# Quantitative Propagation of Chaos for SGD in Wide Neural Networks

**Valentin De Bortoli**
Centre Borelli
ENS Paris Saclay

**Alain Durmus**
Centre Borelli
ENS Paris Saclay

**Xavier Fontaine**
Centre Borelli
ENS Paris Saclay

**Umut Şimşekli**
LTCI, Télécom Paris
Institut Polytechnique de Paris

## Abstract

In this paper, we investigate the limiting behavior of a continuous-time counterpart of the Stochastic Gradient Descent (SGD) algorithm applied to two-layer overparameterized neural networks, as the number or neurons (*i.e.*, the size of the hidden layer) $N \to +\infty$. Following a probabilistic approach, we show 'propagation of chaos' for the particle system defined by this continuous-time dynamics under different scenarios, indicating that the statistical interaction between the particles asymptotically vanishes. In particular, we establish quantitative convergence with respect to $N$ of any particle to a solution of a mean-field McKean-Vlasov equation in the metric space endowed with the Wasserstein distance. In comparison to previous works on the subject, we consider settings in which the sequence of stepsizes in SGD can potentially depend on the number of neurons and the iterations. We then identify two regimes under which different mean-field limits are obtained, one of them corresponding to an implicitly regularized version of the minimization problem at hand. We perform various experiments on real datasets to validate our theoretical results, assessing the existence of these two regimes on classification problems and illustrating our convergence results.

## 1 Introduction

Due to their ability to tackle very challenging problems, neural networks have been extremely popular and keystones in machine learning [1]. Thanks to their practical success, they have become the de facto tool in many application domains, such as image processing [2] and natural language processing [3]. However, the mathematical understanding of these models and their inherent inference mechanism still remains limited.

Among others, one suprising empirical observation about modern neural networks is that increasing the number of neurons in a network often leads to better classification testing and training errors [4], contradicting the classical statistical learning theory [5]. These experimental results suggest that neural network-based methods exhibit a limiting behavior when the number of neurons is large, *i.e.*, when the neural network is *overparameterized*.

In this paper, we contribute to the recent literature on the theoretical analysis of this phenomenon. To this end, we consider a simple two-layer (*i.e.*, one hidden layer) neural network that is parametrized by $N$ weights $w^{1:N} = \{w^{k,N}\}_{k=1}^N$ and trained to minimize the structural risk $\mathscr{R}^N$ by Stochastic Gradient Descent (SGD) using independent and identically distributed (i.i.d.) samples $(X_i, Y_i)_{i\in\mathbb{N}^\star}$. Even in such a simplified setting, the landscape of $\mathscr{R}^N$ is in many cases arduous to be explored, since $\mathscr{R}^N$ is non-convex and might exhibit many local minima and saddle points [6,7]; hence making the minimization of $\mathscr{R}^N$ challenging. However, for large $N$, the analysis of the landscape of $\mathscr{R}^N$ turns out to be much simpler in some situations. For instance [8] has shown that local minima are global minima when the activation function is quadratic as soon as $N$ is larger than twice the size of the original dataset. More generally, relying on approximation or random matrix theory, several works

(*e.g.*, [9–19]) establish favorable properties for the landscape of $\mathscr{R}^N$ as $N \to +\infty$, such as absence of saddle points, poor local minima or connected optima. In addition, minimization by SGD in this setting has also proved to be efficient for some models [20, 21].

In this paper we follow an increasingly popular line of research to analyze the behavior of gradient descent-type algorithms (stochastic or deterministic) used for overparameterized models. This approach consists in establishing a 'continuous mean-field limit' for these algorithms as $N \to +\infty$, and has been successively applied in [22–29]. Based on this result, the *qualitative* long-time behavior of SGD applied to overparameterized neural networks can be deduced: these studies all identify an evolution equation on the limiting probability measure which corresponds to a mean-field ordinary differential equation (ODE), *i.e.*, if the initialization is deterministic, then each hidden unit of the network *independently* evolves along the flow of a specific ODE. This implies that, even though the update step is intrinsically stochastic in SGD, the noise completely vanishes in the limit $N \to +\infty$. In this context, two main strategies have been followed to prove convergence of SGD to this mean-field dynamics. The first one is based on gradient flows in Wasserstein spaces [30–32] and the second one is the 'propagation of chaos' phenomenon [33–35], indicating that the statistical interaction between the individual entries of the network asymptotically vanishes. Both approaches are in fact deeply connected, which stems from the duality between probability and partial differential equation theories [36]. We follow in this paper the second approach and establish that propagation of chaos holds for a continuous counterpart of SGD to a solution of a McKean-Vlasov type diffusion [37] as $N \to +\infty$.

The fact that no noise appears in the mean-field limit of SGD obtained in previous work can seem surprising. Aiming to demystify this matter, we study in this paper the case where the stepsize in SGD can depend on the number of neurons. Our main contribution is to identify two mean-field regimes: The first one is the same as the deterministic mean-field limit obtained in the described literature. The second one is a McKean-Vlasov diffusion for which the covariance matrix is non-zero and depends on the properties of the data distribution. To the best of our knowledge, this limiting diffusion has not been reported in the literature and brings interesting insights on the behavior of neural networks in overparameterized settings. Our results suggest that taking large stepsizes in the stochastic optimization procedure corresponds to an *implicit regularization* of the original problem, which can potentially ease the minimization of the structural risk. In addition, in contrast to previous studies, we establish strong quantitative propagation of chaos and we identify the convergence rate of each neuron to its mean-field limit with respect to $N$. Finally we numerically illustrate the existence of these two regimes and the propagation of chaos phenomenon we derive on several classical classification examples on MNIST and CIFAR-10 datasets. In these experiments, the stochastic regime empirically exhibits slightly better generalization properties compared to the deterministic case identