac{2eL}{a} \left(\frac{2ab + B^2}{da\beta^2} + 1 \right) \right].$$

That is, the SGDm described by (1.4) is expected to approach the minimum point of the function $f(\mathbf{x})$ under appropriate parameters by choosing sufficiently large time n and large sample size N , and sufficiently small β . Moreover, the larger q in Assumption II is, the faster convergence we can obtain.

Remark 6. (i). Note that these results above hold for constant learning rate $\eta_n \equiv \eta$ for all $n \in \mathbb{N}$ as long as the η is a sufficiently small number.

(ii). In the algorithm (1.4), the term $N^{-1} \sum_{i=1}^N \nabla F(x, \xi_{k+1}^i)$ is an unbiased estimator of $\nabla f(x)$. According to the law of large number, it converges to $\nabla f(x)$ almost surely as $N \rightarrow \infty$. This, together with the exponential ergodicity of $(\mathbf{M}_t, \mathbf{X}_t)$, immediately implies an error bound in the order $O(\sqrt{\eta_n})$ in the 1-Wasserstein and total variation distance for the underdamped Langevin sampling as n is large. The rate $O(\sqrt{\eta_n})$ may be not optimal due to the heavy tail effect of ξ_k^i , we conjecture that the rate $O(\sqrt{\eta_n})$ can be improved by assuming that ξ_k^i has a high order moment. We leave this possible improvement to the future research.

3. AUXILIARY LEMMAS

In this section, we list auxiliary lemmas which will be used to prove our main results, and their proofs will be given in Appendix B.

Combining (2.1) and (2.3) implies the following dissipation property of f

$$\langle \mathbf{x}, \nabla f(\mathbf{x}) \rangle \geq \frac{a}{2} |\mathbf{x}|^2 - K, \quad \forall \mathbf{x} \in \mathbb{R}^d,$$

where $K = b + B^2/(2a)$. By (2.2), the following linear growth condition holds

$$(3.1) \quad |\nabla f(\mathbf{x})| \leq L |\mathbf{x}| + B, \quad \forall \mathbf{x} \in \mathbb{R}^d.$$

Recall the Lyapunov function (2.6). Since $\lambda \leq 1/4$ and f is nonnegative, it is easy to see

$$(3.2) \quad 0 \leq \max \left\{ \frac{1-2\lambda}{4(1-\lambda)} |\mathbf{m}|^2, \frac{\gamma^2}{8} (1-2\lambda) |\mathbf{x}|^2 \right\} \leq \mathcal{V}(\mathbf{m}, \mathbf{x}), \quad \forall \mathbf{m}, \mathbf{x} \in \mathbb{R}^d.$$

One can verify that

$$(3.3) \quad \mathcal{A}\mathcal{V}(\mathbf{m}, \mathbf{x}) \leq -\lambda\gamma\mathcal{V}(\mathbf{m}, \mathbf{x}) + \frac{1}{2} (\gamma\mathring{A} + d\beta^2),$$

where constants λ and \mathring{A} satisfy

$$0 < \lambda \leq \min \left\{ \frac{1}{4}, \frac{a}{4L + \gamma^2} \right\}, \quad \mathring{A} \geq K + 2\lambda \left(\frac{B^2}{2L} + A \right).$$

More details about the Lyapunov function can be found in Appendix A.

In order to estimate SDE (1.6), we need to introduce the following auxiliary one step SDE: for $t \in [0, \eta]$ with η being a step size,

$$(3.4) \quad \begin{cases} d\widetilde{\mathbf{M}}_t = -\gamma\widetilde{\mathbf{M}}_0 dt - \frac{1}{N} \sum_{i=1}^N \nabla F(\widetilde{\mathbf{X}}_0, \xi^i) dt + \beta d\mathbf{B}_t, \\ d\widetilde{\mathbf{X}}_t = \widetilde{\mathbf{M}}_0 dt, \end{cases}$$

where $(\widetilde{\mathbf{M}}_0, \widetilde{\mathbf{X}}_0) = (\mathbf{m}, \mathbf{x})$, and ξ^1, \dots, ξ^N are i.i.d. and satisfy Assumption II.

3.1. Moments Estimates. We give in this subsection the moments estimates for \mathbf{M}_t , \mathbf{X}_t , $\widetilde{\mathbf{M}}_{t_n}$ and $\widetilde{\mathbf{X}}_{t_n}$.

Lemma 3.1. *Under Assumption I, for each $p \geq 1$, there exists a positive number C , which also depends on p , such that for any $t \geq 0$, and initial value (\mathbf{m}, \mathbf{x}) ,*

$$\mathbb{E} [\mathcal{V}(\mathbf{M}_t^{\mathbf{m}}, \mathbf{X}_t^{\mathbf{x}})^p] \leq e^{-\lambda\gamma t} \mathcal{V}(\mathbf{m}, \mathbf{x})^p + Cd^p.$$

Combining this lemma with (3.2), for each $p \geq 1$, there exists a positive number C , which also depends on p , such that for all $t \geq 0$ and initial value (\mathbf{m}, \mathbf{x})

$$(3.5) \quad \mathbb{E} |\mathbf{M}_t^{\mathbf{m}}|^{2p} + \mathbb{E} |\mathbf{X}_t^{\mathbf{x}}|^{2p} \leq C (e^{-\lambda\gamma t} \mathcal{V}(\mathbf{m}, \mathbf{x})^p + d^p).$$

Then, using (3.1) and (3.5), we can obtain the following estimates.

Lemma 3.2. *Under Assumption I, for each $p \geq 1$, there exists a positive number C , which also depends on p , such that for all $t \in [0, 1]$ and initial value (\mathbf{m}, \mathbf{x}) :*

$$\mathbb{E} |\mathbf{M}_t^{\mathbf{m}} - \mathbf{m}|^{2p} + \mathbb{E} |\mathbf{X}_t^{\mathbf{x}} - \mathbf{x}|^{2p} \leq Ct^p (\mathcal{V}(\mathbf{m}, \mathbf{x})^p + d^p).$$

What's more, we have the similar consequences for SDE (3.4),

Lemma 3.3. *Consider SDE (3.4). Under Assumptions I and II, for each $1 \leq p \leq q/2$, there exists some positive number C , which also depends on p , such that,*

$$(3.6) \quad \mathbb{E} |\widetilde{\mathbf{M}}_t^{\mathbf{m}} - \mathbf{m}|^{2p} + \mathbb{E} |\widetilde{\mathbf{X}}_t^{\mathbf{x}} - \mathbf{x}|^{2p} \leq Ct^p (\mathcal{V}(\mathbf{m}, \mathbf{x})^p + d^p),$$

for all $t \in [0, \eta]$ with η being a step size and initial value (\mathbf{m}, \mathbf{x}) . If the condition (2.7) in Assumption IV holds additionally, (3.6) holds for any $1 \leq p \leq q'/2$.

On the other hand, using the estimate for SDE (3.4) in an inductive way, SDE (1.6) has moments estimates similar to Lemma 3.1.

Lemma 3.4. *Under Assumptions I, II and III, let $\eta_1 \leq c$ for some positive constant c , and $t_n = \sum_{k=1}^n \eta_k$. For each $1 \leq p \leq q/2$, there exists a positive number C , which also depends on p , such that*

$$(3.7) \quad \mathbb{E} \left[\mathcal{V}(\widetilde{\mathbf{M}}_{t_n}^{\mathbf{m}}, \widetilde{\mathbf{X}}_{t_n}^{\mathbf{x}})^p \right] \leq e^{-\frac{\lambda\gamma}{2} t_n} \mathcal{V}(\mathbf{m}, \mathbf{x})^p + Cd^p,$$

for all $n \in \mathbb{N}$ and initial value (\mathbf{m}, \mathbf{x}) . If the condition (2.7) in Assumption IV holds additionally, (3.7) holds for any $1 \leq p \leq q'/2$.

3.2. Auxiliary Lemmas for 1-Wasserstein Distance. The following lemma provides a bound for the difference between $(\mathbf{M}_\eta^{\mathbf{m}}, \mathbf{X}_\eta^{\mathbf{x}})$ and $(\widetilde{\mathbf{M}}_\eta^{\mathbf{m}}, \widetilde{\mathbf{X}}_\eta^{\mathbf{x}})$.

Lemma 3.5. *Consider SDE (3.4). Under Assumptions I and II, there exists a positive number $C > 0$ such that for any $\eta \in (0, 1)$,*

$$(3.8) \quad d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{M}_\eta^{\mathbf{m}}, \mathbf{X}_\eta^{\mathbf{x}}), \mathcal{L}(\widetilde{\mathbf{M}}_\eta^{\mathbf{m}}, \widetilde{\mathbf{X}}_\eta^{\mathbf{x}})) \leq C\eta^{3/2} \left(1 + \frac{1}{N} \right) \sqrt{\mathcal{V}(\mathbf{m}, \mathbf{x}) + d}.$$

Let the condition (2.7) in Assumption IV hold additionally, we have

$$(3.9) \quad \mathbb{E} |\mathbf{M}_\eta^{\mathbf{m}} - \widetilde{\mathbf{M}}_\eta^{\mathbf{m}}|^4 + \mathbb{E} |\mathbf{X}_\eta^{\mathbf{x}} - \widetilde{\mathbf{X}}_\eta^{\mathbf{x}}|^4 \leq C\eta^6 (\mathcal{V}(\mathbf{m}, \mathbf{x})^2 + d^2) + C \frac{\eta^4}{N^2}.$$

[Sch24] obtained a global contractivity for Langevin dynamics which will help us to prove the main result. Conveniently, we briefly introduce it here.

Recall the condition (2.3)

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq a |\mathbf{x} - \mathbf{y}|^2 - b, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d,$$

then it holds

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \frac{a}{2} |\mathbf{x} - \mathbf{y}|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d \text{ such that } |\mathbf{x} - \mathbf{y}|^2 \geq \frac{2b}{a}.$$

Together with $|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})| \leq L$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, f is a potential function with a L -Lipschitz continuous gradient and that is $a/2$ -strongly convex outside a Euclidean ball of radius $\sqrt{2b/a}$. Then, [Sch24, Theorem 5] immediately implies the following lemma when γ is sufficiently large such that

$$\gamma^2 > \frac{2(2L + a)^2}{a}.$$

Lemma 3.6. *Suppose that $(\mathbf{M}_t^1, \mathbf{X}_t^1)_{t \geq 0}$ and $(\mathbf{M}_t^2, \mathbf{X}_t^2)_{t \geq 0}$ satisfy SDE (1.5) with different initial values $(\mathbf{m}^1, \mathbf{x}^1)$ and $(\mathbf{m}^2, \mathbf{x}^2)$ respectively. Let Assumption I hold and assume $\gamma > \sqrt{2}(2L + a)/\sqrt{a}$ additionally. Then there exist positive constants C and θ (θ does not depend on d) such that*

$$d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{M}_t^1, \mathbf{X}_t^1), \mathcal{L}(\mathbf{M}_t^2, \mathbf{X}_t^2)) \leq C e^{-\theta t} (|\mathbf{m}^1 - \mathbf{x}^1| + |\mathbf{m}^2 - \mathbf{x}^2|),$$

for all $t \geq 0$.

3.3. Auxiliary Lemmas for Total Variation Distance. The following lemma is a direct application of Zhang et al. [Zha10], and we will give the details in Appendix B.3.

Lemma 3.7. *For any fixed $T > 0$, Let process $\mathbf{Y}_t = (\mathbf{M}_t, \mathbf{X}_t)_{t \in [0, T]}$ come from SDE (1.5) and let Assumption I hold. Then, for any function $\phi \in C_b^1(\mathbb{R}^{2d}, \mathbb{R})$, there exists a positive number $C > 0$ such that*

$$|\nabla \mathbb{E} \phi(\mathbf{Y}_T)| \leq C \|\phi\|_\infty (T^{3/2} \vee T^{-3/2}).$$

To estimate the total variation distance, we need to use Malliavin calculus. Let us briefly introduce its preliminary in our setting. More details can be found in [Nua06].

Denote $\mathcal{H} = L^2([0, T]; \mathbb{R}^d)$. Let $\mathbf{W} = (W_t^1, \dots, W_t^d)_{t \geq 0}$ be a d -dimensional Wiener process (a.k.a. Brownian motion) on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, where \mathcal{F} is the σ -field generated by \mathbf{W} . For any $\mathbf{h} = (h_1, \dots, h_d) \in \mathcal{H}$, define the Wiener integral as

$$\mathbf{W}(\mathbf{h}) = \sum_{i=1}^d \int_0^T h_i(t) dW_t^i.$$

We denote by $\mathcal{C}_p^\infty(\mathbb{R}^m, \mathbb{R})$ the set of infinitely differentiable functions $g : \mathbb{R}^m \rightarrow \mathbb{R}$ such that g and all of its partial derivatives have polynomial growth. Denote by \mathcal{S} the set of random variables in $L^2(\Omega)$ with the form:

$$G = g(\mathbf{W}(\mathbf{h}_1), \dots, \mathbf{W}(\mathbf{h}_m)),$$

where $\mathbf{h}_i \in \mathcal{H}$, $i = 1, \dots, m$ for $m \in \mathbb{N}$. Then the first order Malliavin derivative of G is the \mathcal{H} -valued random variable given by

$$D_t G = \sum_{j=1}^m \partial_j g(\mathbf{W}(\mathbf{h}_1), \dots, \mathbf{W}(\mathbf{h}_m)) \mathbf{h}_j(t), \quad 0 \leq t \leq T.$$

So $D_t G \in L^2(\Omega, \mathcal{H})$. We can further define the second order Malliavin derivative of G as the following

$$D_{t_1} D_{t_2} G = \sum_{j_1=1}^m \sum_{j_2=1}^m \partial_{j_1} \partial_{j_2} g(\mathbf{W}(\mathbf{h}_1), \dots, \mathbf{W}(\mathbf{h}_m)) \mathbf{h}_{j_1}(t_1) \mathbf{h}_{j_2}(t_2), \quad 0 \leq t_1, t_2 \leq T.$$

So $D_{t_1} D_{t_2} G \in L^2(\Omega, \mathcal{H} \otimes \mathcal{H})$. Inductively, we can define the k -th order Malliavin derivative $D_{t_1} \dots D_{t_k} G$, which is located in $L^2(\Omega, \mathcal{H}^{\otimes k})$. Define the following norm

$$\|G\|_{k,p} := \left[\mathbb{E} |G|^p + \sum_{j=1}^k \mathbb{E} \|D_{t_1} \dots D_{t_j} G\|_{\mathcal{H}^{\otimes j}}^p \right]^{1/p}.$$

Under this norm, the operator D can be extended from \mathcal{S} to its domain, denoted by $\mathbb{D}^{k,p}$. Specially, $\mathbb{D}^{1,2}$ is also a Hilbert space with product

$$\langle G, H \rangle_{1,2} = \mathbb{E}(GH) + \mathbb{E} \langle DG, DH \rangle_{\mathcal{H}}, \quad \forall G, H \in \mathbb{D}^{1,2}.$$

The following relation is called integration by parts in Malliavin calculus:

$$(3.10) \quad \mathbb{E} [G \delta(\mathbf{h})] = \mathbb{E} [\langle DG, \mathbf{h} \rangle_{\mathcal{H}}], \quad \forall G \in \mathbb{D}^{1,2}, \quad \mathbf{h} \in \mathcal{H},$$

where $\delta(\mathbf{h})$ is called Skorohod integral.

Given $\mathbf{h} \in \mathcal{H}$, we can define the Malliavin derivative along the direction \mathbf{h} , denoted by $D^{\mathbf{h}}G$, as the following:

$$D^{\mathbf{h}}G = \langle DG, \mathbf{h} \rangle_{\mathcal{H}} = \int_0^T \langle \mathbf{h}(s), D_s G \rangle ds.$$

For $\mathbf{G} = (G^1, \dots, G^d)^\top$ with each $G^i \in \mathbb{D}^{1,2}$, we define $D_t \mathbf{G} = (D_t G^1, \dots, D_t G^d)$. And, the norm

$$\|D_{t_1} \dots D_{t_k} \mathbf{G}\|_{\mathcal{H}^{\otimes k}}^2 = \sum_{i=1}^d \|D_{t_1} \dots D_{t_k} G^i\|_{\mathcal{H}^{\otimes k}}^2, \quad \|\mathbf{G}\|_{k,p}^p = \sum_{i=1}^d \|G^i\|_{k,p}^p.$$

The associated Malliavin matrix of \mathbf{G} is defined as the following random semi-definite symmetric matrix

$$\mathbf{\Gamma}(\mathbf{G}) = (\langle DG^i, DG^j \rangle_{\mathcal{H}})_{1 \leq i, j \leq d}.$$

The following abstract lemma will play an important role in the proof of Theorem 2. We leave its proof in Appendix B.3.

Lemma 3.8. *Let $\mathbf{F} = (F^1, \dots, F^d)$ be a random vector such that all of its components $F^i \in \mathbb{D}^{2,8}$, $i = 1, \dots, d$, and its Malliavin matrix $\mathbf{\Gamma}(\mathbf{F})$ is invertible a.s. with $(\det \mathbf{\Gamma}(\mathbf{F}))^{-1} \in \bigcap_{p \geq 1} L^p(\Omega)$. Let $\mathbf{G} = (G^1, \dots, G^d)$ be another random vector with $G^i \in \mathbb{D}^{1,4}$ for all $1 \leq i \leq d$, and $g : \mathbb{R}^d \rightarrow \mathbb{R}$ be a function in $\mathcal{C}_b^1(\mathbb{R}^d, \mathbb{R})$. Then there exists a positive number $C > 0$ such that*

$$(3.11) \quad |\mathbb{E} \langle \nabla g(\mathbf{F}), \mathbf{G} \rangle| \leq C \|g\|_{\infty} \|\mathbf{G}\|_{1,4} \{\mathbb{E}[\mathcal{K}_1 \mathcal{K}_2 \mathcal{K}_3]\}^{1/4},$$

where

$$\mathcal{K}_1 = 1 + \|\mathbf{D}\mathbf{F}\|_{\mathcal{H}}^8, \quad \mathcal{K}_2 = 1 + \|\mathbf{D}^2\mathbf{F}\|_{\mathcal{H}\otimes\mathcal{H}}^4, \quad \mathcal{K}_3 = 1 + \|\mathbf{\Gamma}(\mathbf{F})^{-1}\|_{\text{HS}}^8.$$

In order to apply Lemma 3.8, we need to introduce a series of events to estimate $\mathbb{E}[\mathcal{K}_1\mathcal{K}_2\mathcal{K}_3]$. To this end, let us first recall that for each step $i \in \mathbb{N}$ of SGDm (1.4), ξ_i^r , $r = 1, \dots, N$ are independent copies of random variable ξ satisfying Assumption II, to make notations simple, we denote $\boldsymbol{\xi}_i = (\xi_i^1, \dots, \xi_i^N)$, $i \in \mathbb{N}$, and define

$$(3.12) \quad \varphi(\boldsymbol{\xi}_i) := \frac{1}{N} \sum_{r=1}^N \sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla^2 F(\mathbf{x}, \xi_i^r)\|_{\text{op}}.$$

Define the following events: for $j \leq k$ and $\tau = 1, 2$,

$$(3.13) \quad E_{j,k}^\tau := \left\{ \sum_{i=j+1}^k \eta_i (\varphi(\boldsymbol{\xi}_i)^\tau - \mathbb{E}[\varphi(\boldsymbol{\xi}_i)^\tau]) > t_k - t_j \right\}.$$

Lemma 3.9. *Given Assumptions III and IV, let $t_k = \sum_{i=1}^k \eta_i$. There exist positive constants c and C such that*

$$\mathbb{P}(E_{j,k}^\tau) \leq \frac{C e^{c(t_k - t_j)}}{(t_k - t_j)^2} \eta_k^2, \quad \tau = 1, 2.$$

For notations simplicity, we denote

$$(3.14) \quad \mathbf{Y}_t = (\mathbf{M}_t, \mathbf{X}_t), \quad \widetilde{\mathbf{Y}}_t = (\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t), \quad \forall t \geq 0,$$

where $(\mathbf{M}_t, \mathbf{X}_t)$ and $(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)$ are defined by SDEs (1.5) and (1.6) respectively. For any $0 \leq s \leq t$, denote $\mathbf{Y}_{s,t}^{\mathbf{y}}$ or $\mathbf{Y}_{s,t}(\mathbf{y})$ as the value of the process $(\mathbf{Y}_t)_{t \geq 0}$ at the time t given $\mathbf{Y}_s = \mathbf{y} \in \mathbb{R}^{2d}$. Similarly for $\widetilde{\mathbf{Y}}_{s,t}^{\mathbf{y}}$. When $s = 0$, we will drop the subscript s if there is no ambiguity.

Lemma 3.10. *Let Assumptions I, II, III, and IV hold and recall $t_k = \sum_{i=1}^k \eta_i$. There exists a positive number $C > 0$ such that for any $\mathbf{z} \in \mathbb{R}^{2d}$ and any $\ell < k$ such that $t_k - t_\ell \leq 1/[5(B_1 + \gamma + 2)]$, the followings hold:*

(i) For all $t \in [t_{k-1}, t_k]$

$$\mathbb{E} \left\| \mathbf{D}\widetilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} \right\|_{\mathcal{H}}^8 \leq Cd^4, \quad \mathbb{E} \left\| \mathbf{D}\mathbf{Y}_{t_{k-1}, t}(\widetilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}}^8 \leq Cd^4.$$

(ii) We have

$$\|\mathbf{D}\widetilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}}\|_{\mathcal{H}} \leq C\sqrt{d} \quad \text{on the event } (E_{\ell, k}^1)^c.$$

Lemma 3.11. *Under the same setting as in Lemma 3.10, there exists a positive number $C > 0$ such that for any $\mathbf{z} \in \mathbb{R}^{2d}$ and any $\ell < k$ such that $t_k - t_\ell \leq 1/[5(B_1 + \gamma + 2)]$, there exists a positive number $C > 0$ such that the followings hold:*

$$\mathbb{E} \left\| \mathbf{D}^2 \widetilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} \right\|_{\mathcal{H}\otimes\mathcal{H}}^4 \leq Cd^4, \quad \mathbb{E} \left\| \mathbf{D}^2 \mathbf{Y}_{t_{k-1}, t_k}(\widetilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}\otimes\mathcal{H}}^4 \leq Cd^4.$$

Lemma 3.12. *Under the same setting as in Lemma 3.10, there exists a positive number $C > 0$ such that for any $\mathbf{z} \in \mathbb{R}^{2d}$ and any $\ell < k$ such that $t_k - t_\ell \leq 1/[5(B_1 + \gamma + 2)]$, there exists a positive number $C > 0$ such that the followings hold:*

(i) We have

$$\mathbb{E} \left\| \mathbf{D}\mathbf{Y}_{t_{k-1}, t_k}(\widetilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - \mathbf{D}\widetilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} \right\|_{\mathcal{H}}^4 \leq Cd^2 \eta_k^6 (\mathcal{V}(\mathbf{z})^2 + d^2) + C \frac{d^4 \eta_k^4}{N^2}.$$

(ii) We have

$$\|D\mathbf{Y}_{t_{k-1},t_k}(\tilde{\mathbf{Y}}_{t_\ell,t_{k-1}}^{\mathbf{z}}) - D\tilde{\mathbf{Y}}_{t_\ell,t_k}^{\mathbf{z}}\|_{\mathcal{H}} \leq C\sqrt{d\eta_k} \quad \text{on the event } (E_{\ell,k}^1 \cup E_{\ell,k}^2)^c.$$

Lemma 3.13. *Under the same setting as in Lemma 3.10, for any $\mathbf{z} \in \mathbb{R}^{2d}$ and any $\ell < k$ such that $t_k - t_\ell \leq 1/[5(B_1 + \gamma + 2)]$, there exists a positive number $C > 0$ such that*

$$\lambda_{\min} \left(\Gamma(\tilde{\mathbf{Y}}_{t_\ell,t_k}^{\mathbf{z}}) \right) \geq C(t_k - t_\ell)^3 \quad \text{on the event } (E_{\ell,k}^1 \cup E_{m,k}^1)^c,$$

where $m = \min\{j: t_k - t_j \leq (t_k - t_\ell)/[10(B_1 + \gamma + 2)]\}$ and $\lambda_{\min}(\Gamma(\tilde{\mathbf{Y}}_{t_\ell,t_k}^{\mathbf{z}}))$ is the smallest eigenvalue of the Malliavin matrix $\Gamma(\tilde{\mathbf{Y}}_{t_\ell,t_k}^{\mathbf{z}})$.

4. PROOFS OF THE MAIN RESULTS

In this section, we provide the proofs of Theorems 1, 2, and Corollaries 3, 5 by using the lemmas mentioned above. Employing the Lindeberg principle, we decompose $\mathbb{E}[g(\mathbf{Y}_{t_k})] - \mathbb{E}[g(\tilde{\mathbf{Y}}_{t_k})]$ into k terms and subsequently analyze each term individually, where \mathbf{Y}_t and $\tilde{\mathbf{Y}}_t$ are defined by (3.14).

For any $0 \leq s \leq t$, recall the notations $\mathbf{Y}_{s,t}^{\mathbf{y}}$, $\mathbf{Y}_{s,t}(\mathbf{y})$ and $\tilde{\mathbf{Y}}_{s,t}^{\mathbf{y}}$ immediately below (3.14).

4.1. Proof of Theorem 1.

Proof of Theorem 1. Applying the Lindeberg technique, we have the following decomposition for any Lipschitz function $h: \mathbb{R}^{2d} \rightarrow \mathbb{R}$,

$$(4.1) \quad \mathbb{E}h(\mathbf{Y}_{t_n}^{\mathbf{y}}) - \mathbb{E}h(\tilde{\mathbf{Y}}_{t_n}^{\mathbf{y}}) = \sum_{k=1}^n \mathbb{E} \left[h \left(\mathbf{Y}_{t_{k-1},t_n}(\tilde{\mathbf{Y}}_{t_{k-1}}^{\mathbf{y}}) \right) - h \left(\mathbf{Y}_{t_k,t_n}(\tilde{\mathbf{Y}}_{t_k}^{\mathbf{y}}) \right) \right]$$

with initial value $\mathbf{y} = (\mathbf{m}_0, \mathbf{x}_0)$.

Denote $\mathbf{z} = \mathbf{Y}_{t_{k-1},t_k}(\tilde{\mathbf{Y}}_{t_{k-1}}^{\mathbf{y}})$ and $\tilde{\mathbf{z}} = \tilde{\mathbf{Y}}_{t_k}^{\mathbf{y}}$. We apply Lemma 3.6 on $[t_k, t_n]$ to get

$$\begin{aligned} & \left| \mathbb{E} \left[h \left(\mathbf{Y}_{t_{k-1},t_n}(\tilde{\mathbf{Y}}_{t_{k-1}}^{\mathbf{y}}) \right) - h \left(\mathbf{Y}_{t_k,t_n}(\tilde{\mathbf{Y}}_{t_k}^{\mathbf{y}}) \right) \right] \right| \\ &= \left| \mathbb{E} [h(\mathbf{Y}_{t_k,t_n}(\mathbf{z})) - h(\mathbf{Y}_{t_k,t_n}(\tilde{\mathbf{z}}))] \right| \\ &= \left| \mathbb{E} [h(\mathbf{Y}_{t_k,t_n}(\mathbf{z}^*)) - h(\mathbf{Y}_{t_k,t_n}(\tilde{\mathbf{z}}^*))] \right| \\ &\leq C \|h\|_{\text{Lip}} e^{-\theta(t_n-t_k)} \mathbb{E}|\mathbf{z}^* - \tilde{\mathbf{z}}^*|, \end{aligned}$$

where $(\mathbf{z}^*, \tilde{\mathbf{z}}^*)$ is a random vector whose first and second marginal distributions are the same as \mathbf{z} and $\tilde{\mathbf{z}}$ respectively. Because this inequality also holds for any \mathbf{z}^* and $\tilde{\mathbf{z}}^*$ with this property, we can choose \mathbf{z}^* and $\tilde{\mathbf{z}}^*$ which satisfies $\mathbb{E}|\mathbf{z}^* - \tilde{\mathbf{z}}^*| = d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{z}), \mathcal{L}(\tilde{\mathbf{z}}^*))$ (It is well known that the minimum of the coupling in the 1-Wasserstein can be realized). Hence,

$$\left| \mathbb{E} \left[h \left(\mathbf{Y}_{t_{k-1},t_n}(\tilde{\mathbf{Y}}_{t_{k-1}}^{\mathbf{y}}) \right) - h \left(\mathbf{Y}_{t_k,t_n}(\tilde{\mathbf{Y}}_{t_k}^{\mathbf{y}}) \right) \right] \right| \leq C \|h\|_{\text{Lip}} e^{-\theta(t_n-t_k)} d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{z}), \mathcal{L}(\tilde{\mathbf{z}})),$$

Applying Lemma 3.4, it holds

$$\mathbb{E} \left[\mathcal{V}(\tilde{\mathbf{Y}}_{t_{k-1}}^{\mathbf{y}}) \right] \leq e^{-\lambda\gamma t_{k-1}} \mathcal{V}(\mathbf{y}) + Cd.$$

Furthermore, combining this with (3.8) in Lemma 3.5 immediately yields that

$$\begin{aligned} d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{z}), \mathcal{L}(\tilde{\mathbf{z}})) &\leq C\eta_k^{3/2} \left(1 + \frac{1}{N}\right) \mathbb{E} \left[\sqrt{\mathcal{V}(\tilde{\mathbf{Y}}_{t_{k-1}}^{\mathbf{y}}) + d} \right] \\ &\leq C\eta_k^{3/2} \left(1 + \frac{1}{N}\right) \sqrt{\mathcal{V}(\mathbf{y}) + d}. \end{aligned}$$

Hence,

$$(4.2) \quad \begin{aligned} &\left| \mathbb{E} \left[h \left(\mathbf{Y}_{t_{k-1}, t_n}(\tilde{\mathbf{Y}}_{t_{k-1}}^{\mathbf{y}}) \right) - h \left(\mathbf{Y}_{t_k, t_n}(\tilde{\mathbf{Y}}_{t_k}^{\mathbf{y}}) \right) \right] \right| \\ &\leq C \|h\|_{\text{Lip}} \eta_k^{3/2} e^{-\theta(t_n - t_k)} \left(1 + \frac{1}{N}\right) \sqrt{\mathcal{V}(\mathbf{y}) + d}. \end{aligned}$$

Substituting this formula in (4.1), we have

$$(4.3) \quad \left| \mathbb{E}h(\mathbf{Y}_{t_n}^{\mathbf{y}}) - \mathbb{E}h(\tilde{\mathbf{Y}}_{t_n}^{\mathbf{y}}) \right| \leq C \|h\|_{\text{Lip}} \left(1 + \frac{1}{N}\right) \sqrt{\mathcal{V}(\mathbf{y}) + d} \sum_{k=1}^n \eta_k^{3/2} e^{-\theta(t_n - t_k)}.$$

For the sum term, applying Lemma A.3 with $\varepsilon = 1/2$ therein, we have

$$\sum_{k=1}^n \eta_k^{3/2} e^{-\theta(t_n - t_k)} \leq \frac{4}{2\theta - \omega} \sqrt{\eta_n}.$$

Consequently, the desired result holds since $(\mathbf{m}_n, \mathbf{x}_n)_{n \geq 0}$ and $(\tilde{\mathbf{M}}_n, \tilde{\mathbf{X}}_n)_{n \geq 0}$ have the same distribution. \square

4.2. Proof of Theorem 2. We denote the operator semigroups of $(\mathbf{Y}_t)_{t \geq 0}$ and $(\tilde{\mathbf{Y}}_t)_{t \geq 0}$ by $P_{s,t}$ and $Q_{s,t}$ respectively, defined by

$$P_{s,t} h(\mathbf{y}) = \mathbb{E} [h(\mathbf{Y}_t) | \mathbf{Y}_s = \mathbf{y}], \quad Q_{s,t} h(\mathbf{y}) = \mathbb{E} [h(\tilde{\mathbf{Y}}_t) | \tilde{\mathbf{Y}}_s = \mathbf{y}], \quad 1 \leq s < t,$$

for any $h \in \mathcal{C}_b(\mathbb{R}^{2d}, \mathbb{R})$ and $\mathbf{y} \in \mathbb{R}^{2d}$. As $s = 0$, we will drop the subscript s if no confusions arise.

Proof of Theorem 2. For any $h \in \mathcal{C}_b(\mathbb{R}^{2d}, \mathbb{R})$, by the Lindeberg principle in the form of semigroup, we have the following decomposition

$$\mathbb{E}h(\mathbf{Y}_{t_n}^{\mathbf{y}}) - \mathbb{E}h(\tilde{\mathbf{Y}}_{t_n}^{\mathbf{y}}) = (P_{0,t_n} - Q_{0,t_n}) h(\mathbf{y}) = \sum_{k=1}^n Q_{0,t_{k-1}} \circ (P_{t_{k-1}, t_k} - Q_{t_{k-1}, t_k}) \circ P_{t_k, t_n} h(\mathbf{y}).$$

Define index k_n as

$$k_n = \inf \{k : t_n - t_k < 1\}.$$

Then, the above equation can be divided into the following two parts:

$$(4.4) \quad \begin{aligned} \mathcal{J}_1 &= \sum_{k=1}^{k_n-1} Q_{0,t_{k-1}} \circ (P_{t_{k-1}, t_k} - Q_{t_{k-1}, t_k}) \circ P_{t_k, t_n} h(\mathbf{y}), \\ \mathcal{J}_2 &= \sum_{k=k_n}^n Q_{0,t_{k-1}} \circ (P_{t_{k-1}, t_k} - Q_{t_{k-1}, t_k}) \circ P_{t_k, t_n} h(\mathbf{y}). \end{aligned}$$

For \mathcal{J}_1 and each $1 \leq k < k_n$, we let $P_{t_k, t_n} h(\mathbf{y}) = P_{t_k, t_{k_n}} \circ P_{t_{k_n}, t_n} h(\mathbf{y})$, and define function $g_1(\mathbf{y}) := P_{t_{k_n}, t_n} h(\mathbf{y})$. Then g_1 is a Lipschitz function according to Lemma 3.7 with $T = t_n - t_{k_n} \in [1/2, 1]$ therein, and it holds

$$\|g_1\|_{\text{Lip}} = |\nabla P_{t_{k_n}, t_n} h(\mathbf{y})| \leq C \|h\|_{\infty}.$$

Similar to obtaining (4.2) in the proof of Theorem 1, we have

$$\begin{aligned} & \left| \mathbb{Q}_{0,t_{k-1}} \circ \left(\mathbb{P}_{t_{k-1},t_k} - \mathbb{Q}_{t_{k-1},t_k} \right) \circ \mathbb{P}_{t_k,t_{k_n}} g(\mathbf{y}) \right| \\ & \leq C \|g_1\|_{\text{Lip}} \eta_k^{3/2} e^{-\theta(t_{k_n}-t_k)} \left(1 + \frac{1}{N} \right) \sqrt{\mathcal{V}(\mathbf{y}) + d}. \end{aligned}$$

Since $t_{k_n} > t_n - 1$, we obtain the following estimate of \mathcal{J}_1 :

$$(4.5) \quad |\mathcal{J}_1| \leq C \|h\|_{\infty} \sqrt{\mathcal{V}(\mathbf{y}) + d} \left(1 + \frac{1}{N} \right) \sum_{k=1}^{k_n-1} \eta_k^{3/2} e^{-\theta(t_n-t_k)}.$$

For the term \mathcal{J}_2 , we claim

$$(4.6) \quad |\mathcal{J}_2| \leq C \|h\|_{\infty} d^{7/2} \sqrt{\mathcal{V}(\mathbf{y}) + d} \sum_{k=k_n}^n \left(\eta_k^{3/2} + \frac{\eta_k}{\sqrt{N}} \right) e^{-\theta(t_n-t_k)}.$$

Combining (4.4), (4.5) and (4.6) implies

$$(4.7) \quad \left| \mathbb{E}h(\mathbf{Y}_{t_n}^{\mathbf{y}}) - \mathbb{E}h(\tilde{\mathbf{Y}}_{t_n}^{\mathbf{y}}) \right| \leq C \|h\|_{\infty} d^{7/2} \sqrt{\mathcal{V}(\mathbf{y}) + d} \sum_{k=1}^n \left(\eta_k^{3/2} + \frac{\eta_k}{\sqrt{N}} \right) e^{-\theta(t_n-t_k)}.$$

Then, applying Lemma A.3 yields the desired result.

It remains to prove (4.6). For any fixed $k \in \{k_n, \dots, n\}$ in \mathcal{J}_2 , we can find time t_{ℓ} such that

$$(4.8) \quad \frac{1}{6(B_1 + \gamma + 2)} \leq t_k - t_{\ell} \leq \frac{1}{5(B_1 + \gamma + 2)}.$$

Then we have the following decomposition,

$$\begin{aligned} & \mathbb{Q}_{0,t_{k-1}} \circ \left(\mathbb{P}_{t_{k-1},t_k} - \mathbb{Q}_{t_{k-1},t_k} \right) \circ \mathbb{P}_{t_k,t_n} h(\mathbf{y}) \\ & = \mathbb{Q}_{0,t_{\ell}} \circ \mathbb{Q}_{t_{\ell},t_{k-1}} \circ \left(\mathbb{P}_{t_{k-1},t_k} - \mathbb{Q}_{t_{k-1},t_k} \right) g_2(\mathbf{y}) \\ & = \mathbb{E} \left\{ \left[g_2 \left(\mathbf{Y}_{t_{k-1},t_k} \left(\tilde{\mathbf{Y}}_{t_{\ell},t_{k-1}}^{\mathbf{z}} \right) \right) - g_2 \left(\tilde{\mathbf{Y}}_{t_{\ell},t_k}^{\mathbf{z}} \right) \right] \left(\mathbf{1}_{E_{\ell,k}^1 \cup E_{m,k}^1 \cup E_{\ell,k}^2} + \mathbf{1}_{(E_{\ell,k}^1 \cup E_{m,k}^1 \cup E_{\ell,k}^2)^c} \right) \right\} \\ & = \mathbb{E} \left\{ \left[g_2 \left(\mathbf{Y}_{t_{k-1},t_k} \left(\tilde{\mathbf{Y}}_{t_{\ell},t_{k-1}}^{\mathbf{z}} \right) \right) - g_2 \left(\tilde{\mathbf{Y}}_{t_{\ell},t_k}^{\mathbf{z}} \right) \right] \mathbf{1}_{E_{\ell,k}^1 \cup E_{m,k}^1 \cup E_{\ell,k}^2} \right\} \\ & \quad + \mathbb{E} \left\{ \mathbb{E} \left[g_2 \left(\mathbf{Y}_{t_{k-1},t_k} \left(\tilde{\mathbf{Y}}_{t_{\ell},t_{k-1}}^{\mathbf{z}} \right) \right) - g_2 \left(\tilde{\mathbf{Y}}_{t_{\ell},t_k}^{\mathbf{z}} \right) \middle| \mathcal{F}_{\ell,k} \right] \mathbf{1}_{(E_{\ell,k}^1 \cup E_{m,k}^1 \cup E_{\ell,k}^2)^c} \right\} \\ & =: I_{k,1} + I_{k,2}, \end{aligned}$$

where $g_2(\mathbf{y}) = \mathbb{P}_{t_k,t_n} h(\mathbf{y})$, $\mathbf{z} = \tilde{\mathbf{Y}}_{t_{\ell}}^{\mathbf{y}}$, $\mathcal{F}_{\ell,k} = \sigma(\tilde{\mathbf{Y}}_{t_{\ell}}^{\mathbf{y}}, (\xi_i)_{\ell+1 \leq i \leq k})$, and $E_{\ell,k}^1$, $E_{m,k}^1$, $E_{\ell,k}^2$ are defined in (3.13) with

$$m = \min \left\{ j: t_k - t_j \leq \frac{t_k - t_{\ell}}{10(B_1 + \gamma + 2)} \right\}.$$

Let us estimate $I_{k,1}$ and $I_{k,2}$.

For $I_{k,1}$, since $C_1 \leq t_k - t_m \leq t_k - t_{\ell} \leq C_2$ holds for some constants $C_1, C_2 > 0$, by Lemma 3.9, we have

$$(4.9) \quad \begin{aligned} |I_{k,1}| & \leq 2 \|g_2\|_{\infty} \mathbb{P}(E_{\ell,k}^1 \cup E_{m,k}^1 \cup E_{\ell,k}^2) \\ & \leq 2 \|h\|_{\infty} (\mathbb{P}(E_{\ell,k}^1) + \mathbb{P}(E_{m,k}^1) + \mathbb{P}(E_{\ell,k}^2)) \leq C \|h\|_{\infty} \eta_k^2. \end{aligned}$$

For $I_{k,2}$, observe that

$$(4.10) \quad I_{k,2} = \int_0^1 \mathbb{E} \left\{ \mathbb{E} \left[\left\langle \nabla g_2(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} + r\boldsymbol{\Xi}_{t_\ell, t_k}), \boldsymbol{\Xi}_{t_\ell, t_k} \right\rangle \middle| \mathcal{F}_{\ell, k} \right] \mathbf{1}_{(E_{\ell, k}^1 \cup E_{m, k}^1 \cup E_{\ell, k}^2)^c} \right\} dr,$$

where $\boldsymbol{\Xi}_{t_\ell, t_k} = \mathbf{Y}_{t_{k-1}, t_k}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - \tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}}$. Denote $\mathbf{G} = \boldsymbol{\Xi}_{t_\ell, t_k}$ and $\mathbf{F} = \tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} + r\boldsymbol{\Xi}_{t_\ell, t_k}$, and we shall apply Lemma 3.8 to estimate $\mathbb{E} \left[\left\langle \nabla g_2(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} + r\boldsymbol{\Xi}_{t_\ell, t_k}), \boldsymbol{\Xi}_{t_\ell, t_k} \right\rangle \middle| \mathcal{F}_{\ell, k} \right]$ on the event $(E_{\ell, k}^1 \cup E_{m, k}^1 \cup E_{\ell, k}^2)^c$. To this end, let us consider SDE (1.6) on $[t_\ell, t_k]$. Note that $\mathcal{F}_{\ell, k}$ is independent of $(\mathbf{B}_t)_{t \in [t_\ell, t_k]}$ and that the Malliavin calculus in Lemma 3.8 is associated to $(\mathbf{B}_t)_{t \in [t_\ell, t_k]}$ and has nothing to do with $\mathcal{F}_{\ell, k}$. An advantage is that we can easily bound $\mathcal{K}_1, \mathcal{K}_2, \mathcal{K}_3$ on the right hand of (3.11) on the event $(E_{\ell, k}^1 \cup E_{m, k}^1 \cup E_{\ell, k}^2)^c$.

According to (3.9) in Lemma 3.5 and Lemma 3.12(i), $\|\mathbf{G}\|_{1,4}$ satisfies

$$\|\mathbf{G}\|_{1,4} = (\mathbb{E} |\mathbf{G}|^4 + \mathbb{E} \|\mathbf{D}\mathbf{G}\|_{\mathcal{H}}^4)^{1/4} \leq C \sqrt{d(\mathcal{V}(\mathbf{z}) + d)} \left(\eta_k^{3/2} + \frac{\eta_k}{\sqrt{N}} \right).$$

Next we estimate $\mathbb{E}[\mathcal{K}_1 \mathcal{K}_2 \mathcal{K}_3]$. By Lemmas 3.10(ii) and 3.12(ii),

$$\mathcal{K}_1 = 1 + \|\mathbf{D}\mathbf{F}\|_{\mathcal{H}}^8 \leq Cd^4 \quad \text{on the event } (E_{\ell, k}^1 \cup E_{\ell, k}^2)^c.$$

By Lemma 3.11,

$$\mathbb{E}\mathcal{K}_2 = 1 + \mathbb{E}\|\mathbf{D}^2\mathbf{F}\|_{\mathcal{H} \otimes \mathcal{H}}^4 \leq Cd^4.$$

Let us now bound $\mathcal{K}_3 = 1 + \|\boldsymbol{\Gamma}(\mathbf{F})^{-1}\|_{\text{HS}}^8$. Since $\|\boldsymbol{\Gamma}(\mathbf{F})^{-1}\|_{\text{HS}} \leq \sqrt{2d} [\lambda_{\min}(\boldsymbol{\Gamma}(\mathbf{F}))]^{-1}$, we first show

$$\lambda_{\min}(\boldsymbol{\Gamma}(\mathbf{F})) \geq C > 0 \quad \text{on the event } (E_{\ell, k}^1 \cup E_{m, k}^1 \cup E_{\ell, k}^2)^c.$$

Indeed, observe $\boldsymbol{\Gamma}(\mathbf{F}) = \boldsymbol{\Gamma}(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}}) + r^2\boldsymbol{\Gamma}(\boldsymbol{\Xi}_{t_\ell, t_k}) + r\mathbf{N}$, where the entries of matrix \mathbf{N} are given by

$$N_{ij} = \left\langle \mathbf{D}(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}})^i, \mathbf{D}(\boldsymbol{\Xi}_{t_\ell, t_k})^j \right\rangle_{\mathcal{H}} + \left\langle \mathbf{D}(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}})^j, \mathbf{D}(\boldsymbol{\Xi}_{t_\ell, t_k})^i \right\rangle_{\mathcal{H}}, \quad i, j = 1, \dots, 2d.$$

By Lemmas 3.10(ii) and 3.12(ii), on the event $(E_{\ell, k}^1 \cup E_{\ell, k}^2)^c$, we have $\|\mathbf{D}\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}}\|_{\mathcal{H}} \leq C\sqrt{d}$ and $\|\mathbf{D}\boldsymbol{\Xi}_{t_\ell, t_k}\|_{\mathcal{H}} \leq C\sqrt{d\eta_k}$. Thus

$$\begin{aligned} \|\mathbf{N}\|_{\text{HS}} &\leq 2 \left\{ \sum_{i,j=1}^{2d} \left\| \mathbf{D}(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}})^i \right\|_{\mathcal{H}}^2 \left\| \mathbf{D}(\boldsymbol{\Xi}_{t_\ell, t_k})^j \right\|_{\mathcal{H}}^2 \right\}^{1/2} \\ &\leq 2 \left\| \mathbf{D}\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} \right\|_{\mathcal{H}} \left\| \mathbf{D}\boldsymbol{\Xi}_{t_\ell, t_k} \right\|_{\mathcal{H}} \leq Cd\sqrt{\eta_k}, \quad \text{on the event } (E_{\ell, k}^1 \cup E_{\ell, k}^2)^c. \end{aligned}$$

This, together with Lemma 3.13 and (4.8), implies

$$\begin{aligned} \lambda_{\min}(\boldsymbol{\Gamma}(\mathbf{F})) &\geq \lambda_{\min}(\boldsymbol{\Gamma}(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}})) + r^2\lambda_{\min}(\boldsymbol{\Gamma}(\boldsymbol{\Xi}_{t_\ell, t_k})) - r\|\mathbf{N}\|_{\text{HS}} \\ &\geq C + 0 - Cd\sqrt{\eta_k} \\ &> C/2, \quad \text{on the event } (E_{\ell, k}^1 \cup E_{m, k}^1 \cup E_{\ell, k}^2)^c, \end{aligned}$$

for any $r \in [0, 1]$ and η_k sufficiently small such that $\eta_k \leq cd^{-2}$ for some positive constant c . Then, we have

$$\mathcal{K}_3 = 1 + \|\boldsymbol{\Gamma}(\mathbf{F})^{-1}\|_{\text{HS}}^8 \leq Cd^4, \quad \text{on the event } (E_{\ell, k}^1 \cup E_{m, k}^1 \cup E_{\ell, k}^2)^c.$$

Hence, the following holds on event $(E_{\ell,k}^1 \cup E_{m,k}^1 \cup E_{\ell,k}^2)^c$

$$(4.11) \quad \mathbb{E} \left[\left\langle \nabla g_2(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} + r \Xi_{t_\ell, t_k}), \Xi_{t_\ell, t_k} \right\rangle \middle| \mathcal{F}_{\ell, k} \right] \leq C \|g_2\|_\infty d^{7/2} \sqrt{\mathcal{V}(\mathbf{z}) + d} \left(\eta_k^{3/2} + \frac{\eta_k}{\sqrt{N}} \right).$$

Recalling $\mathbf{z} = \tilde{\mathbf{Y}}_{t_\ell}^{\mathbf{y}}$, Lemma 3.4 yields that

$$\mathbb{E} \left[\sqrt{\mathcal{V}(\tilde{\mathbf{Y}}_{t_\ell}^{\mathbf{y}}) + d} \right] \leq \sqrt{\mathbb{E} \left[\mathcal{V}(\tilde{\mathbf{Y}}_{t_\ell}^{\mathbf{y}}) \right] + d} \leq C \sqrt{\mathcal{V}(\mathbf{y}) + d}.$$

This, together with (4.10) and (4.11), implies

$$(4.12) \quad |I_{k,2}| \leq C \|g_2\|_\infty d^{7/2} \sqrt{\mathcal{V}(\mathbf{y}) + d} \left(\eta_k^{3/2} + \frac{\eta_k}{\sqrt{N}} \right).$$

It follows from (4.9), (4.12) and $\|g_2\|_\infty \leq \|h\|_\infty$ that

$$\begin{aligned} |\mathcal{J}_2| &\leq C \|h\|_\infty d^{7/2} \sqrt{\mathcal{V}(\mathbf{y}) + d} \sum_{k=k_n}^n \left(\eta_k^{3/2} + \frac{\eta_k}{\sqrt{N}} \right) \\ &\leq C e^\theta \|h\|_\infty d^{7/2} \sqrt{\mathcal{V}(\mathbf{y}) + d} \sum_{k=k_n}^n \left(\eta_k^{3/2} + \frac{\eta_k}{\sqrt{N}} \right) e^{-\theta(t_n - t_k)}, \end{aligned}$$

where in the second inequality we use $t_n - t_k < 1$ such that $e^{-\theta(t_n - t_k)} \geq e^{-\theta}$ for each $k_n \leq k \leq n$. Hence, (4.6) is verified.

We complete the proof. \square

4.3. Proofs of Corollaries 3 and 4.

Proof of Corollary 3. Actually, by (4.3), we only need to analyze

$$\sum_{k=1}^n \eta_k^{3/2} e^{-\theta(t_n - t_k)}, \quad \eta_k = \frac{\eta}{k^\alpha}, \quad \alpha \in (0, 1).$$

Observe that

$$\eta_{k-1} - \eta_k = \eta \cdot \frac{k^\alpha - (k-1)^\alpha}{k^\alpha(k-1)^\alpha}.$$

Notice that $k^\alpha - (k-1)^\alpha \rightarrow 0$ and $\sqrt{\eta_k} e^{\theta t_k} \rightarrow +\infty$ as $k \rightarrow \infty$, so there exists $n_0 \in \mathbb{N}$ such that for all $k \geq n_0$

$$\eta_k \leq \frac{2\theta - \omega}{2\theta^2}, \quad \eta_{k-1} - \eta_k \leq \omega \eta_k^2 \quad \text{and} \quad \sqrt{\eta_k} e^{\theta t_k} \geq 1.$$

It is obvious that the desired result holds for $n \leq n_0$. And for $n > n_0$, we split the sum term in (4.3) into two parts:

$$(4.13) \quad \begin{aligned} \sum_{k=1}^{n_0} \eta_k^{3/2} e^{-\theta(t_n - t_k)} &\leq C e^{-\theta t_n} \leq C \sqrt{\eta_n} = C \frac{\eta^{1/2}}{n^{\alpha/2}}, \\ \sum_{k=n_0+1}^n \eta_k^{3/2} e^{-\theta(t_n - t_k)} &\leq C \sqrt{\eta_n} = C \frac{\eta^{1/2}}{n^{\alpha/2}}, \end{aligned}$$

where we use Lemma A.3 to estimate the second part. Substituting (4.13) into (4.3) implies the desired result. \square

Proof of Corollary 4. By (4.7), the proof is similar to the proof of Corollary 3. We only need to analyze the following sum additionally,

$$\sum_{k=1}^n \eta_k e^{-\theta(t_n - t_k)}, \quad \eta_k = \frac{\eta}{k^\alpha}, \quad \alpha \in (0, 1).$$

Using the same notations in the proof of Corollary 3, we have

$$(4.14) \quad \sum_{k=1}^n \eta_k e^{-\theta(t_n - t_k)} = \sum_{k=1}^{n_0} \eta_k e^{-\theta(t_n - t_k)} + \sum_{k=n_0+1}^n \eta_k e^{-\theta(t_n - t_k)} \leq C.$$

Substituting (4.13) and (4.14) into (4.7) implies the desired result. \square

4.4. Proof of Corollary 5.

Proof of Corollary 5. Recall $f(\mathbf{x}) = \mathbb{E}F(\mathbf{x}, \xi)$ for any $\mathbf{x} \in \mathbb{R}^d$. As ξ is independent of $(\tilde{\mathbf{X}}_{t_n})_{n \geq 1}$, we have the following decomposition:

$$\begin{aligned} & \mathbb{E}F(\tilde{\mathbf{X}}_{t_n}, \xi) - \min_{\mathbf{x} \in \mathbb{R}^d} \mathbb{E}F(\mathbf{x}, \xi) \\ &= [\mathbb{E}f(\tilde{\mathbf{X}}_{t_n}) - \mathbb{E}f(\mathbf{X}_{t_n})] + [\mathbb{E}f(\mathbf{X}_{t_n}) - \mathbb{E}_\pi f(\mathbf{X})] + [\mathbb{E}_\pi f(\mathbf{X}) - f(\mathbf{x}^*)], \end{aligned}$$

where $\mathbf{x}^* = \arg \min f(\mathbf{x})$. According to the conclusion of exponential ergodicity [MT93], the second term on the right-hand side can be estimated by

$$|\mathbb{E}f(\mathbf{X}_t) - \mathbb{E}_\pi f(\mathbf{X})| \leq B_e (\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0) + 1) e^{-\kappa t},$$

for some positive constants B_e and κ . As for the last term, Proposition 11 in [RRT17] yields that

$$\int_{\mathbb{R}^d} f(\boldsymbol{\omega}) \boldsymbol{\pi}_{\mathbf{x}}(d\boldsymbol{\omega}) - \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \leq \frac{d\beta^2}{4} \log \left[\frac{2eL}{a} \left(\frac{2K}{d\beta^2} + 1 \right) \right]$$

as $\beta^2 \leq a/2$, where $\boldsymbol{\pi}_{\mathbf{x}}$ is the marginal distribution of $\boldsymbol{\pi}$. So it remains to estimate the first term by Theorem 1.

Pick a localization function $\psi \in C^\infty(\mathbb{R}^d, \mathbb{R})$ satisfying $0 \leq \psi(\mathbf{x}) \leq 1$ and

$$\psi(\mathbf{x}) = \begin{cases} 1, & |\mathbf{x}| \leq 1, \\ 0, & |\mathbf{x}| \geq 2. \end{cases}$$

For $n \geq 1$, denote

$$g_n(\mathbf{x}) = f(\mathbf{x}) \psi \left(\frac{\mathbf{x}}{R_n} \right),$$

with some constants $R_n > 0$ determined later, then we have

$$\begin{aligned} |\nabla g_n(\mathbf{x})| &= \left| \psi \left(\frac{\mathbf{x}}{R_n} \right) \nabla f(\mathbf{x}) + \frac{1}{R_n} f(\mathbf{x}) \nabla \psi \left(\frac{\mathbf{x}}{R_n} \right) \right| \\ &\leq \left(|\nabla f(\mathbf{x})| + \frac{\|\nabla \psi\|_\infty}{R_n} |f(\mathbf{x})| \right) \mathbf{1}_{[0, 2R_n]}(|\mathbf{x}|). \end{aligned}$$

Recall that $|f(\mathbf{x})| \leq C|\mathbf{x}|^2$ and $|\nabla f(\mathbf{x})| \leq C|\mathbf{x}|$ for all $\mathbf{x} \in \mathbb{R}^d$ satisfying $|\mathbf{x}| \geq 1$, so $|\nabla g_n(\mathbf{x})| \leq CR_n$ for any $\mathbf{x} \in \mathbb{R}^d$, i.e., $\|g_n\|_{\text{Lip}} \leq CR_n$. It follows from Lemma 3.4 that

$$\left| \mathbb{E}f(\tilde{\mathbf{X}}_{t_n}) - \mathbb{E}g_n(\tilde{\mathbf{X}}_{t_n}) \right| \leq \mathbb{E} \left| f(\tilde{\mathbf{X}}_{t_n}) \mathbf{1}_{(R_n, +\infty)}(|\tilde{\mathbf{X}}_{t_n}|) \right|$$

$$\begin{aligned}
&\leq C\mathbb{E}\left[|\tilde{\mathbf{X}}_{t_n}|^2\mathbf{1}_{(R_n,+\infty)}(|\tilde{\mathbf{X}}_{t_n}|)\right] \\
&\leq CR_n^{2-2p}\mathbb{E}\left[|\tilde{\mathbf{X}}_{t_n}|^{2p}\mathbf{1}_{(R_n,+\infty)}(|\tilde{\mathbf{X}}_{t_n}|)\right] \\
&\leq CR_n^{2-2p}(\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0)^p + d^p).
\end{aligned}$$

Similarly, the following holds by Lemma 3.1

$$|\mathbb{E}f(\mathbf{X}_{t_n}) - \mathbb{E}g_n(\mathbf{X}_{t_n})| \leq CR_n^{2-2p}(\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0)^p + d^p).$$

Thus, by Theorem 1 we have

$$\begin{aligned}
&\left|\mathbb{E}f(\tilde{\mathbf{X}}_{t_n}) - \mathbb{E}f(\mathbf{X}_{t_n})\right| \\
&\leq \left|\mathbb{E}f(\tilde{\mathbf{X}}_{t_n}) - \mathbb{E}g_n(\tilde{\mathbf{X}}_{t_n})\right| + \left|\mathbb{E}g_n(\tilde{\mathbf{X}}_{t_n}) - \mathbb{E}g_n(\mathbf{X}_{t_n})\right| + \left|\mathbb{E}g_n(\mathbf{X}_{t_n}) - \mathbb{E}f(\mathbf{X}_{t_n})\right| \\
&\leq CR_n^{2-2p}(\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0)^p + d^p) + C\|g_n\|_{\text{Lip}}\left(1 + \frac{1}{N}\right)\sqrt{\eta_n(\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0) + d)} \\
&\leq CR_n^{2-2p}(\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0)^p + d^p) + CR_n\left(1 + \frac{1}{N}\right)\sqrt{\eta_n(\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0) + d)}.
\end{aligned}$$

Since the above equation holds with arbitrary $R_n > 0$, we can take

$$R_n = \eta_n^{1/(2-4p)}\sqrt{\mathcal{V}(\mathbf{m}_0, \mathbf{x}_0) + d}.$$

Let $2p = q$, we obtain the desired. □

REFERENCES

- [AAMN24] Dallas Albritton, Scott Armstrong, Jean-Christophe Mourrat, and Matthew Novack. Variational methods for the kinetic Fokker-Planck equation. *Anal. PDE*, 17(6):1953–2010, 2024.
- [BGM10] François Bolley, Arnaud Guillin, and Florent Malrieu. Trend to equilibrium and particle approximation for a weakly selfconsistent Vlasov-Fokker-Planck equation. *M2AN Math. Model. Numer. Anal.*, 44(5):867–884, 2010.
- [BLB17] Aleksandar Botev, Guy Lever, and David Barber. Nesterov’s accelerated gradient and momentum as approximations to regularised update descent. In *2017 International joint conference on neural networks (IJCNN)*, pages 1899–1903. IEEE, 2017.
- [CCAY⁺18] Xiang Cheng, Niladri S. Chatterji, Yasin Abbasi-Yadkori, Peter L. Bartlett, and Michael I. Jordan. Sharp convergence rates for langevin dynamics in the nonconvex setting, 2018.
- [CCBJ18] Xiang Cheng, Niladri S. Chatterji, Peter L. Bartlett, and Michael I. Jordan. Underdamped langevin mcmc: A non-asymptotic analysis. In Sbastien Bubeck, Vianney Perchet, and Philippe Rigollet, editors, *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pages 300–323. PMLR, 06–09 Jul 2018.
- [CDC15] Changyou Chen, Nan Ding, and Lawrence Carin. On the convergence of stochastic gradient MCMC algorithms with high-order integrators. *Advances in neural information processing systems*, 28, 2015.
- [CHS87] Tzue-Shuh Chiang, Chii-Ruey Hwang, and Shuenn Jyi Sheu. Diffusion for global optimization in \mathbb{R}^n . *SIAM Journal on Control and Optimization*, 25(3):737–753, 1987.
- [CL20] Jerry Chee and Ping Li. Understanding and detecting convergence for stochastic gradient descent with momentum. In *2020 IEEE International Conference on Big Data (Big Data)*, pages 133–140. IEEE, 2020.
- [CLTZ20] Xi Chen, Jason D. Lee, Xin T. Tong, and Yichen Zhang. Statistical inference for model parameters in stochastic gradient descent. *Ann. Statist.*, 48(1):251–273, 2020.
- [CLW23] Yu Cao, Jianfeng Lu, and Lihan Wang. On explicit L^2 -convergence rate estimate for underdamped Langevin dynamics. *Arch. Ration. Mech. Anal.*, 247(5):Paper No. 90, 34, 2023.
- [CR22] Huy N Chau and Miklós Rásonyi. Stochastic gradient Hamiltonian Monte Carlo for non-convex learning. *Stochastic Processes and their Applications*, 149:341–368, 2022.
- [CSX23] Peng Chen, Qi-Man Shao, and Lihu Xu. A probability approximation framework: Markov process approach. *The Annals of Applied Probability*, 33(2):1419–1459, 2023.
- [DP14] Zhao Dong and Xuhui Peng. Malliavin matrix of degenerate SDE and gradient estimate. *Electron. J. Probab.*, 19:no. 73, 26, 2014.
- [Ebe11] Andreas Eberle. Reflection coupling and Wasserstein contractivity without convexity. *Comptes Rendus Mathématique*, 349(19-20):1101–1104, 2011.
- [Ebe16] Andreas Eberle. Reflection couplings and contraction rates for diffusions. *Probab. Theory Related Fields*, 166(3-4):851–886, 2016.
- [EGZ19] Andreas Eberle, Arnaud Guillin, and Raphael Zimmer. Couplings and quantitative contraction rates for Langevin dynamics. *The Annals of Probability*, 47(4):1982 – 2010, 2019.
- [FDBD21] Xavier Fontaine, Valentin De Bortoli, and Alain Durmus. Convergence rates and approximation results for SGD and its continuous-time counterpart. In *Conference on Learning Theory*, pages 1965–2058. PMLR, 2021.
- [GGZ22] Xuefeng Gao, Mert Gürbüzbalaban, and Lingjiong Zhu. Global convergence of stochastic gradient Hamiltonian Monte Carlo for nonconvex stochastic optimization: Nonasymptotic performance bounds and momentum-based acceleration. *Operations Research*, 70(5):2931–2947, 2022.
- [GM91] Saul B Gelfand and Sanjoy K Mitter. Recursive stochastic algorithms for global optimization in \mathbb{R}^d . *SIAM Journal on Control and Optimization*, 29(5):999–1018, 1991.

- [GM16] Arnaud Guillin and Pierre Monmarché. Optimal linear drift for the speed of convergence of an hypoelliptic diffusion. *Electron. Commun. Probab.*, 21:Paper No. 74, 14, 2016.
- [KF16] Donghwan Kim and Jeffrey A Fessler. Optimized first-order methods for smooth convex minimization. *Mathematical programming*, 159:81–107, 2016.
- [KLRT15] Jakub Konečný, Jie Liu, Peter Richtárik, and Martin Takáč. Mini-batch semi-stochastic gradient descent in the proximal setting. *IEEE Journal of Selected Topics in Signal Processing*, 10(2):242–255, 2015.
- [LR20] Nicolas Loizou and Peter Richtárik. Momentum and stochastic momentum for stochastic gradient, newton, proximal point and subspace descent methods. *Computational Optimization and Applications*, 77(3):653–710, 2020.
- [LTW17] Qianxiao Li, Cheng Tai, and E Weinan. Stochastic modified equations and adaptive stochastic gradient algorithms. In *International Conference on Machine Learning*, pages 2101–2110. PMLR, 2017.
- [LTW19] Qianxiao Li, Cheng Tai, and E Weinan. Stochastic modified equations and dynamics of stochastic gradient algorithms I: Mathematical foundations. *The Journal of Machine Learning Research*, 20(1):1474–1520, 2019.
- [LZCS14] Mu Li, Tong Zhang, Yuqiang Chen, and Alexander J Smola. Efficient mini-batch training for stochastic optimization. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 661–670, 2014.
- [MJ19] Michael Muehlebach and Michael Jordan. A dynamical systems perspective on nesterov acceleration. In *International Conference on Machine Learning*, pages 4656–4662. PMLR, 2019.
- [MT93] Sean P. Meyn and R. L. Tweedie. Stability of Markovian processes. III. Foster-Lyapunov criteria for continuous-time processes. *Adv. in Appl. Probab.*, 25(3):518–548, 1993.
- [Nea11] Radford M. Neal. MCMC using Hamiltonian dynamics. In *Handbook of Markov chain Monte Carlo*, Chapman & Hall/CRC Handb. Mod. Stat. Methods, pages 113–162. CRC Press, Boca Raton, FL, 2011.
- [Nor86] James Norris. Simplified Malliavin calculus. In *Séminaire de Probabilités, XX, 1984/85*, volume 1204 of *Lecture Notes in Math.*, pages 101–130. Springer, Berlin, 1986.
- [Nua06] David Nualart. *The Malliavin calculus and related topics*. Probability and its Applications (New York). Springer-Verlag, Berlin, second edition, 2006.
- [Pav14] Grigorios A. Pavliotis. *Stochastic processes and applications*, volume 60 of *Texts in Applied Mathematics*. Springer, New York, 2014. Diffusion processes, the Fokker-Planck and Langevin equations.
- [PP23] Gilles Pages and Fabien Panloup. Unadjusted Langevin algorithm with multiplicative noise: Total variation and Wasserstein bounds. *The Annals of Applied Probability*, 33(1):726–779, 2023.
- [RRT17] Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. Non-convex learning via stochastic gradient Langevin dynamics: a nonasymptotic analysis. In Satyen Kale and Ohad Shamir, editors, *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pages 1674–1703. PMLR, 07–10 Jul 2017.
- [Sch24] Katharina Schuh. Global contractivity for Langevin dynamics with distribution-dependent forces and uniform in time propagation of chaos. *Ann. Inst. Henri Poincaré Probab. Stat.*, 60(2):753–789, 2024.
- [SW24] Katharina Schuh and Peter A. Whalley. Convergence of kinetic langevin samplers for non-convex potentials, 2024.
- [Wu01] Liming Wu. Large and moderate deviations and exponential convergence for stochastic damping Hamiltonian systems. *Stochastic processes and their applications*, 91(2):205–238, 2001.
- [Zha10] Xicheng Zhang. Stochastic flows and Bismut formulas for stochastic Hamiltonian systems. *Stochastic Process. Appl.*, 120(10):1929–1949, 2010.

APPENDIX

APPENDIX A. SUPPORTING LEMMAS

The first lemma shows that $f(\mathbf{x})$ admits lower and upper bounds that are quadratic functions.

Lemma A.1. [RRT17, Lemma 2] *If function $f(\mathbf{x})$ satisfies Assumption 2.1, then the following quadratic bounds hold,*

$$\frac{a}{2} |\mathbf{x}|^2 - \frac{K}{2} \log 3 \leq f(\mathbf{x}) \leq L |\mathbf{x}|^2 + \frac{B^2}{2L} + A, \quad \forall \mathbf{x} \in \mathbb{R}^d,$$

where $K = b + B^2/(2a)$.

Applying Assumption I (iii), it is easy to verify that

$$(A.1) \quad \langle \mathbf{x}, \nabla f(\mathbf{x}) \rangle \geq \frac{a}{2} |\mathbf{x}|^2 - K \geq 2\lambda \left(f(\mathbf{x}) + \frac{1}{4} \gamma^2 |\mathbf{x}|^2 \right) - \mathring{A}, \quad \forall \mathbf{x} \in \mathbb{R}^d,$$

where the constants λ and \mathring{A} satisfy

$$(A.2) \quad 0 < \lambda \leq \min \left\{ \frac{1}{4}, \frac{a}{4L + \gamma^2} \right\}, \quad \mathring{A} \geq K + 2\lambda \left(\frac{B^2}{2L} + A \right).$$

Applying (A.1), the next lemma shows that the Lyapunov function $\mathcal{V}(\mathbf{m}, \mathbf{x})$ defined in (2.6) satisfies (3.3).

Lemma A.2. *Under Assumption 2.1, and λ and \mathring{A} are in (A.2), we have that*

$$\mathscr{A}\mathcal{V}(\mathbf{m}, \mathbf{x}) \leq -\lambda\gamma\mathcal{V}(\mathbf{m}, \mathbf{x}) + (\gamma\mathring{A} + d\beta^2)/2.$$

Proof. Recall that

$$\mathcal{V}(\mathbf{m}, \mathbf{x}) := f(\mathbf{x}) + \frac{\gamma^2}{4} \left(\left| \mathbf{x} + \frac{1}{\gamma} \mathbf{m} \right|^2 + \left| \frac{1}{\gamma} \mathbf{m} \right|^2 - \lambda |\mathbf{x}|^2 \right),$$

and we have

$$\begin{aligned} \nabla_{\mathbf{m}} \mathcal{V}(\mathbf{m}, \mathbf{x}) &= \mathbf{m} + \frac{\gamma}{2} \mathbf{x}, & \Delta_{\mathbf{m}} \mathcal{V}(\mathbf{m}, \mathbf{x}) &= d, \\ \nabla_{\mathbf{x}} \mathcal{V}(\mathbf{m}, \mathbf{x}) &= \nabla f(\mathbf{x}) + \frac{\gamma^2(1-\lambda)}{2} \mathbf{x} + \frac{\gamma}{2} \mathbf{m}. \end{aligned}$$

Hence, by inequality (A.1), we have

$$\begin{aligned} \mathscr{A}\mathcal{V}(\mathbf{m}, \mathbf{x}) &= -\langle \nabla_{\mathbf{m}} \mathcal{V}(\mathbf{m}, \mathbf{x}), \gamma \mathbf{m} + \nabla f(\mathbf{x}) \rangle + \langle \nabla_{\mathbf{x}} \mathcal{V}(\mathbf{m}, \mathbf{x}), \mathbf{m} \rangle + \frac{1}{2} \beta^2 \Delta_{\mathbf{m}} \mathcal{V}(\mathbf{m}, \mathbf{x}) \\ &\leq -\left[\lambda\gamma \left(f(\mathbf{x}) + \frac{\gamma^2}{4} |\mathbf{x}|^2 \right) \right] + \frac{\gamma}{2} \mathring{A} - \frac{\lambda\gamma^2}{2} \langle \mathbf{m}, \mathbf{x} \rangle - \frac{\gamma}{2} |\mathbf{m}|^2 + \frac{1}{2} d\beta^2 \\ &= -\lambda\gamma\mathcal{V}(\mathbf{m}, \mathbf{x}) - \frac{\lambda^2\gamma^3}{4} |\mathbf{x}|^2 - \frac{(1-\lambda)\gamma}{2} |\mathbf{m}|^2 + \frac{1}{2} (\gamma\mathring{A} + d\beta^2), \end{aligned}$$

where the inequality comes from (A.1). Since $\lambda \leq 1/4$, we complete the proof. \square

The process $(\mathbf{M}_t, \mathbf{X}_t)_{t \geq 0}$ admits a unique stationary distribution as

$$\pi(d\mathbf{m}, d\mathbf{x}) = \frac{1}{Z} \exp \left\{ -\frac{\gamma}{\beta^2} (|\mathbf{m}|^2 + 2f(\mathbf{x})) \right\} d\mathbf{m} d\mathbf{x},$$

where Z is a normalized constant given by

$$Z = \left(\frac{\beta^2 \pi}{\gamma} \right)^{d/2} \int_{\mathbb{R}^d} e^{-2\gamma f(\mathbf{x})/\beta^2} d\mathbf{x},$$

see e.g. [Pav14]. Thanks to [MT93], Lemma A.2 yields that process $(\mathbf{M}_t, \mathbf{X}_t)_{t \geq 0}$ is exponentially ergodic. That is, there exist some positive constants B_e and κ , such that for all $\mathbf{m}, \mathbf{x} \in \mathbb{R}^d$,

$$\sup_{g \leq \mathcal{V}+1} \left| \mathbb{E}[g(\mathbf{M}_t^{\mathbf{m}}, \mathbf{X}_t^{\mathbf{x}})] - \mathbb{E}_{\pi}[g(\mathbf{M}, \mathbf{X})] \right| \leq B_e(\mathcal{V}(\mathbf{m}, \mathbf{x}) + 1)e^{-\kappa t}, \quad \forall t \geq 0,$$

where (\mathbf{M}, \mathbf{X}) has the distribution π . By the definition of $\mathcal{V}(\mathbf{m}, \mathbf{x})$, this immediately implies that for all $\mathbf{x} \in \mathbb{R}^d$,

$$(A.3) \quad |\mathbb{E}f(\mathbf{X}_t^{\mathbf{x}}) - \mathbb{E}_{\pi}[f(\mathbf{X})]| \leq B_e(\mathcal{V}(\mathbf{m}, \mathbf{x}) + 1)e^{-\kappa t}, \quad \forall t \geq 0.$$

Lemma A.3. *Let step size $(\eta_k)_{k \geq 1}$ satisfy Assumption III and $t_n = \sum_{i=1}^n \eta_i$. For any $\varepsilon \in [0, 1/2]$, the following holds*

$$\sum_{k=1}^n \eta_k^{1+\varepsilon} e^{-\theta(t_n - t_k)} \leq \frac{4}{2\theta - (4\varepsilon - 1)\omega} \eta_n^{\varepsilon}.$$

Proof. Notice that for each $k \leq n$

$$(A.4) \quad \eta_k^{1+\varepsilon} e^{-\theta(t_n - t_k)} = e^{-\theta t_n} \cdot \eta_k^{1+\varepsilon} \frac{e^{\theta t_k} - e^{\theta t_{k-1}}}{1 - e^{-\theta \eta_k}} \leq \frac{4e^{-\theta t_n}}{2\theta + \omega} \cdot \eta_k^{\varepsilon} (e^{\theta t_k} - e^{\theta t_{k-1}}),$$

where in the last inequality, we use the fact

$$e^{-x} \leq 1 - x + \frac{x^2}{2} \leq 1 - x + \frac{2\theta - \omega}{4\theta} x = 1 - \frac{2\theta + \omega}{4\theta} x, \quad \forall 0 \leq x \leq 1 - \frac{\omega}{2\theta},$$

to obtain that

$$1 - e^{-\theta \eta_k} \geq \frac{2\theta + \omega}{4} \eta_k, \quad \forall \eta_k \leq \frac{2\theta - \omega}{2\theta^2}.$$

Besides, according to condition $\eta_{k-1} - \eta_k \leq \omega \eta_k^2$ and Bernoulli's inequality, it holds that

$$\eta_{k-1}^{\varepsilon} - \eta_k^{\varepsilon} = \eta_k^{\varepsilon} \left[\left[\frac{\eta_{k-1}}{\eta_k} \right]^{\varepsilon} - 1 \right] \leq \varepsilon \eta_k^{\varepsilon} \left[\frac{\eta_{k-1}}{\eta_k} - 1 \right] \leq \varepsilon \omega \eta_k^{1+\varepsilon},$$

which derives that

$$(A.5) \quad \begin{aligned} \sum_{k=1}^n \eta_k^{\varepsilon} (e^{\theta t_k} - e^{\theta t_{k-1}}) &= \sum_{k=1}^n (\eta_k^{\varepsilon} e^{\theta t_k} - \eta_{k-1}^{\varepsilon} e^{\theta t_{k-1}}) + \sum_{k=1}^n (\eta_{k-1}^{\varepsilon} - \eta_k^{\varepsilon}) e^{\theta t_{k-1}} \\ &\leq \eta_n^{\varepsilon} e^{\theta t_n} + \varepsilon \omega \sum_{k=1}^n \eta_k^{1+\varepsilon} e^{\theta t_k}. \end{aligned}$$

Combining (A.4) and (A.5) shows that

$$\sum_{k=1}^n \eta_k^{1+\varepsilon} e^{-\theta(t_n - t_k)} \leq \frac{4}{2\theta + \omega} \eta_n^{\varepsilon} + \frac{4\varepsilon\omega}{2\theta + \omega} \sum_{k=1}^n \eta_k^{1+\varepsilon} e^{-\theta(t_n - t_k)},$$

and a simply calculation yields the desired. \square

APPENDIX B. PROOFS OF AUXILIARY LEMMAS

Lemma B.1. *Let $\mathbf{U}^i, i = 1, \dots, N$ be i.i.d. random vectors in \mathbb{R}^d satisfying $\mathbb{E}\mathbf{U}^1 = \mathbf{0}$ and $\mathbb{E}[|\mathbf{U}^1|^{2p}] < +\infty$ for some $p \in \mathbb{N}$. Then there exists a constant $C > 0$ only depending on p , such that*

$$\mathbb{E} \left[\left| \frac{1}{N} \sum_{i=1}^N \mathbf{U}^i \right|^{2p} \right] \leq \frac{C}{N^p} \mathbb{E} [|\mathbf{U}^1|^{2p}].$$

Proof. Notice that

$$\left| \sum_{i=1}^N \mathbf{U}^i \right|^{2p} = \prod_{j=1}^p \left\langle \sum_{i_{2j-1}=1}^N \mathbf{U}^{i_{2j-1}}, \sum_{i_{2j}=1}^N \mathbf{U}^{i_{2j}} \right\rangle = \sum_{i_1, \dots, i_{2p}=1}^N \prod_{j=1}^p \langle \mathbf{U}^{i_{2j-1}}, \mathbf{U}^{i_{2j}} \rangle,$$

which implies that

$$\mathbb{E} \left[\left| \frac{1}{N} \sum_{i=1}^N \mathbf{U}^i \right|^{2p} \right] = \frac{1}{N^{2p}} \sum_{i_1, \dots, i_{2p}=1}^N \mathbb{E} \left[\prod_{j=1}^p \langle \mathbf{U}^{i_{2j-1}}, \mathbf{U}^{i_{2j}} \rangle \right].$$

Since $\mathbf{U}^1, \dots, \mathbf{U}^N$ are independent and $\mathbb{E}\mathbf{U}^1 = \mathbf{0}$, we have $\mathbb{E}[\prod_{j=1}^p \langle \mathbf{U}^{i_{2j-1}}, \mathbf{U}^{i_{2j}} \rangle] = 0$ whenever some $i \in \{1, \dots, N\}$ appears only once among i_1, \dots, i_{2p} . For the rest terms, Hölder's inequality shows that

$$\mathbb{E} \left[\prod_{j=1}^p \langle \mathbf{U}^{i_{2j-1}}, \mathbf{U}^{i_{2j}} \rangle \right] \leq \mathbb{E} \left[\prod_{j=1}^{2p} |\mathbf{U}^{i_j}| \right] \leq \prod_{j=1}^{2p} \left\{ \mathbb{E} [|\mathbf{U}^{i_j}|^{2p}] \right\}^{1/2p} = \mathbb{E} [|\mathbf{U}^1|^{2p}].$$

Furthermore, for $l = 1, \dots, p$, the index set $\{1, \dots, 2p\}$ can be decomposed into l parts such that each part consists of at least two indices, and c_l denotes the number of such partitions. Then we have

$$\sum_{i_1, \dots, i_{2p}=1}^N \mathbb{E} \left[\prod_{j=1}^p \langle \mathbf{U}^{i_{2j-1}}, \mathbf{U}^{i_{2j}} \rangle \right] \leq \sum_{l=1}^p c_l N^l \mathbb{E} [|\mathbf{U}^1|^{2p}].$$

The desired result follows from

$$\mathbb{E} \left[\left| \frac{1}{N} \sum_{i=1}^N \mathbf{U}^i \right|^{2p} \right] \leq \frac{1}{N^{2p}} \sum_{l=1}^p c_l N^l \mathbb{E} [|\mathbf{U}^1|^{2p}] \leq \frac{C}{N^p} \mathbb{E} [|\mathbf{U}^1|^{2p}].$$

□

B.1. Proofs of Lemmas for Moments Estimations.

Proof of Lemma 3.1. Recalling the Lyapunov function $\mathcal{V}(\mathbf{m}, \mathbf{x})$ in (2.6), we define function F_t as

$$F_t := \mathbb{E} \mathcal{V}(\mathbf{M}_t^{\mathbf{m}}, \mathbf{X}_t^{\mathbf{x}}) = \mathbb{E} \left[f(\mathbf{X}_t^{\mathbf{x}}) + \frac{\gamma^2}{4} \left(\left| \mathbf{X}_t^{\mathbf{x}} + \frac{1}{\gamma} \mathbf{M}_t^{\mathbf{m}} \right|^2 + \left| \frac{1}{\gamma} \mathbf{M}_t^{\mathbf{m}} \right|^2 - \lambda |\mathbf{X}_t^{\mathbf{x}}|^2 \right) \right].$$

Itô's formula implies that F_t satisfies the following differential inequality

$$\frac{dF_t}{dt} \leq -\lambda \gamma F_t + \frac{1}{2} (\gamma \dot{A} + d\beta^2),$$

with initial condition $F_0 = \mathcal{V}(\mathbf{m}, \mathbf{x})$. Hence, we can obtain that

$$(B.1) \quad F_t \leq e^{-\lambda\gamma t} F_0 + \frac{1}{2\lambda\gamma} \left(\gamma \dot{A} + d\beta^2 \right),$$

which yields the result with respect to of $p = 1$ immediately. As for the case of $p > 1$, we can apply similar argument. By Itô's formula, it holds

$$\begin{aligned} \mathcal{A}(\mathcal{V}(\mathbf{m}, \mathbf{x})^p) &\leq p\mathcal{V}(\mathbf{m}, \mathbf{x})^{p-1} \mathcal{A}\mathcal{V}(\mathbf{m}, \mathbf{x}) + \frac{p(p-1)\beta^2}{2} \mathcal{V}(\mathbf{m}, \mathbf{x})^{p-2} \left| \mathbf{m} + \frac{\gamma}{2} \mathbf{x} \right|^2 \\ &\leq -p\lambda\gamma \mathcal{V}(\mathbf{m}, \mathbf{x})^p + C\mathcal{V}(\mathbf{m}, \mathbf{x})^{p-1}, \end{aligned}$$

where the last inequality is because of (3.2) and (3.3). Then, by Young's inequality, the function $G_t := \mathbb{E}[\mathcal{V}(\mathbf{M}_t^{\mathbf{m}}, \mathbf{X}_t^{\mathbf{x}})^p]$ satisfies

$$\frac{dG_t}{dt} \leq -\lambda\gamma G_t + C,$$

with initial value $G_0 = \mathcal{V}(\mathbf{m}, \mathbf{x})^p$, where $C > 0$ here also depends on p . This implies the desired immediately. \square

Proof of Lemma 3.2. Under Assumption I, we have

$$\begin{aligned} \mathbb{E} |\mathbf{M}_t^{\mathbf{m}} - \mathbf{m}|^{2p} &= \mathbb{E} \left| \int_0^t [-\gamma \mathbf{M}_s - \nabla f(\mathbf{X}_s)] ds + \beta \mathbf{B}_t \right|^{2p} \\ &\leq Ct^{2p-1} \mathbb{E} \int_0^t |\gamma \mathbf{M}_s + \nabla f(\mathbf{X}_s)|^{2p} ds + C\mathbb{E} |\mathbf{B}_t|^{2p} \\ &\leq Ct^{2p-1} \int_0^t (\gamma^{2p} \mathbb{E} |\mathbf{M}_s|^{2p} + L^{2p} \mathbb{E} |\mathbf{X}_s|^{2p} + B^{2p}) ds + Ct^p d^p, \\ &\leq Ct^{2p} (\mathcal{V}(\mathbf{m}, \mathbf{x})^p + d^p) + Ct^p d^p, \end{aligned}$$

where the first inequality is by the Hölder's inequality, the second inequality is by (3.1) and the last one is due to (3.5). Similarly, it holds

$$\mathbb{E} |\mathbf{X}_t^{\mathbf{x}} - \mathbf{x}|^{2p} = \mathbb{E} \left| \int_0^t \mathbf{M}_s ds \right|^{2p} \leq t^{2p-1} \int_0^t \mathbb{E} |\mathbf{M}_s|^{2p} ds \leq Ct^{2p} (\mathcal{V}(\mathbf{m}, \mathbf{x})^p + d^p).$$

Combining above two inequalities, we can obtain the desired. \square

Proof of Lemma 3.3. Recall SDE (3.4) and rewrite it as

$$\begin{cases} \widetilde{\mathbf{M}}_t^{\mathbf{m}} - \mathbf{m} = -\gamma t \mathbf{m} - t \nabla f(\mathbf{x}) + t \left(\nabla f(\mathbf{x}) - \frac{1}{N} \sum_{i=1}^N \nabla F(\mathbf{x}, \xi^i) \right) + \beta \mathbf{B}_t, \\ \widetilde{\mathbf{X}}_t^{\mathbf{x}} - \mathbf{x} = t \mathbf{m}. \end{cases}$$

for all $t \in [0, \eta]$ with η being a step size. To make notations simple, denote

$$(B.2) \quad \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) = \nabla f(\mathbf{x}) - \nabla F(\mathbf{x}, \xi^i), \quad i = 1, \dots, N,$$

which are i.i.d. random vectors. Then, Assumption II and Lemma B.1 imply that

$$(B.3) \quad \mathbb{E} \left| \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) \right|^{2p} \leq \frac{C}{N^p} \mathbb{E} |\boldsymbol{\Lambda}(\mathbf{x}, \xi^i)|^{2p} \leq \frac{C}{N^p}, \quad \forall 1 \leq p \leq \frac{q}{2}.$$

If the condition (2.7) in Assumption IV holds additionally, (B.3) holds for all $1 \leq p \leq q'/2$. Thus, we have that

$$\begin{aligned} & \mathbb{E} \left| \widetilde{\mathbf{M}}_t^{\mathbf{m}} - \mathbf{m} \right|^{2p} \\ & \leq C \left[t^{2p} \mathbb{E} |\gamma \mathbf{m} + \nabla f(\mathbf{x})|^{2p} + t^{2p} \mathbb{E} |N^{-1} \sum_{i=1}^d \Lambda(\mathbf{x}, \xi^i)|^{2p} + \beta^{2p} \mathbb{E} |\mathbf{B}_t|^{2p} \right] \\ & \leq C \left[|\mathbf{m}|^{2p} + |\mathbf{x}|^{2p} + 1 \right] t^{2p} + Ct^p d^p. \end{aligned}$$

Using (3.2), we can obtain

$$\mathbb{E} \left| \widetilde{\mathbf{M}}_t^{\mathbf{m}} - \mathbf{m} \right|^{2p} \leq Ct^p (\mathcal{V}(\mathbf{m}, \mathbf{x})^p + d^p).$$

Similarly, we have

$$\mathbb{E} \left| \widetilde{\mathbf{X}}_t^{\mathbf{x}} - \mathbf{x} \right|^{2p} = t^{2p} \mathbb{E} |\mathbf{m}|^{2p} \leq Ct^{2p} \mathcal{V}(\mathbf{m}, \mathbf{x})^p.$$

Hence, by $0 \leq t \leq \eta < 1$, combining all above yields the desired. \square

Proof of Lemma 3.4. Consider SDE (1.6), let us prove the case of $p = 1$ firstly. By Itô's formula, $(\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t))_{0 \leq t \leq \eta}$ satisfies

$$\begin{aligned} \text{(B.4)} \quad d\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t) &= -\frac{\gamma}{2} \langle \widetilde{\mathbf{M}}_t, \mathbf{m} \rangle dt - \frac{\lambda\gamma^2}{2} \langle \widetilde{\mathbf{X}}_t, \mathbf{m} \rangle dt + \langle \nabla f(\widetilde{\mathbf{X}}_t), \mathbf{m} \rangle dt + \frac{d\beta^2}{2} dt \\ &\quad - \left\langle \widetilde{\mathbf{M}}_t + \frac{\gamma}{2} \widetilde{\mathbf{X}}_t, \frac{1}{N} \sum_{i=1}^N \nabla F(\mathbf{x}, \xi^i) \right\rangle dt + \beta \left\langle \widetilde{\mathbf{M}}_t + \frac{\gamma}{2} \widetilde{\mathbf{X}}_t, d\mathbf{B}_t \right\rangle \\ &=: [I_1 + I_2] dt + \beta \left\langle \widetilde{\mathbf{M}}_t + \frac{\gamma}{2} \widetilde{\mathbf{X}}_t, d\mathbf{B}_t \right\rangle, \end{aligned}$$

where I_1 and I_2 are given by

$$\begin{aligned} I_1 &= -\frac{\gamma}{2} |\widetilde{\mathbf{M}}_t|^2 - \frac{\lambda\gamma^2}{2} \langle \widetilde{\mathbf{X}}_t, \widetilde{\mathbf{M}}_t \rangle - \frac{\gamma}{2} \langle \widetilde{\mathbf{X}}_t, \nabla f(\widetilde{\mathbf{X}}_t) \rangle + \frac{d\beta^2}{2}, \\ I_2 &= \frac{\gamma}{2} \left\langle \widetilde{\mathbf{M}}_t + \lambda\gamma \widetilde{\mathbf{X}}_t, \widetilde{\mathbf{M}}_t - \mathbf{m} \right\rangle - \left\langle \nabla f(\widetilde{\mathbf{X}}_t), \widetilde{\mathbf{M}}_t - \mathbf{m} \right\rangle \\ &\quad + \left\langle \widetilde{\mathbf{M}}_t + \frac{\gamma}{2} \widetilde{\mathbf{X}}_t, \nabla f(\widetilde{\mathbf{X}}_t) - \nabla f(\mathbf{x}) \right\rangle + \left\langle \widetilde{\mathbf{M}}_t + \frac{\gamma}{2} \widetilde{\mathbf{X}}_t, \frac{1}{N} \sum_{i=1}^N \Lambda(\mathbf{x}, \xi^i) \right\rangle. \end{aligned}$$

Here, $\Lambda(\mathbf{x}, \xi^i)$, $i = 1, \dots, N$ are defined by (B.2).

For I_1 , by the definition of \mathcal{V} in (2.6) and (A.1), we have

$$\begin{aligned} \text{(B.5)} \quad I_1 &\leq - \left[\lambda\gamma \left(f(\widetilde{\mathbf{X}}_t) + \frac{\gamma^2}{4} |\widetilde{\mathbf{X}}_t|^2 \right) \right] - \frac{\lambda\gamma^2}{2} \langle \widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t \rangle - \frac{\gamma}{2} |\widetilde{\mathbf{M}}_t|^2 + \frac{1}{2} (\gamma \dot{A} + d\beta^2) \\ &\leq -\lambda\gamma \mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t) + \frac{1}{2} (\gamma \dot{A} + d\beta^2). \end{aligned}$$

For I_2 , the Young's inequality yields that for any $\varepsilon > 0$

$$\begin{aligned} I_2 &\leq \varepsilon \left[\frac{\gamma}{2} \left| \widetilde{\mathbf{M}}_t + \lambda\gamma \widetilde{\mathbf{X}}_t \right|^2 + 2 \left(L^2 \left| \widetilde{\mathbf{X}}_t \right|^2 + B^2 \right) + 2 \left| \widetilde{\mathbf{M}}_t + \frac{\gamma}{2} \widetilde{\mathbf{X}}_t \right|^2 \right] \\ &\quad + \frac{1}{4\varepsilon} \left[\frac{\gamma+2}{2} \left| \widetilde{\mathbf{M}}_t - \mathbf{m} \right|^2 + L^2 \left| \widetilde{\mathbf{X}}_t - \mathbf{x} \right|^2 \right] + \frac{1}{4\varepsilon} \left| \frac{1}{N} \sum_{i=1}^N \Lambda(\mathbf{x}, \xi^i) \right|^2. \end{aligned}$$

Thus, by choosing a special ε and (3.2), we can obtain

$$(B.6) \quad I_2 \leq \frac{\lambda\gamma}{3} \mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t) + C \left[\left| \widetilde{\mathbf{M}}_t - \mathbf{m} \right|^2 + \left| \widetilde{\mathbf{X}}_t - \mathbf{x} \right|^2 \right] + C \left| \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) \right|^2 + C.$$

Combining (B.4), (B.5) and (B.6), we have

$$(B.7) \quad \begin{aligned} d\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t) &\leq -\frac{2\lambda\gamma}{3} \mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t) dt + \beta \left\langle \widetilde{\mathbf{M}}_t + \frac{\gamma}{2} \widetilde{\mathbf{X}}_t, d\mathbf{B}_t \right\rangle \\ &\quad + C \left\{ \left[\left| \widetilde{\mathbf{M}}_t - \mathbf{m} \right|^2 + \left| \widetilde{\mathbf{X}}_t - \mathbf{x} \right|^2 \right] + \left| \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) \right|^2 + d \right\} dt \end{aligned}$$

Then, by Lemmas 3.3 and B.1, we have the following differential inequality

$$\frac{d}{dt} \mathbb{E} \mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t) \leq -\frac{2\lambda\gamma}{3} \mathbb{E} \mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t) + Ct \mathcal{V}(\mathbf{m}, \mathbf{x}) + Cd,$$

for all $0 \leq t \leq \eta < 1$. Then, solving above inequality, we have

$$\mathbb{E} \mathcal{V}(\widetilde{\mathbf{M}}_\eta, \widetilde{\mathbf{X}}_\eta) \leq e^{-2\lambda\gamma\eta/3} \mathcal{V}(\mathbf{m}, \mathbf{x}) + C\eta^2 \mathcal{V}(\mathbf{m}, \mathbf{x}) + Cd\eta.$$

Furthermore, the fact $1 - x \leq e^{-x} \leq 1 - x + x^2/2$ for all $x \in \mathbb{R}$ yields that one can find a positive constant η_0 such that the following inequality holds for all $\eta \leq \eta_0$,

$$(B.8) \quad \mathbb{E} \mathcal{V}(\widetilde{\mathbf{M}}_\eta, \widetilde{\mathbf{X}}_\eta) \leq e^{-\gamma\lambda\eta/2} \mathcal{V}(\mathbf{m}, \mathbf{x}) + Cd\eta.$$

By Markov property, (B.8) implies the following recursive inequality,

$$\mathbb{E} \mathcal{V}(\widetilde{\mathbf{M}}_{t_k}, \widetilde{\mathbf{X}}_{t_k}) \leq e^{-\lambda\gamma\eta_k/2} \mathbb{E} \mathcal{V}(\widetilde{\mathbf{M}}_{t_{k-1}}, \widetilde{\mathbf{X}}_{t_{k-1}}) + Cd\eta_k,$$

as long as $\eta_k \leq \eta_0$. Then, (3.7) with $p = 1$ follows from that

$$(B.9) \quad \begin{aligned} &\mathbb{E} \mathcal{V}(\widetilde{\mathbf{M}}_{t_n}, \widetilde{\mathbf{X}}_{t_n}) - e^{-\lambda\gamma t_n/2} \mathcal{V}(\mathbf{m}, \mathbf{x}) \\ &= \sum_{k=1}^n e^{-\lambda\gamma(t_n - t_k)/2} \left[\mathbb{E} \mathcal{V}(\widetilde{\mathbf{M}}_{t_k}, \widetilde{\mathbf{X}}_{t_k}) - e^{-\lambda\gamma\eta_k/2} \mathbb{E} \mathcal{V}(\widetilde{\mathbf{M}}_{t_{k-1}}, \widetilde{\mathbf{X}}_{t_{k-1}}) \right] \\ &\leq Cd \sum_{k=1}^n \eta_k e^{-\lambda\gamma(t_n - t_k)/2} \\ &= Cd \sum_{k=1}^n \frac{\eta_k}{1 - e^{-\lambda\gamma\eta_k/2}} \left[e^{-\lambda\gamma(t_n - t_k)/2} - e^{-\lambda\gamma(t_n - t_{k-1})/2} \right] \leq Cd, \end{aligned}$$

where the last inequality holds as $\eta_k \leq 2/(\gamma\lambda)$.

As for $p > 1$, by Itô's formula and (B.4), we have that for $0 \leq t \leq \eta$,

$$\begin{aligned} d\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^p &= p\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^{p-1} [I_1 + I_2] dt + p\beta \mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^{p-1} \left\langle \widetilde{\mathbf{M}}_t + \frac{\gamma}{2} \widetilde{\mathbf{X}}_t, d\mathbf{B}_t \right\rangle \\ &\quad + \frac{p(p-1)\beta^2}{2} \left| \widetilde{\mathbf{M}}_t + \frac{\gamma}{2} \widetilde{\mathbf{X}}_t \right|^2 \mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^{p-2} dt. \end{aligned}$$

This, together with (B.5), (B.6) and $|\mathbf{m} + \gamma\mathbf{x}/2|^2 \leq C\mathcal{V}(\mathbf{m}, \mathbf{x})$ by (3.2), it holds that

$$(B.10) \quad \begin{aligned} d\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^p &\leq \left[p\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^{p-1} \left(-\frac{2\lambda\gamma}{3}\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t) + C\mathcal{R}_0 \right) \right] dt \\ &\quad + p\beta\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^{p-1} \left\langle \widetilde{\mathbf{M}}_t + \frac{\gamma}{2}\widetilde{\mathbf{X}}_t, d\mathbf{B}_t \right\rangle, \end{aligned}$$

where the reserved item \mathcal{R}_0 is of the form

$$\mathcal{R}_0 = \left| \widetilde{\mathbf{M}}_t - \mathbf{m} \right|^2 + \left| \widetilde{\mathbf{X}}_t - \mathbf{x} \right|^2 + \left| \frac{1}{N} \sum_{i=1}^N \Lambda(\mathbf{x}, \xi^i) \right|^2 + d,$$

which satisfies the following by Lemmas 3.3 and B.1:

$$(B.11) \quad \mathbb{E}[\mathcal{R}_0^p] \leq Ct^p\mathcal{V}(\mathbf{m}, \mathbf{x})^p + Cd^p.$$

By Young's inequality again, we can obtain that

$$\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^{p-1} \cdot (C\mathcal{R}_0) \leq \left(1 - \frac{1}{p}\right) \frac{2\lambda\gamma}{3}\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^p + C\mathcal{R}_0^p.$$

Substituting this into (B.10), there is

$$d\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^p \leq -\frac{2\lambda\gamma}{3}\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^p dt + C\mathcal{R}_0^p dt + p\beta\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^{p-1} \left\langle \widetilde{\mathbf{M}}_t + \frac{\gamma}{2}\widetilde{\mathbf{X}}_t, d\mathbf{B}_t \right\rangle.$$

Combining (B.11), this implies the follow differential inequality

$$\frac{d}{dt}\mathbb{E}\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^p \leq -\frac{2\lambda\gamma}{3}\mathbb{E}\mathcal{V}(\widetilde{\mathbf{M}}_t, \widetilde{\mathbf{X}}_t)^p + Ct^p\mathcal{V}(\mathbf{m}, \mathbf{x})^p + Cd^p.$$

By Gronwall's inequality, we have

$$\mathbb{E}\mathcal{V}(\widetilde{\mathbf{M}}_\eta, \widetilde{\mathbf{X}}_\eta)^p \leq e^{-2\lambda\gamma\eta/3}\mathcal{V}(\mathbf{m}, \mathbf{x})^p + C\eta^p\mathcal{V}(\mathbf{m}, \mathbf{x})^p + Cd^p\eta.$$

And one can also find a positive constant η'_0 such that for all $\eta \leq \eta'_0$, the following holds,

$$\mathbb{E}\mathcal{V}(\widetilde{\mathbf{M}}_\eta, \widetilde{\mathbf{X}}_\eta)^p \leq e^{-\lambda\gamma\eta/2}\mathcal{V}(\mathbf{m}, \mathbf{x})^p + Cd^p\eta.$$

Hence, for any $k \geq 1$, as $\eta_k \leq \eta'_0$, by Markov property, we have

$$\mathbb{E}\mathcal{V}(\widetilde{\mathbf{M}}_{t_k}, \widetilde{\mathbf{X}}_{t_k})^p \leq e^{-\lambda\gamma\eta_k/2}\mathbb{E}\mathcal{V}(\widetilde{\mathbf{M}}_{t_{k-1}}, \widetilde{\mathbf{X}}_{t_{k-1}})^p + Cd^p\eta_k.$$

which implies (3.7) for $1 < p \leq q/2$ under Assumption II (or $1 < p \leq q'/2$ if the condition (2.7) in Assumption IV holds) by recursion similar to (B.9).

The proof is complete. \square

B.2. Proofs of Auxiliary Lemmas for 1-Wasserstein Distance.

Proof of Lemma 3.5. (i) Suppose that $(\widehat{\mathbf{M}}_t, \widehat{\mathbf{X}}_t)_{0 \leq t \leq \eta}$ is the solution of SDE

$$\begin{cases} d\widehat{\mathbf{M}}_t = -\gamma\mathbf{m} dt - \nabla f(\mathbf{x}) dt + \beta d\mathbf{B}_t, \\ d\widehat{\mathbf{X}}_t = \mathbf{m} dt, \end{cases}$$

with initial value $(\widehat{\mathbf{M}}_0, \widehat{\mathbf{X}}_0) = (\mathbf{m}, \mathbf{x})$. We claim that

$$(B.12) \quad d_{\mathcal{W}_1}(\mathcal{L}(\mathbf{M}_\eta, \mathbf{X}_\eta), \mathcal{L}(\widehat{\mathbf{M}}_\eta, \widehat{\mathbf{X}}_\eta)) \leq C\eta^{3/2}\sqrt{\mathcal{V}(\mathbf{m}, \mathbf{x}) + d},$$

$$(B.13) \quad d_{\mathcal{W}_1}(\mathcal{L}(\widehat{\mathbf{M}}_\eta, \widehat{\mathbf{X}}_\eta), \mathcal{L}(\widetilde{\mathbf{M}}_\eta, \widetilde{\mathbf{X}}_\eta)) \leq C\frac{\sqrt{d}}{N}\eta^{3/2},$$

which immediately implies the desired result by the triangle inequality.

Let us now show (B.12). Observe that $(\mathbf{M}_t - \widehat{\mathbf{M}}_t, \mathbf{X}_t - \widehat{\mathbf{X}}_t)_{0 \leq t \leq \eta}$ satisfies

$$\begin{cases} d(\mathbf{M}_t - \widehat{\mathbf{M}}_t) = -\gamma(\mathbf{M}_t - \mathbf{m}) dt - (\nabla f(\mathbf{X}_t) - \nabla f(\mathbf{x})) dt, \\ d(\mathbf{X}_t - \widehat{\mathbf{X}}_t) = (\mathbf{M}_t - \mathbf{m}) dt, \end{cases}$$

with initial value $(\mathbf{M}_t - \widehat{\mathbf{M}}_t, \mathbf{X}_t - \widehat{\mathbf{X}}_t) = (\mathbf{0}, \mathbf{0})$. Then, it holds that

$$\begin{aligned} \mathbb{E} \left| \mathbf{M}_t - \widehat{\mathbf{M}}_t \right|^2 &\leq 2\gamma^2 \mathbb{E} \left| \int_0^t (\mathbf{M}_s - \mathbf{m}) ds \right|^2 + 2\mathbb{E} \left| \int_0^t (\nabla f(\mathbf{X}_s) - \nabla f(\mathbf{x})) ds \right|^2 \\ &\leq 2\gamma^2 t \mathbb{E} \int_0^t |\mathbf{M}_s - \mathbf{m}|^2 ds + 2t \mathbb{E} \int_0^t |\nabla f(\mathbf{X}_s) - \nabla f(\mathbf{x})|^2 ds, \\ \mathbb{E} \left| \mathbf{X}_t - \widehat{\mathbf{X}}_t \right|^2 &= \mathbb{E} \left| \int_0^t (\mathbf{M}_s - \mathbf{m}) ds \right|^2 \leq t \int_0^t \mathbb{E} |\mathbf{M}_s - \mathbf{m}|^2 ds. \end{aligned}$$

These, together with Lemma 3.2, we have

$$\mathbb{E} \left| \mathbf{M}_\eta - \widehat{\mathbf{M}}_\eta \right|^2 + \mathbb{E} \left| \mathbf{X}_\eta - \widehat{\mathbf{X}}_\eta \right|^2 \leq C\eta^3(\mathcal{V}(\mathbf{m}, \mathbf{x}) + d).$$

By definition of 1-Wasserstein distance, we immediately know (B.12) holds true.

It remains to prove (B.13). Due to $\widehat{\mathbf{X}}_\eta = \widetilde{\mathbf{X}}_\eta = \mathbf{x} + \eta\mathbf{m}$, it is equivalent to prove

$$d_{W_1}(\mathcal{L}(\widehat{\mathbf{M}}_\eta), \mathcal{L}(\widetilde{\mathbf{M}}_\eta)) \leq C \frac{\sqrt{d}}{N} \eta^{3/2},$$

or by Kantorovich-Rubinstein Theorem,

$$\left| \mathbb{E}h(\widehat{\mathbf{M}}_\eta) - \mathbb{E}h(\widetilde{\mathbf{M}}_\eta) \right| \leq C \frac{\sqrt{d}}{N} \eta^{3/2} \|h\|_{\text{Lip}},$$

for any Lipschitz function h . Since h can be pointwise approximated by a sequence of $h_i \in \mathcal{C}_b^2(\mathbb{R}^d, \mathbb{R})$, $i = 1, 2, \dots$, satisfying $\|\nabla h_i\|_\infty \leq 2\|h\|_{\text{Lip}}$, we just need to show that

$$(B.14) \quad \left| \mathbb{E}h(\widehat{\mathbf{M}}_\eta) - \mathbb{E}h(\widetilde{\mathbf{M}}_\eta) \right| \leq C \frac{\sqrt{d}}{N} \eta^{3/2} \|\nabla h\|_\infty, \quad \forall h \in \mathcal{C}_b^2(\mathbb{R}^d).$$

For any fixed $h \in \mathcal{C}_b^2(\mathbb{R}^d)$, we define

$$h_0(\mathbf{z}) := \mathbb{E} [h(\mathbf{z} + \beta\mathbf{B}_\eta)] = \int_{\mathbb{R}^d} h(\mathbf{z} + \mathbf{y}) \cdot (2\pi\beta^2\eta)^{-d/2} \exp\{-|\mathbf{y}|^2/(2\beta^2\eta)\} d\mathbf{y}.$$

Straightforward calculations derive that

$$\begin{aligned} \frac{\partial^2 h_0}{\partial z_i \partial z_j}(\mathbf{z}) &= \int_{\mathbb{R}^d} \frac{\partial^2 h}{\partial z_i \partial z_j}(\mathbf{z} + \mathbf{y}) \cdot (2\pi\beta^2\eta)^{-d/2} \exp\{-|\mathbf{y}|^2/(2\beta^2\eta)\} d\mathbf{y} \\ &= \int_{\mathbb{R}^d} \frac{\partial h}{\partial z_i}(\mathbf{z} + \mathbf{y}) \cdot \frac{y_j}{\beta^2\eta} \cdot (2\pi\beta^2\eta)^{-d/2} \exp\{-|\mathbf{y}|^2/(2\beta^2\eta)\} d\mathbf{y} \\ &= \frac{1}{\beta^2\eta} \mathbb{E} \left[\frac{\partial h}{\partial z_i}(\mathbf{z} + \beta\mathbf{B}_\eta) \cdot \beta\mathbf{B}_\eta^{(j)} \right], \end{aligned}$$

where $\mathbf{B}_\eta^{(j)}$ is the j -th component of \mathbf{B}_η . That is

$$\nabla^2 h_0(\mathbf{z}) = \frac{1}{\beta\eta} \mathbb{E} [\nabla h(\mathbf{z} + \beta\mathbf{B}_\eta) \mathbf{B}_\eta^\top],$$

which implies that

$$(B.15) \quad \|\nabla^2 h_0\|_\infty \leq \frac{\|\nabla h\|_\infty \mathbb{E} |\mathbf{B}_\eta|}{\beta \eta} \leq \frac{\|\nabla h\|_\infty}{\beta} \sqrt{\frac{d}{\eta}}.$$

Taylor's expansion gives that

$$\begin{aligned} & \mathbb{E} h(\widetilde{\mathbf{M}}_\eta) - \mathbb{E} h(\widehat{\mathbf{M}}_\eta) \\ &= \mathbb{E} \left[h_0 \left((1 - \eta\gamma) \mathbf{m} - \eta \frac{1}{N} \sum_{i=1}^N \nabla F(\mathbf{x}, \xi^i) \right) \right] - \mathbb{E} [h_0((1 - \eta\gamma) \mathbf{m} - \eta \nabla f(\mathbf{x}))] \\ &= \eta \mathbb{E} \left[\left\langle \nabla h_0((1 - \eta\gamma) \mathbf{m} - \eta \nabla f(\mathbf{x})), \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) \right\rangle \right] \\ & \quad + \frac{1}{2} \eta^2 \mathbb{E} \left[\left\langle \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i), \nabla^2 h_0(\boldsymbol{\zeta}) \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) \right\rangle \right] \\ &= \frac{1}{2} \eta^2 \mathbb{E} \left[\left\langle \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i), \nabla^2 h_0(\boldsymbol{\zeta}) \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) \right\rangle \right], \end{aligned}$$

holds for some random vector $\boldsymbol{\zeta}$, where $\boldsymbol{\Lambda}(\mathbf{x}, \xi^i)$, $i = 1, \dots, N$ are defined in (B.2) and we use the independence between \mathbf{B}_η and ξ . Together with (B.15) and Lemma B.1, we have

$$\left| \mathbb{E} h(\widetilde{\mathbf{M}}_\eta) - \mathbb{E} h(\widehat{\mathbf{M}}_\eta) \right| \leq \frac{1}{2} \eta^2 \|\nabla^2 h_0\|_\infty \mathbb{E} \xi \left| \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) \right|^2 \leq C \frac{\sqrt{d}}{N} \eta^{3/2} \|\nabla h\|_\infty,$$

which is (B.14), and leads to (B.13).

(ii) Notice that $(\mathbf{M}_t - \widetilde{\mathbf{M}}_t, \mathbf{X}_t - \widetilde{\mathbf{X}}_t)_{0 \leq t \leq \eta}$ satisfies

$$\begin{cases} d(\mathbf{M}_t - \widetilde{\mathbf{M}}_t) = -\gamma(\mathbf{M}_t - \mathbf{m}) dt - (\nabla f(\mathbf{X}_t) - \nabla f(\mathbf{x})) dt - \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) dt, \\ d(\mathbf{X}_t - \widetilde{\mathbf{X}}_t) = (\mathbf{M}_t - \mathbf{m}) dt, \end{cases}$$

with initial value $(\mathbf{M}_0 - \widetilde{\mathbf{M}}_0, \mathbf{X}_0 - \widetilde{\mathbf{X}}_0) = (\mathbf{0}, \mathbf{0})$. Combining Hölder's inequality and Lemma B.1 derives that

$$\begin{aligned} \mathbb{E} \left| \mathbf{M}_\eta - \widetilde{\mathbf{M}}_\eta \right|^4 &\leq C \mathbb{E} \left| \int_0^\eta (\mathbf{M}_s - \mathbf{m}) ds \right|^4 + C \mathbb{E} \left| \int_0^\eta (\nabla f(\mathbf{X}_s) - \nabla f(\mathbf{x})) ds \right|^4 \\ & \quad + C \eta^4 \mathbb{E} \left| \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Lambda}(\mathbf{x}, \xi^i) \right|^4 \\ &\leq C \eta^3 \int_0^\eta \mathbb{E} |\mathbf{M}_s - \mathbf{m}|^4 ds + C \eta^3 \int_0^\eta \mathbb{E} |\mathbf{X}_s - \mathbf{x}|^4 ds + C \frac{\eta^4}{N^2}, \\ \mathbb{E} \left| \mathbf{X}_\eta - \widetilde{\mathbf{X}}_\eta \right|^4 &= \mathbb{E} \left| \int_0^\eta (\mathbf{M}_s - \mathbf{m}) ds \right|^4 \leq \eta^3 \int_0^\eta \mathbb{E} |\mathbf{M}_s - \mathbf{m}|^4 ds. \end{aligned}$$

Then Lemma 3.2 with $p = 2$ implies

$$\mathbb{E} \left| \mathbf{M}_\eta - \widetilde{\mathbf{M}}_\eta \right|^4 + \mathbb{E} \left| \mathbf{X}_\eta - \widetilde{\mathbf{X}}_\eta \right|^4 \leq C \eta^6 (\mathcal{V}(\mathbf{m}, \mathbf{x})^2 + d^2) + C \frac{\eta^4}{N^2}.$$

We complete the proof. \square

B.3. Proofs of Auxiliary Lemmas in the Total Variation Distance. Before proving Lemma 3.7, we shall introduce the well known Bismut formula in Malliavin calculus. To this end, let us briefly recall the Malliavin calculus by Bismut in our setting. Let $T > 0$ be an arbitrary number, for any $t \in [0, T]$, we take \mathbf{Y}_t as a functional of Brownian motion, i.e., $\mathbf{Y}_t(\mathbf{B})$, let \mathbf{h} be an adaptive stochastic process in $L^2([0, T] \times \Omega, \mathbb{R}^d)$, and define a general Malliavin derivative along \mathbf{h} as the following

$$(B.16) \quad D^{\mathbf{h}}\mathbf{Y}_t(\mathbf{B}) = \lim_{\varepsilon \rightarrow 0} \frac{\mathbf{Y}_t(\mathbf{B} + \varepsilon \mathbf{H}) - \mathbf{Y}_t(\mathbf{B})}{\varepsilon},$$

where $\mathbf{H}(t) = \int_0^t \mathbf{h}(s) ds$ for $t \in [0, T]$, as long as the above limit exists in $L^2(\Omega)$. For $\phi \in C_b^1(\mathbb{R}^d, \mathbb{R})$, by the chain rule,

$$D^{\mathbf{h}}\phi(\mathbf{Y}_t(\mathbf{B})) = \nabla\phi(\mathbf{Y}_t(\mathbf{B})) D^{\mathbf{h}}\mathbf{Y}_t(\mathbf{B}) = \sum_{i=1}^d \partial_i \phi(\mathbf{Y}_t(\mathbf{B})) D^{\mathbf{h}}F_t^i(\mathbf{B})$$

The following Bismut's formula, which can be taken as an integration by parts, is an important property in Malliavin calculus:

$$(B.17) \quad \mathbb{E} [\nabla\phi(\mathbf{Y}_t) D^{\mathbf{h}}\mathbf{Y}_t] = \mathbb{E} \left[\phi(\mathbf{Y}_t) \int_0^t \langle \mathbf{h}(s), d\mathbf{B}_s \rangle \right].$$

Proof of Lemma 3.7. Let $\mathbf{Y}_t^y = (\mathbf{M}_t^{\mathbf{m}}, \mathbf{X}_t^{\mathbf{x}})$ come from SDE (1.5). Let \mathbf{h} be an adaptive stochastic process in $L^2([0, T] \times \Omega, \mathbb{R}^{2d})$ to be chosen later, by [Nor86], the Malliavin derivative of \mathbf{Y}_t along \mathbf{h} satisfies

$$\begin{cases} D^{\mathbf{h}}\mathbf{M}_t^{\mathbf{m}} = -\gamma \int_0^t D^{\mathbf{h}}\mathbf{M}_s^{\mathbf{m}} ds - \int_0^t \nabla^2 f(\mathbf{X}_s^{\mathbf{x}}) D^{\mathbf{h}}\mathbf{X}_s^{\mathbf{x}} ds + \beta \int_0^t \mathbf{h}(s) ds, \\ D^{\mathbf{h}}\mathbf{X}_t^{\mathbf{x}} = \int_0^t D^{\mathbf{h}}\mathbf{M}_s^{\mathbf{m}} ds. \end{cases}$$

Besides, for any $\mathbf{w} = (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^{2d}$, the Jacobian flow of \mathbf{Y}_t with respect to \mathbf{w} is defined as

$$\nabla\mathbf{Y}_t^y \cdot \mathbf{w} = \lim_{\varepsilon \rightarrow 0} \frac{\mathbf{Y}_t^{y+\varepsilon\mathbf{w}} - \mathbf{Y}_t^y}{\varepsilon}.$$

More precisely, we have that

$$\begin{aligned} \nabla_{\mathbf{m}}\mathbf{Y}_t \cdot \mathbf{u} &= \begin{bmatrix} \mathbf{u} \\ \mathbf{0} \end{bmatrix} + \int_0^t \begin{bmatrix} -\gamma\mathbf{I} & \nabla^2 f(\mathbf{X}_t) \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \nabla_{\mathbf{m}}\mathbf{Y}_t \cdot \mathbf{u} dt, \\ \nabla_{\mathbf{x}}\mathbf{Y}_t \cdot \mathbf{v} &= \begin{bmatrix} \mathbf{0} \\ \mathbf{v} \end{bmatrix} + \int_0^t \begin{bmatrix} -\gamma\mathbf{I} & \nabla^2 f(\mathbf{X}_t) \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \nabla_{\mathbf{x}}\mathbf{Y}_t \cdot \mathbf{v} dt. \end{aligned}$$

Observe that for any $\phi \in C_b^1(\mathbb{R}^{2d}, \mathbb{R})$

$$(B.18) \quad \nabla\mathbb{E}\phi(\mathbf{Y}_T) \cdot \mathbf{w} = \mathbb{E} [\nabla\phi(\mathbf{Y}_T) \nabla_{\mathbf{m}}\mathbf{Y}_T \cdot \mathbf{u}] + \mathbb{E} [\nabla\phi(\mathbf{Y}_T) \nabla_{\mathbf{x}}\mathbf{Y}_T \cdot \mathbf{v}].$$

We aim to find some special \mathbf{h} and $\tilde{\mathbf{h}}$ such that

$$\nabla_{\mathbf{m}}\mathbf{Y}_T \cdot \mathbf{u} = D^{\mathbf{h}}\mathbf{Y}_T \quad \text{and} \quad \nabla_{\mathbf{x}}\mathbf{Y}_T \cdot \mathbf{v} = D^{\tilde{\mathbf{h}}}\mathbf{Y}_T,$$

and consequently

$$\nabla\mathbb{E}\phi(\mathbf{Y}_T) \cdot \mathbf{w} = \mathbb{E} [\nabla\phi(\mathbf{Y}_T) D^{\mathbf{h}}\mathbf{Y}_T] + \mathbb{E} [\nabla\phi(\mathbf{Y}_T) D^{\tilde{\mathbf{h}}}\mathbf{Y}_T],$$

which enables us to apply Bismut's formula (B.17).

Define $\mathbf{H}(t) = \int_0^t \mathbf{h}(s) ds$ and $\tilde{\mathbf{H}}(t) = \int_0^t \tilde{\mathbf{h}}(s) ds$, by Zhang [Zha10, Theorem 3.3], for any $T > 0$ and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$, we can choose $\mathbf{H}(t)$ and $\tilde{\mathbf{H}}(t)$ specified as below. Firstly, for any $\mathbf{u} \in \mathbb{R}^d$, let

$$\mathbf{H}(t) = -\frac{1}{\beta} \left(\mathbf{u} - \mathbf{g}(t) - \gamma \int_0^t \mathbf{g}(s) ds - \int_0^t ds \int_0^s \nabla^2 f(\mathbf{X}_s) \mathbf{g}(r) dr \right),$$

for $t \in [0, T]$, where $\mathbf{g}(t)$ is given by

$$\mathbf{g}(t) = \begin{cases} (T - 3t) \mathbf{u}/T, & 0 \leq t < T/2, \\ (t - T) \mathbf{u}/T, & T/2 \leq t \leq T, \end{cases}$$

satisfying $\mathbf{g}(0) = \mathbf{u}$, $\mathbf{g}(T) = \mathbf{0}$ and $\int_0^T \mathbf{g}(t) dt = \mathbf{0}$. One can verify that

$$\begin{cases} \langle \nabla_{\mathbf{m}} \mathbf{M}_t(\mathbf{m}, \mathbf{x}), \mathbf{u} \rangle = D^{\mathbf{h}} \mathbf{M}_t(\mathbf{m}, \mathbf{x}) + \mathbf{g}(t), \\ \langle \nabla_{\mathbf{m}} \mathbf{X}_t(\mathbf{m}, \mathbf{x}), \mathbf{u} \rangle = D^{\mathbf{h}} \mathbf{X}_t(\mathbf{m}, \mathbf{x}) + \int_0^t \mathbf{g}(s) ds, \end{cases}$$

and $\nabla_{\mathbf{m}} \mathbf{Y}_T \cdot \mathbf{u} = D^{\mathbf{h}} \mathbf{Y}_T$. Then, applying Bismut formula to $\mathbb{E} [\nabla \phi(\mathbf{Y}_T) D^{\mathbf{h}} \mathbf{Y}_T]$ yields that

$$\langle \nabla_{\mathbf{m}} \mathbb{E} \phi(\mathbf{Y}_T), \mathbf{u} \rangle = \mathbb{E} \left[-\phi(\mathbf{Y}_T) \frac{1}{\beta} \int_0^T \left(\nabla^2 f(\mathbf{X}_s) \int_0^s \mathbf{g}(r) dr + \gamma \mathbf{g}(s) + \dot{\mathbf{g}}(s) \right) d\mathbf{B}_s \right].$$

Hence, there exists a constant C only depending on L , γ and β such that

$$(B.19) \quad |\langle \nabla_{\mathbf{m}} \mathbb{E} \phi(\mathbf{Y}_T), \mathbf{u} \rangle| \leq C \|\phi\|_{\infty} |\mathbf{u}| (T^{3/2} \vee T^{-1/2}).$$

Analogously, for any $\mathbf{v} \in \mathbb{R}^d$, let $\tilde{\mathbf{H}}(t)$

$$\tilde{\mathbf{H}}(t) = \frac{1}{\beta} \left[-\tilde{\mathbf{g}}(t) - \gamma \int_0^t \tilde{\mathbf{g}}(s) ds - \int_0^t \nabla^2 f(\mathbf{X}_s) \left(\mathbf{v} + \int_0^s \tilde{\mathbf{g}}(r) dr \right) ds \right],$$

for $t \in [0, T]$, where $\tilde{\mathbf{g}}(t)$ is given by

$$\tilde{\mathbf{g}}(t) = \begin{cases} -4t/T^2 \mathbf{v} & , \quad 0 \leq t \leq T/2, \\ 4(t - T)/T^2 \mathbf{v}, & T/2 \leq t \leq T, \end{cases}$$

satisfying $\tilde{\mathbf{g}}(0) = \tilde{\mathbf{g}}(T) = \mathbf{0}$ and $\int_0^T \tilde{\mathbf{g}}(t) dt = -\mathbf{v}$. Then, it is easy to verify that

$$\begin{cases} \langle \nabla_{\mathbf{x}} \mathbf{M}_t(\mathbf{m}, \mathbf{x}), \mathbf{v} \rangle = D^{\tilde{\mathbf{h}}} \mathbf{M}_t(\mathbf{m}, \mathbf{x}) + \tilde{\mathbf{g}}(t), \\ \langle \nabla_{\mathbf{x}} \mathbf{X}_t(\mathbf{m}, \mathbf{x}), \mathbf{v} \rangle = D^{\tilde{\mathbf{h}}} \mathbf{X}_t(\mathbf{m}, \mathbf{x}) + \mathbf{v} + \int_0^t \tilde{\mathbf{g}}(s) ds, \end{cases}$$

and $\nabla_{\mathbf{x}} \mathbf{Y}_T \cdot \mathbf{v} = D^{\tilde{\mathbf{h}}} \mathbf{Y}_T$. Applying Bismut's formula to $\mathbb{E} [\nabla \phi(\mathbf{Y}_T) D^{\tilde{\mathbf{h}}} \mathbf{Y}_T]$ yields

$$\langle \nabla_{\mathbf{x}} \mathbb{E} \phi(\mathbf{Y}_T), \mathbf{v} \rangle = \mathbb{E} \left\{ -\phi(\mathbf{Y}_T) \frac{1}{\beta} \int_0^T \left[\nabla^2 f(\mathbf{X}_s) \left(\mathbf{v} + \int_0^s \tilde{\mathbf{g}}(r) dr \right) + \gamma \tilde{\mathbf{g}}(s) + \dot{\tilde{\mathbf{g}}}(s) \right] d\mathbf{B}_s \right\},$$

which immediately implies that there exists a constant $C > 0$ only depending on L , γ and β such that

$$(B.20) \quad |\langle \nabla_{\mathbf{x}} \mathbb{E} \phi(\mathbf{Y}_T), \mathbf{v} \rangle| \leq C \|\phi\|_{\infty} |\mathbf{v}| (T^{1/2} \vee T^{-3/2}).$$

Combining (B.18), (B.19) and (B.20), we have

$$|\nabla \mathbb{E} \phi(\mathbf{Y}_T) \cdot \mathbf{w}| \leq C \|\phi\|_\infty (|\mathbf{u}| + |\mathbf{v}|) (T^{3/2} \vee T^{-3/2}) \leq C \|\phi\|_\infty (T^{3/2} \vee T^{-3/2}) |\mathbf{w}|$$

for any $\mathbf{w} = (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^{2d}$, and this implies the desired. \square

Proof of Lemma 3.8. Firstly, the chain rule of Malliavin derivative yields that

$$\langle Dg(\mathbf{F}), DF^j \rangle_{\mathcal{H}} = \sum_{i=1}^d \partial_i g(\mathbf{F}) \langle DF^i, DF^j \rangle_{\mathcal{H}} = \sum_{i=1}^d \partial_i g(\mathbf{F}) (\Gamma(\mathbf{F}))_{i,j}, \quad j = 1, \dots, d,$$

which implies that

$$(B.21) \quad \partial_i g(\mathbf{F}) = \sum_{j=1}^d \langle Dg(\mathbf{F}), DF^j \rangle_{\mathcal{H}} (\Gamma(\mathbf{F})^{-1})_{i,j}, \quad i = 1, \dots, d.$$

Then, we have

$$\begin{aligned} \mathbb{E} \langle \nabla g(\mathbf{F}), \mathbf{G} \rangle &= \mathbb{E} \left\langle Dg(\mathbf{F}), \sum_{i=1}^d \sum_{j=1}^d (\Gamma(\mathbf{F}))_{i,j}^{-1} G^i DF^j \right\rangle_{\mathcal{H}} \\ &= \mathbb{E} \left[g(\mathbf{F}) \delta \left(\sum_{j=1}^d (\Gamma(\mathbf{F})^{-1} \mathbf{G})_j DF^j \right) \right]. \end{aligned}$$

where the first equality is by (B.21), and the last one is due to the duality (3.10) between derivative operator and divergence operator. Then, the Jensen's inequality yields that

$$(B.22) \quad |\mathbb{E} \langle \nabla g(\mathbf{F}), \mathbf{G} \rangle| \leq \|g\|_\infty \sqrt{\mathbb{E} \left[\left| \delta \left(\sum_{j=1}^d (\Gamma(\mathbf{F})^{-1} \mathbf{G})_j DF^j \right) \right|^2 \right]}.$$

Applying the continuity of divergence operator δ from $\mathbb{D}^{1,2}(\mathcal{H})$ into $L^2(\Omega)$ [Nua06, Proposition 1.5.7], there exists a constant $C > 0$ such that

$$(B.23) \quad \mathbb{E} \left[\left| \delta \left(\sum_{j=1}^d (\Gamma(\mathbf{F})^{-1} \mathbf{G})_j DF^j \right) \right|^2 \right] \leq C(\mathcal{I}_1 + \mathcal{I}_2),$$

where \mathcal{I}_1 and \mathcal{I}_2 are given by

$$\mathcal{I}_1 = \mathbb{E} \left[\left\| \sum_{j=1}^d (\Gamma(\mathbf{F})^{-1} \mathbf{G})_j DF^j \right\|_{\mathcal{H}}^2 \right], \quad \mathcal{I}_2 = \mathbb{E} \left[\left\| \sum_{j=1}^d D \left((\Gamma(\mathbf{F})^{-1} \mathbf{G})_j DF^j \right) \right\|_{\mathcal{H} \otimes \mathcal{H}}^2 \right]$$

For \mathcal{I}_1 , it follows from Cauchy-Schwarz inequality that

$$\begin{aligned} (B.24) \quad \mathcal{I}_1 &\leq \mathbb{E} \left[\left| \sum_{j=1}^d |(\Gamma(\mathbf{F})^{-1} \mathbf{G})_j| \|DF_j\|_{\mathcal{H}} \right|^2 \right] \\ &\leq \mathbb{E} \left[\left(\sum_{j=1}^d (\Gamma(\mathbf{F})^{-1} \mathbf{G})_j^2 \right) \left(\sum_{j=1}^d \|DF^j\|_{\mathcal{H}}^2 \right) \right] \\ &\leq \mathbb{E} \left\{ \left[\|\Gamma(\mathbf{F})^{-1}\|_{\text{op}} |\mathbf{G}| \|D\mathbf{F}\|_{\mathcal{H}} \right]^2 \right\}. \end{aligned}$$

For \mathcal{I}_2 , according to [Nua06, Lemma 2.1.6], it holds that

$$D(\mathbf{\Gamma}(\mathbf{F})^{-1}\mathbf{G})_j = \sum_{i=1}^d (\mathbf{\Gamma}(\mathbf{F})^{-1})_{j,i} DG^i - \sum_{i,k,l=1}^d G^i (\mathbf{\Gamma}(\mathbf{F})^{-1})_{j,k} (\mathbf{\Gamma}(\mathbf{F})^{-1})_{l,i} D(\mathbf{\Gamma}(\mathbf{F})_{k,l}).$$

This, together with representation $D(\mathbf{\Gamma}(\mathbf{F})_{k,l}) = D\langle DF^k, DF^l \rangle_{\mathcal{H}}$ for any $1 \leq k, l \leq d$, we have

$$(B.25) \quad \left\| \sum_{j=1}^d D\left((\mathbf{\Gamma}(\mathbf{F})^{-1}\mathbf{G})_j DF^j\right) \right\|_{\mathcal{H} \otimes \mathcal{H}} \leq \mathcal{I}_{2,1} + \mathcal{I}_{2,2} + \mathcal{I}_{2,3},$$

where $\mathcal{I}_{2,1}$, $\mathcal{I}_{2,2}$ and $\mathcal{I}_{2,3}$ are given by

$$\begin{aligned} \mathcal{I}_{2,1} &= \sum_{j=1}^d \left| (\mathbf{\Gamma}(\mathbf{F})^{-1}\mathbf{G})_j \right| \|D^2 F^j\|_{\mathcal{H} \otimes \mathcal{H}} \\ \mathcal{I}_{2,2} &= \sum_{i,j=1}^d \left| (\mathbf{\Gamma}(\mathbf{F})^{-1})_{j,i} \right| \|DG^i\|_{\mathcal{H}} \|DF^j\|_{\mathcal{H}} \\ \mathcal{I}_{2,3} &= \sum_{i,j,k,l=1}^d \left| G^i (\mathbf{\Gamma}(\mathbf{F})^{-1})_{j,k} (\mathbf{\Gamma}(\mathbf{F})^{-1})_{l,i} \right| \|D\langle DF^k, DF^l \rangle_{\mathcal{H}}\|_{\mathcal{H}} \|DF^j\|_{\mathcal{H}}. \end{aligned}$$

It is easy to estimate $\mathcal{I}_{2,1}$ and $\mathcal{I}_{2,2}$ which satisfy

$$(B.26) \quad \mathcal{I}_{2,1} \leq \|\mathbf{\Gamma}(\mathbf{F})^{-1}\|_{\text{op}} |\mathbf{G}| \|D^2 \mathbf{F}\|_{\mathcal{H} \otimes \mathcal{H}}, \quad \mathcal{I}_{2,2} \leq \|\mathbf{\Gamma}(\mathbf{F})^{-1}\|_{\text{HS}} \|\mathbf{DG}\|_{\mathcal{H}} \|\mathbf{DF}\|_{\mathcal{H}}.$$

For $\mathcal{I}_{2,3}$, we claim that

$$(B.27) \quad \mathcal{I}_{2,3} \leq 2 \|\mathbf{\Gamma}(\mathbf{F})^{-1}\|_{\text{HS}}^2 |\mathbf{G}| \|\mathbf{DF}\|_{\mathcal{H}}^2 \|D^2 \mathbf{F}\|_{\mathcal{H} \otimes \mathcal{H}}.$$

Combining (B.25), (B.26) and (B.27), we can obtain that

$$(B.28) \quad \begin{aligned} \mathcal{I}_2 &= \mathbb{E} \left[\left\| \sum_{j=1}^d D\left((\mathbf{\Gamma}(\mathbf{F})^{-1}\mathbf{G})_j DF^j\right) \right\|_{\mathcal{H} \otimes \mathcal{H}}^2 \right] \\ &\leq C \left\{ \mathbb{E} \left[\|\mathbf{\Gamma}(\mathbf{F})^{-1}\|_{\text{op}}^2 |\mathbf{G}|^2 \|D^2 \mathbf{F}\|_{\mathcal{H} \otimes \mathcal{H}}^2 \right] \right. \\ &\quad \left. + \mathbb{E} \left[\|\mathbf{\Gamma}(\mathbf{F})^{-1}\|_{\text{HS}}^2 \|\mathbf{DG}\|_{\mathcal{H}}^2 \|\mathbf{DF}\|_{\mathcal{H}}^2 \right] \right. \\ &\quad \left. + \mathbb{E} \left[\|\mathbf{\Gamma}(\mathbf{F})^{-1}\|_{\text{HS}}^4 |\mathbf{G}|^2 \|\mathbf{DF}\|_{\mathcal{H}}^4 \|D^2 \mathbf{F}\|_{\mathcal{H} \otimes \mathcal{H}}^2 \right] \right\}. \end{aligned}$$

Combining (B.23), (B.24) and (B.28), by Hölder's inequality, we can obtain that

$$\begin{aligned} &\mathbb{E} \left[\left| \delta \left(\sum_{j=1}^d (\mathbf{\Gamma}(\mathbf{F})^{-1}\mathbf{G})_j DF^j \right) \right|^2 \right] \\ &\leq C \|\mathbf{G}\|_{1,4}^2 \left\{ \mathbb{E} \left[(1 + \|\mathbf{DF}\|_{\mathcal{H}}^8) (1 + \|D^2 \mathbf{F}\|_{\mathcal{H} \otimes \mathcal{H}}^4) (1 + \|\mathbf{\Gamma}(\mathbf{F})^{-1}\|_{\text{HS}}^8) \right] \right\}^{1/2}. \end{aligned}$$

Substituting this into (B.22), we can obtain the desired.

It remains to prove (B.27). Observe that for any $1 \leq k, l \leq d$

$$\|D\langle DF^k, DF^l \rangle_{\mathcal{H}}\|_{\mathcal{H}} \leq \|D^2 F^k\|_{\mathcal{H} \otimes \mathcal{H}} \|DF^l\|_{\mathcal{H}} + \|DF^k\|_{\mathcal{H}} \|D^2 F^l\|_{\mathcal{H} \otimes \mathcal{H}}.$$

Then, for each $j = 1, \dots, d$, it holds that

$$\begin{aligned} & \sum_{i,k,l=1}^d \left| G^i(\mathbf{\Gamma}(\mathbf{F})^{-1})_{j,k}(\mathbf{\Gamma}(\mathbf{F})^{-1})_{l,i} \right| \left\| \mathbf{D} \langle \mathbf{D}F^k, \mathbf{D}F^l \rangle_{\mathcal{H}} \right\|_{\mathcal{H}} \\ & \leq \sum_{k=1}^d \left| (\mathbf{\Gamma}(\mathbf{F})^{-1})_{j,k} \right| \left\| \mathbf{D}^2 F^k \right\|_{\mathcal{H} \otimes \mathcal{H}} \sum_{i,l=1}^d \left| G^i(\mathbf{\Gamma}(\mathbf{F})^{-1})_{l,i} \right| \left\| \mathbf{D}F^l \right\|_{\mathcal{H}} \\ & \quad + \sum_{k=1}^d \left| (\mathbf{\Gamma}(\mathbf{F})^{-1})_{j,k} \right| \left\| \mathbf{D}F^k \right\|_{\mathcal{H}} \sum_{i,l=1}^d \left| G^i(\mathbf{\Gamma}(\mathbf{F})^{-1})_{l,i} \right| \left\| \mathbf{D}^2 F^l \right\|_{\mathcal{H} \otimes \mathcal{H}}. \end{aligned}$$

The Cauchy-Schwarz inequality implies that

$$\begin{aligned} \left\{ \sum_{i,l=1}^d \left| G^i(\mathbf{\Gamma}(\mathbf{F})^{-1})_{l,i} \right| \left\| \mathbf{D}F^l \right\|_{\mathcal{H}} \right\}^2 & \leq \left[\sum_{i=1}^d (G^i)^2 \right] \left[\sum_{i,l=1}^d (\mathbf{\Gamma}(\mathbf{F})^{-1})_{l,i}^2 \right] \left[\sum_{l=1}^d \left\| \mathbf{D}F^l \right\|_{\mathcal{H}}^2 \right] \\ & \leq |\mathbf{G}|^2 \left\| \mathbf{\Gamma}(\mathbf{F})^{-1} \right\|_{\text{HS}}^2 \left\| \mathbf{D}\mathbf{F} \right\|_{\mathcal{H}}^2, \end{aligned}$$

and similarly

$$\left\{ \sum_{i,l=1}^d \left| G^i(\mathbf{\Gamma}(\mathbf{F})^{-1})_{l,i} \right| \left\| \mathbf{D}^2 F^l \right\|_{\mathcal{H} \otimes \mathcal{H}} \right\}^2 \leq |\mathbf{G}|^2 \left\| \mathbf{\Gamma}(\mathbf{F})^{-1} \right\|_{\text{HS}}^2 \left\| \mathbf{D}^2 \mathbf{F} \right\|_{\mathcal{H}}^2.$$

Combining inequalities above, we can obtain

$$\begin{aligned} \mathcal{I}_{2,3} & \leq |\mathbf{G}| \left\| \mathbf{\Gamma}(\mathbf{F})^{-1} \right\|_{\text{HS}} \left\| \mathbf{D}\mathbf{F} \right\|_{\mathcal{H}} \sum_{j,k=1}^d \left| (\mathbf{\Gamma}(\mathbf{F})^{-1})_{j,k} \right| \left\| \mathbf{D}^2 F^k \right\|_{\mathcal{H} \otimes \mathcal{H}} \left\| \mathbf{D}F^j \right\|_{\mathcal{H}} \\ & \quad + |\mathbf{G}| \left\| \mathbf{\Gamma}(\mathbf{F})^{-1} \right\|_{\text{HS}} \left\| \mathbf{D}^2 \mathbf{F} \right\|_{\mathcal{H} \otimes \mathcal{H}} \sum_{j,k=1}^d \left| (\mathbf{\Gamma}(\mathbf{F})^{-1})_{j,k} \right| \left\| \mathbf{D}F^k \right\|_{\mathcal{H}} \left\| \mathbf{D}F^j \right\|_{\mathcal{H}} \\ & \leq 2 \left\| \mathbf{\Gamma}(\mathbf{F})^{-1} \right\|_{\text{HS}}^2 |\mathbf{G}| \left\| \mathbf{D}\mathbf{F} \right\|_{\mathcal{H}}^2 \left\| \mathbf{D}^2 \mathbf{F} \right\|_{\mathcal{H} \otimes \mathcal{H}}. \end{aligned}$$

The proof is complete. \square

Proof of Lemma 3.9. According to Markov's inequality, it holds

$$\begin{aligned} \mathbb{P}(E_{j,k}^1) & \leq \mathbb{P} \left(\eta_j \sum_{i=j+1}^k (\varphi(\boldsymbol{\xi}_i) - \mathbb{E}\varphi(\boldsymbol{\xi}_i)) > (k-j)\eta_k \right) \\ & \leq \left(\frac{\eta_j}{\eta_k} \right)^4 \mathbb{E} \left| \frac{1}{k-j} \sum_{i=j+1}^k (\varphi(\boldsymbol{\xi}_i) - \mathbb{E}\varphi(\boldsymbol{\xi}_i)) \right|^4, \end{aligned}$$

where in the first inequality we use $t_k - t_j \geq (k-j)\eta_k$ due to $\eta_i \leq \eta_j$ for any $i > j$. Assumption III implies that

$$\frac{\eta_j}{\eta_k} = \prod_{i=j+1}^k \frac{\eta_{i-1}}{\eta_i} \leq \prod_{i=j+1}^k (1 + \omega\eta_i) \leq \prod_{i=j+1}^k e^{\omega\eta_i} = e^{\omega(t_k - t_j)},$$

where the first inequality follows from the condition $\eta_{i-1} \leq \eta_i(1 + \omega\eta_i)$. Assumption IV guarantees that $\mathbb{E}|\varphi(\boldsymbol{\xi}_i) - \mathbb{E}[\varphi(\boldsymbol{\xi}_i)]|^4 < \infty$, then together with Lemma B.1, the

following holds for constant $C > 0$,

$$\mathbb{P}(E_{j,k}^1) \leq \frac{Ce^{4\omega(t_k-t_j)}}{(k-j)^2} \leq \frac{Ce^{4\omega(t_k-t_j)}}{(t_k-t_j)^2} \eta_j^2 \leq \frac{Ce^{6\omega(t_k-t_j)}}{(t_k-t_j)^2} \eta_k^2.$$

Since Assumption IV implies $\mathbb{E} |\varphi(\boldsymbol{\xi}_i)^2 - \mathbb{E}[\varphi(\boldsymbol{\xi}_i)^2]|^4 < +\infty$, $\mathbb{P}(E_{j,k}^2)$ can be analogously estimated only with $\varphi(\boldsymbol{\xi}_i)$ replaced by $\varphi(\boldsymbol{\xi}_i)^2$. \square

Recall SDE (1.6), we have

$$\begin{cases} \widetilde{\mathbf{M}}_{t_j} = (1 - \gamma\eta_j)\widetilde{\mathbf{M}}_{t_{j-1}} - \frac{1}{N} \sum_{r=1}^N \nabla F(\widetilde{\mathbf{X}}_{t_{j-1}}, \xi_j^r) \eta_j + \beta(\mathbf{B}_{t_j} - \mathbf{B}_{t_{j-1}}), \\ \widetilde{\mathbf{X}}_{t_j} = \eta_j \widetilde{\mathbf{M}}_{t_{j-1}} + \widetilde{\mathbf{X}}_{t_{j-1}}, \end{cases}$$

for $j \geq 1$ and initial value $(\widetilde{\mathbf{M}}_0, \widetilde{\mathbf{X}}_0) = (\mathbf{m}, \mathbf{x})$. By [Nua06, Section 2.2.2], the Malliavin derivative satisfies

$$(B.29) \quad \begin{bmatrix} D_s \widetilde{\mathbf{M}}_{t_j} \\ D_s \widetilde{\mathbf{X}}_{t_j} \end{bmatrix} = \mathbf{E}_j^\top \begin{bmatrix} D_s \widetilde{\mathbf{M}}_{t_{j-1}} \\ D_s \widetilde{\mathbf{X}}_{t_{j-1}} \end{bmatrix} + \beta \begin{bmatrix} \mathbf{I}_d \\ \mathbf{0} \end{bmatrix} \mathbf{1}_{[t_{j-1}, t_j)}(s),$$

where the matrix \mathbf{E}_j is defined for simplicity as

$$(B.30) \quad \mathbf{E}_j \equiv \begin{bmatrix} (1 - \gamma\eta_j)\mathbf{I}_d & \eta_j \mathbf{I}_d \\ -N^{-1} \sum_{r=1}^N \nabla^2 F(\widetilde{\mathbf{X}}_{t_{j-1}}, \xi_j^r) \eta_j & \mathbf{I}_d \end{bmatrix}.$$

Besides, define \mathbf{S}_j as

$$\mathbf{S}_k \equiv \mathbf{I}_{2d}, \quad \mathbf{S}_j \equiv \mathbf{E}_{j+1} \dots \mathbf{E}_k, \quad j = 1, \dots, k-1.$$

Given $(\widetilde{\mathbf{M}}_{t_\ell}, \widetilde{\mathbf{X}}_{t_\ell}) = \mathbf{z}$, the Malliavin derivative of $\widetilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} = (\widetilde{\mathbf{M}}_{t_k}, \widetilde{\mathbf{X}}_{t_k})$ has the following recursion:

$$(B.31) \quad \begin{bmatrix} D_s \widetilde{\mathbf{M}}_{t_k} \\ D_s \widetilde{\mathbf{X}}_{t_k} \end{bmatrix} = \beta \sum_{j=\ell+1}^k \mathbf{S}_j^\top \begin{bmatrix} \mathbf{I}_d \\ \mathbf{0} \end{bmatrix} \mathbf{1}_{[t_{j-1}, t_j)}(s),$$

and the corresponding Malliavin matrix is given by

$$(B.32) \quad \Gamma(\widetilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}}) = \int_{t_\ell}^{t_k} \begin{bmatrix} D_s \widetilde{\mathbf{M}}_{t_k} \\ D_s \widetilde{\mathbf{X}}_{t_k} \end{bmatrix} \begin{bmatrix} D_s \widetilde{\mathbf{M}}_{t_k} \\ D_s \widetilde{\mathbf{X}}_{t_k} \end{bmatrix}^\top ds = \beta^2 \sum_{j=\ell+1}^k \eta_j \mathbf{S}_j^\top \begin{bmatrix} \mathbf{I}_d \\ \mathbf{0} \end{bmatrix} [\mathbf{I}_d \quad \mathbf{0}] \mathbf{S}_j.$$

Before estimating the lower bound of eigenvalues of $\Gamma(\widetilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}})$, let us derive an estimate of $\|\mathbf{S}_j - \mathbf{I}_{2d}\|_{\text{op}}$ first.

Lemma B.2. *For any fixed $k \geq 1$, the following holds for each $1 \leq j \leq k$,*

$$\|\mathbf{S}_j - \mathbf{I}_{2d}\|_{\text{op}} \leq \exp \left\{ \sum_{i=j+1}^k \eta_i (\varphi(\boldsymbol{\xi}_i) + \gamma + 1) \right\} - 1,$$

where $\varphi(\boldsymbol{\xi}_i) = N^{-1} \sum_{r=1}^N \sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla^2 F(\mathbf{x}, \xi_i^r)\|_{\text{op}}$, $i = j+1, \dots, k$, are in (3.12).

Proof. We prove this by a backward induction. By definition, $\mathbf{S}_k = \mathbf{I}_{2d}$, so this result holds for $j = k$. Assume that

$$\|\mathbf{S}_\ell - \mathbf{I}_{2d}\|_{\text{op}} \leq \exp \left\{ \sum_{i=\ell+1}^k \eta_i (\varphi(\boldsymbol{\xi}_i) + \gamma + 1) \right\} - 1,$$

holds for some $\ell \geq 2$, then

$$\|\mathbf{S}_{\ell-1} - \mathbf{I}_{2d}\|_{\text{op}} = \|\mathbf{E}_\ell \mathbf{S}_\ell - \mathbf{E}_\ell + \mathbf{E}_\ell - \mathbf{I}_{2d}\|_{\text{op}} \leq \|\mathbf{E}_\ell\|_{\text{op}} \|\mathbf{S}_\ell - \mathbf{I}_{2d}\|_{\text{op}} + \|\mathbf{E}_\ell - \mathbf{I}_{2d}\|_{\text{op}}.$$

For any matrix $\mathbf{K} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{2d \times 2d}$ and $\mathbf{w} = (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^{2d}$, it holds

$$|\mathbf{K}\mathbf{w}| = |\mathbf{A}\mathbf{u} + \mathbf{B}\mathbf{v} + \mathbf{C}\mathbf{u}| \leq [\|\mathbf{A}\|_{\text{op}} + \|\mathbf{B}\|_{\text{op}} + \|\mathbf{C}\|_{\text{op}}] |\mathbf{w}|,$$

which implies that $\|\mathbf{K}\|_{\text{op}} \leq \|\mathbf{A}\|_{\text{op}} + \|\mathbf{B}\|_{\text{op}} + \|\mathbf{C}\|_{\text{op}}$. Then, we have

$$\|\mathbf{E}_\ell - \mathbf{I}_{2d}\|_{\text{op}} = \eta_\ell \left\| \begin{bmatrix} -\gamma \mathbf{I}_d & \mathbf{I}_d \\ -\frac{1}{N} \sum_{r=1}^N \nabla^2 F(\tilde{\mathbf{X}}_{t_{\ell-1}}, \xi_\ell^r) & \mathbf{0} \end{bmatrix} \right\|_{\text{op}} \leq \eta_\ell (\varphi(\xi_\ell) + \gamma + 1).$$

Hence, we obtain the desired result by

$$\begin{aligned} \|\mathbf{S}_{\ell-1} - \mathbf{I}_{2d}\|_{\text{op}} &\leq [\|\mathbf{E}_\ell - \mathbf{I}_{2d}\|_{\text{op}} + 1] \|\mathbf{S}_\ell - \mathbf{I}_{2d}\|_{\text{op}} - 1 \\ &\leq \exp \left\{ \sum_{i=\ell}^k \eta_i (\varphi(\xi_i) + \gamma + 1) \right\} - 1. \end{aligned}$$

The proof is complete. \square

Next, we prove Lemmas 3.10 – 3.13. To make notations simple, we define

$$(B.33) \quad \mathbf{b}(\mathbf{y}) = \begin{bmatrix} -\gamma \mathbf{m} - \nabla f(\mathbf{x}) \\ \mathbf{m} \end{bmatrix}, \quad \tilde{\mathbf{b}}(\mathbf{y}, \xi) = \begin{bmatrix} -\gamma \mathbf{m} - N^{-1} \sum_{i=1}^N \nabla F(\mathbf{x}, \xi^i) \\ \mathbf{m} \end{bmatrix},$$

for $\mathbf{y} = (\mathbf{m}, \mathbf{x})$ and $\xi = (\xi^1, \dots, \xi^N)$.

Proof of Lemma 3.10. (i) By the definition and (B.32),

$$\left\| D\tilde{\mathbf{Y}}_{t_\ell, t_k}^z \right\|_{\mathcal{H}}^2 = \text{tr} \left(\Gamma(\tilde{\mathbf{Y}}_{t_\ell, t_k}^z) \right) = \beta^2 \sum_{j=\ell+1}^k \eta_j \text{tr} \left(\mathbf{S}_j^\top \begin{bmatrix} \mathbf{I}_d \\ \mathbf{0} \end{bmatrix} [\mathbf{I}_d \quad \mathbf{0}] \mathbf{S}_j \right) \leq 2d\beta^2 \sum_{j=\ell+1}^k \eta_j \|\mathbf{S}_j\|_{\text{op}}^2.$$

Together with Hölder's inequality, we have

$$(B.34) \quad \left\| D\tilde{\mathbf{Y}}_{t_\ell, t_k}^z \right\|_{\mathcal{H}}^8 \leq Cd^4 (t_k - t_\ell)^3 \sum_{j=\ell+1}^k \eta_j \|\mathbf{S}_j\|_{\text{op}}^8.$$

Notice that $\mathbf{S}_j = \mathbf{E}_{j+1} \dots \mathbf{E}_k$ with \mathbf{E}_j defined in (B.30) implies

$$\mathbb{E} \|\mathbf{S}_j\|_{\text{op}}^8 \leq \mathbb{E} \left\{ \mathbb{E} \left[\left(1 + \|\mathbf{E}_k - \mathbf{I}_{2d}\|_{\text{op}} \right)^8 \left| (\tilde{\mathbf{X}}_{t_{i-1}}, \xi_i)_{j+1 \leq i \leq k-1} \right. \right] \|\mathbf{E}_{j+1} \dots \mathbf{E}_{k-1}\|_{\text{op}}^8 \right\},$$

and

$$\begin{aligned} \left(1 + \|\mathbf{E}_k - \mathbf{I}_{2d}\|_{\text{op}} \right)^8 &\leq \left[1 + \eta_k \left(\frac{1}{N} \sum_{r=1}^N \sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla^2 F(\mathbf{x}, \xi_k^r)\|_{\text{op}} + \gamma + 1 \right) \right]^8 \\ &\leq \frac{1}{1 - 7\eta_k} \left[1 + \eta_k \left(\frac{1}{N} \sum_{r=1}^N \sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla^2 F(\mathbf{x}, \xi_k^r)\|_{\text{op}} + \gamma + 1 \right) \right]^8 \\ &\leq \left(1 + \frac{7\eta_k}{1 - 7\eta_1} \right) \left[1 + C\eta_k \left(\frac{1}{N} \sum_{r=1}^N \sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla^2 F(\mathbf{x}, \xi_k^r)\|_{\text{op}}^8 + 1 \right) \right], \end{aligned}$$

where in the second inequality we use $(1 + \eta_k x)^8 \leq 1 + \eta_k x^8 + 7\eta_k(1 + \eta_k x)^8$ for $x \geq 0$. According to Assumption IV, we have

$$\mathbb{E} \|\mathbf{S}_j\|_{\text{op}}^8 \leq (1 + C\eta_k) \mathbb{E} \|\mathbf{E}_{j+1} \dots \mathbf{E}_{k-1}\|_{\text{op}}^8 \leq e^{C\eta_k} \mathbb{E} \|\mathbf{E}_{j+1} \dots \mathbf{E}_{k-1}\|_{\text{op}}^8.$$

Then, applying this method recursively shows that $\mathbb{E} \|\mathbf{S}_j\|_{\text{op}}^8 \leq e^{C(t_k - t_j)}$ for $j = \ell + 1, \dots, k$, so (B.34) and condition $t_k - t_\ell \leq 1/[5(B_1 + \gamma + 2)]$ imply that

$$\mathbb{E} \left\| D\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} \right\|_{\mathcal{H}}^8 \leq Cd^4 (t_k - t_\ell)^3 \sum_{j=\ell+1}^k \eta_j e^{C(t_k - t_\ell)} \leq Cd^4.$$

Analogously, it can be shown that $\mathbb{E} \|D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}\|_{\mathcal{H}}^8 \leq Cd^4$.

For any $t \in [t_{k-1}, t_k]$ and $\mathbf{h} \in \mathcal{H}$ with $\|\mathbf{h}\|_{\mathcal{H}} \leq 1$, by [Nor86]

$$(B.35) \quad \begin{aligned} D^{\mathbf{h}} \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) &= D^{\mathbf{h}} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} + \beta \int_{t_{k-1}}^t \begin{bmatrix} \mathbf{I}_d \\ \mathbf{0} \end{bmatrix} \mathbf{h}_s ds \\ &+ \int_{t_{k-1}}^t \nabla \mathbf{b}(\mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}})) D^{\mathbf{h}} \mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) du. \end{aligned}$$

Since $\|\nabla \mathbf{b}(\mathbf{y})\|_{\text{op}} \leq L + \gamma + 1$ where $\mathbf{b}(\mathbf{y})$ is in (B.33), it follows that

$$\left\| D\mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}} \leq \left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} + \beta\sqrt{\eta_k} + C \int_{t_{k-1}}^t \left\| D\mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}} du.$$

Then Grönwall's inequality implies

$$\mathbb{E} \left\| D\mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}}^8 \leq C \left(\mathbb{E} \left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}}^8 + 1 \right) \leq Cd^4, \quad t \in [t_{k-1}, t_k].$$

(ii) Given event $(E_{\ell, k}^1)^c$, i.e., $\sum_{i=\ell+1}^k \eta_i (\varphi(\boldsymbol{\xi}_i) - \mathbb{E}\varphi(\boldsymbol{\xi}_i)) \leq t_k - t_\ell$, Lemma B.2 and Assumption IV derive the following equation for some constant $C > 0$,

$$\|\mathbf{S}_j\|_{\text{op}} \leq \|\mathbf{S}_j - \mathbf{I}_{2d}\|_{\text{op}} + 1 \leq e^{\sum_{i=j+1}^k \eta_i (\varphi(\boldsymbol{\xi}_i) + \gamma + 1)} \leq e^{(t_k - t_\ell)(\mathbb{E}\varphi(\boldsymbol{\xi}_i) + \gamma + 2)} \leq C,$$

for $j = \ell + 1, \dots, k$. According to (B.34), the following holds ,

$$\left\| D\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} \right\|_{\mathcal{H}}^8 \leq Cd^4 (t_k - t_\ell)^3 \sum_{j=\ell+1}^k \eta_j \|\mathbf{S}_j\|_{\text{op}}^8 \leq Cd^4 \quad \text{on } (E_{\ell, k}^1)^c.$$

The proof is complete. \square

Proof of Lemma 3.11. For any $\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}$ with $\|\mathbf{h}_1\|_{\mathcal{H}}, \|\mathbf{h}_2\|_{\mathcal{H}} \leq 1$, by [Nor86], we have

$$\begin{aligned} D^{\mathbf{h}_1} D^{\mathbf{h}_2} \tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} &= D^{\mathbf{h}_1} D^{\mathbf{h}_2} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} + \eta_k \nabla \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}, \boldsymbol{\xi}_k) D^{\mathbf{h}_1} D^{\mathbf{h}_2} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \\ &+ \eta_k \nabla^2 \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}, \boldsymbol{\xi}_k) D^{\mathbf{h}_1} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} D^{\mathbf{h}_2} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}, \end{aligned}$$

where $\tilde{\mathbf{b}}(\mathbf{y}, \boldsymbol{\xi})$ is defined in (B.33). Hence,

$$\begin{aligned} \left| D^{\mathbf{h}_1} D^{\mathbf{h}_2} \tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} \right| &\leq \left| D^{\mathbf{h}_1} D^{\mathbf{h}_2} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right| + \eta_k \left\| \nabla \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}, \boldsymbol{\xi}_k) \right\|_{\text{op}} \left| D^{\mathbf{h}_1} D^{\mathbf{h}_2} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right| \\ &+ \eta_k \left\| \nabla^2 \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}, \boldsymbol{\xi}_k) \right\|_{\text{op}} \left| D^{\mathbf{h}_1} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right| \left| D^{\mathbf{h}_2} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right|, \end{aligned}$$

which implies that

$$\begin{aligned} \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_k}^z \right\|_{\mathcal{H} \otimes \mathcal{H}} &\leq \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H} \otimes \mathcal{H}} + \eta_k \left\| \nabla \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z, \boldsymbol{\xi}_k) \right\|_{\text{op}} \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H} \otimes \mathcal{H}} \\ &\quad + \eta_k \left\| \nabla^2 \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z, \boldsymbol{\xi}_k) \right\|_{\text{op}} \left\| D \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H}}^2. \end{aligned}$$

Notice that $\eta_k \leq \eta_1$ and $(x + \eta_k y)^4 \leq x^4 + \eta_k y^4 + 3\eta_k(x + \eta_k y)^4$ for $x, y \geq 0$, so

$$\begin{aligned} \mathbb{E} \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_k}^z \right\|_{\mathcal{H} \otimes \mathcal{H}}^4 &\leq \left(1 + \frac{3\eta_k}{1 - 3\eta_1} \right) \mathbb{E} \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H} \otimes \mathcal{H}}^4 \\ (B.36) \quad &\quad + \frac{8\eta_k}{1 - 3\eta_1} \left\{ \mathbb{E} \left[\left\| \nabla \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z, \boldsymbol{\xi}_k) \right\|_{\text{op}}^4 \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H} \otimes \mathcal{H}}^4 \right] \right. \\ &\quad \left. + \mathbb{E} \left[\left\| \nabla^2 \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z, \boldsymbol{\xi}_k) \right\|_{\text{op}}^4 \left\| D \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H}}^8 \right] \right\}. \end{aligned}$$

Combining Assumption IV, Lemma 3.10, and the independence between $\boldsymbol{\xi}_k$ and other random variables on the right-hand side, we have

$$\begin{aligned} &\mathbb{E} \left[\left\| \nabla^2 \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z, \boldsymbol{\xi}_k) \right\|_{\text{op}}^4 \left\| D \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H}}^8 \right] \\ (B.37) \quad &= \mathbb{E} \left\{ \mathbb{E}_{\boldsymbol{\xi}_k} \left[\left\| \nabla^2 \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z, \boldsymbol{\xi}_k) \right\|_{\text{op}}^4 \right] \left\| D \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H}}^8 \right\} \\ &\leq C \mathbb{E} \left\| D \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H}}^8 \leq C d^4. \end{aligned}$$

Similarly, by Assumption IV, and the independence between $\boldsymbol{\xi}_k$ and other random variables, we also have

$$(B.38) \quad \mathbb{E} \left[\left\| \nabla \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z, \boldsymbol{\xi}_k) \right\|_{\text{op}}^4 \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H} \otimes \mathcal{H}}^4 \right] \leq C \mathbb{E} \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H} \otimes \mathcal{H}}^4.$$

Combining (B.36), (B.37) and (B.38), since $D^2 \tilde{\mathbf{Y}}_{t_\ell, t_\ell}^z = \mathbf{0}$, we have by an induction

$$\begin{aligned} \mathbb{E} \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_k}^z \right\|_{\mathcal{H} \otimes \mathcal{H}}^4 &\leq (1 + C\eta_k) \mathbb{E} \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \right\|_{\mathcal{H} \otimes \mathcal{H}}^4 + C d^4 \eta_k \\ &\leq C d^4 \sum_{i=\ell+1}^k \eta_i \prod_{j=i+1}^k (1 + C\eta_j) \leq C d^4 (t_k - t_\ell) e^{C(t_k - t_\ell)}. \end{aligned}$$

Due to $t_k - t_\ell \leq 1/[5(B_1 + \gamma + 1)]$, we can obtain the desired.

For $\mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z)$, $t \in [t_{k-1}, t_k]$ and any $\mathbf{h}_i \in \mathcal{H}$ with $\|\mathbf{h}_i\|_{\mathcal{H}} \leq 1$, $i = 1, 2$, by [Nor86], the second order Malliavin derivation satisfies

$$\begin{aligned} D^{\mathbf{h}_1} D^{\mathbf{h}_2} \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z) &= D^{\mathbf{h}_1} D^{\mathbf{h}_2} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z \\ &\quad + \int_{t_{k-1}}^t \nabla \mathbf{b}(\mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z)) D^{\mathbf{h}_1} D^{\mathbf{h}_2} \mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z) du \\ &\quad + \int_{t_{k-1}}^t \nabla^2 b(\mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z)) D^{\mathbf{h}_1} \mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z) D^{\mathbf{h}_2} \mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^z) du. \end{aligned}$$

Since $\|\nabla \mathbf{b}(\mathbf{y})\|_{\text{op}} \leq L + \gamma + 1$ and $\|\nabla^2 b(\mathbf{y})\|_{\text{op}} \leq B_2$ due to Assumption IV, we have

$$\begin{aligned} \left\| D^2 \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H} \otimes \mathcal{H}} &\leq \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H} \otimes \mathcal{H}} + C \int_{t_{k-1}}^t \left\| D^2 \mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H} \otimes \mathcal{H}} du \\ &\quad + C \int_{t_{k-1}}^{t_k} \left\| D \mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}}^2 du. \end{aligned}$$

By Grönwall's inequality,

$$\left\| D^2 \mathbf{Y}_{t_{k-1}, t_k}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H} \otimes \mathcal{H}} \leq C \left\| D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H} \otimes \mathcal{H}} + C \int_{t_{k-1}}^{t_k} \left\| D \mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}}^2 du.$$

We have shown that $\mathbb{E} \|D^2 \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}\|_{\mathcal{H} \otimes \mathcal{H}}^4 \leq Cd^4$, together with Lemma 3.10 and Hölder's inequality, we get

$$\mathbb{E} \left\| D^2 \mathbf{Y}_{t_{k-1}, t_k}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H} \otimes \mathcal{H}}^4 \leq Cd^4.$$

The proof is complete. \square

Proof of Lemma 3.12. Recall the SDEs (1.5) and (1.6), and define

$$\Xi_{t_\ell, t_k} := \mathbf{Y}_{t_{k-1}, t_k}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - \tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}} = \int_{t_{k-1}}^{t_k} \left(\mathbf{b}(\mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}})) - \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}, \boldsymbol{\xi}_k) \right) dt,$$

where functions \mathbf{b} and $\tilde{\mathbf{b}}$ are defined in (B.33). For any $\mathbf{h} \in \mathcal{H}$ with $\|\mathbf{h}\|_{\mathcal{H}} \leq 1$, by [Nor86], we have

$$(B.39) \quad D^{\mathbf{h}} \Xi_{t_\ell, t_k} = \int_{t_{k-1}}^{t_k} \left(\nabla \mathbf{b}(\mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}})) D^{\mathbf{h}} \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - \nabla \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}, \boldsymbol{\xi}_k) D^{\mathbf{h}} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right) dt.$$

(i) Observe that $D^{\mathbf{h}} \Xi_{t_\ell, t_k} = \mathcal{I}_1 + \mathcal{I}_2 + \mathcal{I}_3$, where

$$(B.40) \quad \begin{aligned} \mathcal{I}_1 &= \int_{t_{k-1}}^{t_k} \nabla \mathbf{b}(\mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}})) \left(D^{\mathbf{h}} \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - D^{\mathbf{h}} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right) dt, \\ \mathcal{I}_2 &= \int_{t_{k-1}}^{t_k} \left(\nabla \mathbf{b}(\mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}})) - \nabla \mathbf{b}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right) D^{\mathbf{h}} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} dt, \\ \mathcal{I}_3 &= \eta_k \left(\nabla \mathbf{b}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - \nabla \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}, \boldsymbol{\xi}_k) \right) D^{\mathbf{h}} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}. \end{aligned}$$

We estimate them one by one.

For \mathcal{I}_1 , one has

$$\begin{aligned} &D^{\mathbf{h}} \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - D^{\mathbf{h}} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \\ &= \beta \int_{t_{k-1}}^t \begin{bmatrix} \mathbf{I}_d \\ \mathbf{0} \end{bmatrix} \mathbf{h}_s ds + \int_{t_{k-1}}^t \nabla \mathbf{b}(\mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}})) D^{\mathbf{h}} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} du \\ &\quad + \int_{t_{k-1}}^t \nabla \mathbf{b}(\mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}})) \left[D^{\mathbf{h}} \mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - D^{\mathbf{h}} \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right] du. \end{aligned}$$

Together with $\|\nabla \mathbf{b}(\mathbf{y})\|_{\text{op}} \leq L + \gamma + 1$, the following holds,

$$\left\| D \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - D \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} \leq C \int_{t_{k-1}}^t \left\| D \mathbf{Y}_{t_{k-1}, u}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - D \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} du$$

$$+ C\sqrt{\eta_k} \left(\left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} + 1 \right),$$

for all $t \in [t_{k-1}, t_k]$. And Grönwall's inequality derives

$$(B.41) \quad \left\| D\mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} \leq C\sqrt{\eta_k} \left(\left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} + 1 \right).$$

which leads to

$$(B.42) \quad |\mathcal{I}_1|^4 \leq C\eta_k^6 \left(\left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}}^4 + 1 \right) \|\mathbf{h}\|_{\mathcal{H}}^4.$$

For \mathcal{I}_2 , we let $\Delta = \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}$, then

$$\mathcal{I}_2 = D^{\mathbf{h}}\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \int_{t_{k-1}}^{t_k} \int_0^1 \nabla^2 \mathbf{b}(r\Delta + \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \cdot \Delta \, dr \, dt.$$

By $\|\nabla^2 \mathbf{b}\|_{\text{op}} \leq B_2$ due to Assumption IV, we have

$$(B.43) \quad |\mathcal{I}_2|^4 \leq C \left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}}^4 \eta_k^3 \int_{t_{k-1}}^{t_k} \left| \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right|^4 dt \|\mathbf{h}\|_{\mathcal{H}}^4.$$

Besides, by Lemmas 3.2 and 3.4, we have for any $t \in [t_{k-1}, t_k]$,

$$\begin{aligned} \mathbb{E} \left[\left| \mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) - \tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right|^4 \right] &\leq C(t - t_{k-1})^2 \left(\mathbb{E}[\mathcal{V}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}})^2] + d^2 \right) \\ &\leq C\eta_k^2 (\mathcal{V}(\mathbf{z})^2 + d^2). \end{aligned}$$

Then, by Markov property, it holds that

$$(B.44) \quad \mathbb{E} |\mathcal{I}_2|^4 \leq C\eta_k^6 (\mathcal{V}(\mathbf{z})^2 + d^2) \|\mathbf{h}\|_{\mathcal{H}}^4 \mathbb{E} \left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}}^4.$$

For \mathcal{I}_3 , since $\|\mathbf{A}\|_{\text{op}} \leq \|\mathbf{A}\|_{\text{HS}} \leq \sqrt{d} \|\mathbf{A}\|_{\text{op}}$ for any matrix \mathbf{A} , it follows from Lemma B.1 that

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\xi}} \left\| \nabla \mathbf{b}(\mathbf{y}) - \nabla \tilde{\mathbf{b}}(\mathbf{y}, \boldsymbol{\xi}) \right\|_{\text{op}}^4 &= \mathbb{E}_{\boldsymbol{\xi}} \left\| \frac{1}{N} \sum_{r=1}^N (\nabla^2 F(\mathbf{x}, \xi^r) - \nabla^2 f(\mathbf{x})) \right\|_{\text{op}}^4 \\ &\leq \frac{C}{N^2} \mathbb{E}_{\boldsymbol{\xi}} \left\| \nabla^2 F(\mathbf{x}, \xi^r) - \nabla^2 f(\mathbf{x}) \right\|_{\text{HS}}^4 \leq \frac{Cd^2}{N^2}. \end{aligned}$$

Due $\boldsymbol{\xi}$ is independent of other random variables, this immediately leads to

$$(B.45) \quad \mathbb{E} |\mathcal{I}_3|^4 \leq C\eta_k^4 \|\mathbf{h}\|_{\mathcal{H}}^4 \frac{d^2}{N^2} \mathbb{E} \left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}}^4.$$

Since $\mathbb{E} \left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}}^4 \leq Cd^2$ by Lemma 3.10, then Combining (B.40), (B.42), (B.44) and (B.45), we get

$$\mathbb{E} \left\| D\boldsymbol{\Xi}_{t_\ell, t_k} \right\|_{\mathcal{H}}^4 \leq Cd^2 \eta_k^6 (\mathcal{V}(\mathbf{z})^2 + d^2) + C \frac{d^4 \eta_k^4}{N^2}.$$

(ii) According to (B.39), we have that

$$(B.46) \quad \begin{aligned} \left\| D\boldsymbol{\Xi}_{t_\ell, t_k} \right\|_{\mathcal{H}} &\leq \int_{t_{k-1}}^{t_k} \left\| \nabla \mathbf{b} \right\|_{\text{op}} \left\| D\mathbf{Y}_{t_{k-1}, t}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}} dt \\ &\quad + \eta_k \left\| \nabla \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}}, \boldsymbol{\xi}_k) \right\|_{\text{op}} \left\| D\tilde{\mathbf{Y}}_{t_\ell, t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} \end{aligned}$$

(B.41) together with a triangle inequality, it holds

$$\left\| \mathbf{D}\mathbf{Y}_{t_{k-1},t}(\tilde{\mathbf{Y}}_{t_\ell,t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}} - \left\| \mathbf{D}\tilde{\mathbf{Y}}_{t_\ell,t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} \leq C\sqrt{\eta_k} \left(\left\| \mathbf{D}\tilde{\mathbf{Y}}_{t_\ell,t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} + 1 \right)$$

for all $t \in [t_{k-1}, t_k]$. Besides, by Lemma 3.10, $\left\| \mathbf{D}\tilde{\mathbf{Y}}_{t_\ell,t_{k-1}}^{\mathbf{z}} \right\|_{\mathcal{H}} \leq C\sqrt{d}$ on the event $(E_{\ell,k}^1)^c$. Then, we have

$$(B.47) \quad \left\| \mathbf{D}\mathbf{Y}_{t_{k-1},t}(\tilde{\mathbf{Y}}_{t_\ell,t_{k-1}}^{\mathbf{z}}) \right\|_{\mathcal{H}} \leq C\sqrt{d}, \quad \text{on the event } (E_{\ell,k}^1)^c.$$

On the other hand, given $(E_{\ell,k}^2)^c$, i.e., $\sum_{i=\ell+1}^k \eta_i(\varphi(\boldsymbol{\xi}_i)^2 - \mathbb{E}[\varphi(\boldsymbol{\xi}_i)^2]) \leq t_k - t_\ell$ with $\varphi(\boldsymbol{\xi}_i)$ defined by (3.12), we have

$$(B.48) \quad \eta_k \varphi(\boldsymbol{\xi}_k)^2 \leq \sum_{i=\ell+1}^k \eta_i \varphi(\boldsymbol{\xi}_i)^2 \leq (t_k - t_\ell)(\mathbb{E}[\varphi(\boldsymbol{\xi}_k)^2] + 1) \leq C,$$

where the last inequality is due to Assumption IV.

Combining $\|\nabla \mathbf{b}\|_{\text{op}} \leq L + \gamma + 1$, (B.46), (B.47) and (B.48), we have

$$\begin{aligned} \|\mathbf{D}\boldsymbol{\Xi}_{t_\ell,t_k}\|_{\mathcal{H}} &\leq C\sqrt{d}\eta_k \left(\left\| \nabla \tilde{\mathbf{b}}(\tilde{\mathbf{Y}}_{t_\ell,t_{k-1}}^{\mathbf{z}}, \boldsymbol{\xi}_k) \right\|_{\text{op}} + 1 \right) \\ &\leq C\sqrt{d}\eta_k(\varphi(\boldsymbol{\xi}_k) + \gamma + 2) \leq C\sqrt{d\eta_k}, \quad \text{on the event } (E_{\ell,k}^1 \cup E_{\ell,k}^2)^c. \end{aligned}$$

The proof is complete. \square

Proof of Lemma 3.13. For any vector $\mathbf{w} \in \mathbb{R}^{2d}$ with $|\mathbf{w}| = 1$, the relation (B.32) yields that

$$(B.49) \quad \mathbf{w}^\top \boldsymbol{\Gamma}(\tilde{\mathbf{Y}}_{t_\ell,t_k}^{\mathbf{z}}) \mathbf{w} = \beta^2 \sum_{j=\ell+1}^k \eta_j \left| [\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_j \mathbf{w} \right|^2.$$

Inspired by [DP14], we divide \mathbf{w} and \mathbf{S}_j into

$$\mathbf{w} = \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}, \quad \mathbf{S}_j = \begin{bmatrix} \mathbf{A}_j & \mathbf{B}_j \\ \mathbf{C}_j & \mathbf{D}_j \end{bmatrix},$$

with vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ and matrices $\mathbf{A}_j, \mathbf{B}_j, \mathbf{C}_j, \mathbf{D}_j \in \mathbb{R}^{d \times d}$.

We split the proof into two cases, one being $|\mathbf{u}| > (t_k - t_\ell)/4$ and the other $|\mathbf{u}| \leq (t_k - t_\ell)/4$.

(i) The case $|\mathbf{u}| > (t_k - t_\ell)/4$. Observe that

$$(B.50) \quad \begin{aligned} \left| [\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_j \mathbf{w} \right|^2 &\geq \frac{1}{2} \left| [\mathbf{I}_d \ \mathbf{0}] \mathbf{w} \right|^2 - \left| [\mathbf{I}_d \ \mathbf{0}] (\mathbf{S}_j - \mathbf{I}_{2d}) \mathbf{w} \right|^2 \\ &\geq \frac{1}{2} |\mathbf{u}|^2 - \|\mathbf{S}_j - \mathbf{I}_{2d}\|_{\text{op}}^2. \end{aligned}$$

Assumption IV and the definition of m yield that

$$(B.51) \quad (t_k - t_m)(\mathbb{E}\varphi(\boldsymbol{\xi}_k) + \gamma + 2) \leq \frac{t_k - t_\ell}{10}.$$

Given $(E_{m,k}^1)^c$, i.e., $\sum_{i=m+1}^k \eta_i(\varphi(\boldsymbol{\xi}_i) - \mathbb{E}\varphi(\boldsymbol{\xi}_i)) \leq t_k - t_m$, we have

$$(B.52) \quad \sum_{i=m+1}^k \eta_i(\varphi(\boldsymbol{\xi}_i) + \gamma + 1) \leq (t_k - t_m)(\mathbb{E}\varphi(\boldsymbol{\xi}_k) + \gamma + 2).$$

Combining (B.51), (B.52) and Lemma B.2, we have that for $j = m, \dots, k$,

$$\|\mathbf{S}_j - \mathbf{I}_{2d}\|_{\text{op}} \leq e^{\sum_{i=m+1}^k \eta_i(\varphi(\xi_i) + \gamma + 1)} - 1 \leq e^{(t_k - t_\ell)/10} - 1 \leq \frac{t_k - t_\ell}{8},$$

where the last inequality follows from $e^{x/10} - 1 \leq x/8$ for $x \in (0, 1/10)$. Substituting this into (B.50) and $|\mathbf{u}| > (t_k - t_\ell)/4$ imply

$$|[\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_j \mathbf{w}|^2 \geq \frac{(t_k - t_\ell)^2}{64}.$$

Then, due to $t_k - t_{m-1} \geq (t_k - t_\ell)/[10(\gamma + B_1 + 2)]$ by the definition of m , we obtain

$$(B.53) \quad \mathbf{w}^\top \Gamma(\tilde{\mathbf{Y}}_{t_\ell, t_k}^z) \mathbf{w} \geq \beta^2 \sum_{j=m}^k \eta_j |[\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_j \mathbf{w}|^2 \geq C(t_k - t_\ell)^3, \text{ on } (E_{m,k}^1)^c.$$

(ii) The case $|\mathbf{u}| \leq (t_k - t_\ell)/4$. In this case we clearly have $|\mathbf{v}| = (1 - |\mathbf{u}|^2)^{1/2} \geq 3/4$. Notice that $\mathbf{S}_{i-1} = \mathbf{E}_i \mathbf{S}_i$ implies

$$[\mathbf{A}_{i-1} \ \mathbf{B}_{i-1}] = [(1 - \gamma\eta_i)\mathbf{I}_d \ \eta_i\mathbf{I}_d] \begin{bmatrix} \mathbf{A}_i & \mathbf{B}_i \\ \mathbf{C}_i & \mathbf{D}_i \end{bmatrix},$$

i.e.

$$[\mathbf{A}_{i-1} \ \mathbf{B}_{i-1}] - [\mathbf{A}_i \ \mathbf{B}_i] = -\gamma\eta_i [\mathbf{A}_i \ \mathbf{B}_i] + \eta_i [\mathbf{C}_i \ \mathbf{D}_i].$$

Then, summing this over i , we can obtain that

$$[\mathbf{A}_j \ \mathbf{B}_j] - [\mathbf{A}_k \ \mathbf{B}_k] = -\gamma \sum_{i=j+1}^k \eta_i [\mathbf{A}_i \ \mathbf{B}_i] + \sum_{i=j+1}^k \eta_i [\mathbf{C}_i \ \mathbf{D}_i].$$

Due to $\mathbf{A}_k = \mathbf{I}_d$ and $\mathbf{B}_k = \mathbf{0}$, we get the following relation,

$$(B.54) \quad \mathbf{L}_j := [\mathbf{A}_j \ \mathbf{B}_j] + \gamma \sum_{i=j+1}^k \eta_i [\mathbf{A}_i \ \mathbf{B}_i] = \sum_{i=j+1}^k \eta_i [\mathbf{C}_i \ \mathbf{D}_i] + [\mathbf{I}_d \ \mathbf{0}] =: \mathbf{R}_j,$$

where we denote the matrices on the left-hand side and right-hand side by \mathbf{L}_j and \mathbf{R}_j respectively.

For \mathbf{L}_j , observe that $[\mathbf{A}_j \ \mathbf{B}_j] = [\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_j$, then we have

$$\begin{aligned} |\mathbf{L}_j \mathbf{w}|^2 &\leq 2 |[\mathbf{A}_j \ \mathbf{B}_j] \mathbf{w}|^2 + 2\gamma^2 \left(\sum_{i=\ell+1}^k \eta_i |[\mathbf{A}_i \ \mathbf{B}_i] \mathbf{w}| \right)^2 \\ &\leq 2 |[\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_j \mathbf{w}|^2 + 2\gamma^2 (t_k - t_\ell) \sum_{i=\ell+1}^k \eta_i |[\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_i \mathbf{w}|^2, \end{aligned}$$

which immediately implies that

$$(B.55) \quad \sum_{j=\ell+1}^k \eta_j |\mathbf{L}_j \mathbf{w}|^2 \leq 2[1 + \gamma^2(t_k - t_\ell)^2] \sum_{j=\ell+1}^k \eta_j |[\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_j \mathbf{w}|^2.$$

For \mathbf{R}_j , observe that $[\mathbf{C}_j \ \mathbf{D}_j] = [\mathbf{0} \ \mathbf{I}_d] \mathbf{S}_j$, then

$$\mathbf{R}_j \mathbf{w} = \sum_{i=j+1}^k \eta_i [\mathbf{0} \ \mathbf{I}_d] \mathbf{w} + \sum_{i=j+1}^k \eta_i ([\mathbf{C}_i \ \mathbf{D}_i] - [\mathbf{0} \ \mathbf{I}_d]) \mathbf{w} + [\mathbf{I}_d \ \mathbf{0}] \mathbf{w}$$

$$= (t_k - t_j)\mathbf{v} + \sum_{i=j+1}^k \eta_i \left([\mathbf{C}_i \ \mathbf{D}_i] - [\mathbf{0} \ \mathbf{I}_d] \right) \mathbf{w} + \mathbf{u}.$$

Given $(E_{\ell,k}^1)^c$, together with Lemma B.2 implies

$$\begin{aligned} \left\| [\mathbf{C}_i \ \mathbf{D}_i] - [\mathbf{0} \ \mathbf{I}_d] \right\|_{\text{op}} &\leq \|\mathbf{S}_i - \mathbf{I}_{2d}\|_{\text{op}} \leq e^{\sum_{i=\ell+1}^k \eta_i (\varphi(\boldsymbol{\xi}_i) + \gamma + 1)} - 1 \\ &\leq e^{(t_k - t_\ell)(\mathbb{E}\varphi(\boldsymbol{\xi}_k) + \gamma + 2)} - 1 \leq e^{1/5} - 1 \leq \frac{1}{4}. \end{aligned}$$

Together with $|\mathbf{u}| \leq (t_k - t_\ell)/4$ and $|\mathbf{v}| \geq 3/4$, we have

$$\begin{aligned} |\mathbf{R}_j \mathbf{w}| &\geq (t_k - t_j) |\mathbf{v}| - \sum_{i=j+1}^k \eta_i \left| \left([\mathbf{C}_i \ \mathbf{D}_i] - [\mathbf{0} \ \mathbf{I}_d] \right) \mathbf{w} \right| - |\mathbf{u}| \\ &\geq \frac{3}{4}(t_k - t_j) - \sum_{i=j+1}^k \eta_i \left\| [\mathbf{C}_i \ \mathbf{D}_i] - [\mathbf{0} \ \mathbf{I}_d] \right\|_{\text{op}} - \frac{1}{4}(t_k - t_\ell) \\ &\geq \frac{1}{2}(t_k - t_j) - \frac{1}{4}(t_k - t_\ell). \end{aligned}$$

Take m' such that $t_k - t_{m'} \geq 2(t_k - t_\ell)/3$ and $t_{m'} - t_\ell \geq (t_k - t_\ell)/4$, then

$$\begin{aligned} \sum_{j=\ell+1}^k \eta_j |\mathbf{R}_j \mathbf{w}|^2 &\geq \sum_{j=\ell+1}^{m'} \eta_j |\mathbf{R}_j \mathbf{w}|^2 \\ \text{(B.56)} \quad &\geq \left[\frac{1}{2}(t_k - t_{m'}) - \frac{1}{4}(t_k - t_\ell) \right]^2 \sum_{j=\ell+1}^{m'} \eta_j \geq \frac{(t_k - t_\ell)^3}{576}. \end{aligned}$$

Combining (B.54), (B.55) and (B.56) leads to

$$2[1 + \gamma^2(t_k - t_\ell)^2] \sum_{j=\ell+1}^k \eta_j \left| [\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_j \mathbf{w} \right|^2 \geq \frac{(t_k - t_\ell)^3}{576}.$$

This, together with (B.49) and $t_k - t_\ell \leq 1/[5(B_1 + \gamma + 2)]$, immediately implies

$$\text{(B.57)} \quad \mathbf{w}^\top \boldsymbol{\Gamma}(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}}) \mathbf{w} = \beta^2 \sum_{j=1}^k \eta_j \left| [\mathbf{I}_d \ \mathbf{0}] \mathbf{S}_j \mathbf{w} \right|^2 \geq C(t_k - t_\ell)^3, \quad \text{on } (E_{\ell,k}^1)^c.$$

Consequently, (B.53) and (B.57) yield the desired,

$$\lambda_{\min} \left(\boldsymbol{\Gamma}(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}}) \right) = \min_{|\mathbf{w}|=1} \mathbf{w}^\top \boldsymbol{\Gamma}(\tilde{\mathbf{Y}}_{t_\ell, t_k}^{\mathbf{z}}) \mathbf{w} \geq C(t_k - t_\ell)^3, \quad \text{on } (E_{\ell,k}^1 \cup E_{m,k}^1)^c.$$

The proof is complete. \square

(A. Guillin) LABORATOIRE DE MATHÉMATIQUES BLAISE PASCAL, CNRS-UMR 6620, UNIVERSITÉ CLERMONT-AUVERGNE, AVENUE DES LANDAIS, 63177, AUBIERE, CEDEX

E-mail address: `arnaud.guillin@uca.fr`

(Y. Wang) 1. DEPARTMENT OF MATHEMATICS, FACULTY OF SCIENCE AND TECHNOLOGY, UNIVERSITY OF MACAU, TAIPA, MACAU, 999078, CHINA; 2. UM ZHUHAI RESEARCH INSTITUTE, ZHUHAI, GUANGDONG, 519000, CHINA.

E-mail address: `yc17447@um.edu.mo`

(L. Xu) DEPARTMENT OF MATHEMATICS, FACULTY OF SCIENCE AND TECHNOLOGY, UNIVERSITY OF MACAU, MACAU S.A.R., CHINA.

E-mail address: `lihuxu@um.edu.mo`

(H. Yang) 1. SCHOOL OF MATHEMATICAL SCIENCES, PEKING UNIVERSITY, BEIJING, CHINA. 2. BEIJING INTERNATIONAL CENTER FOR MATHEMATICAL RESEARCH (BICMR), PEKING UNIVERSITY, BEIJING, CHINA.

E-mail address: `yanghr@pku.edu.cn`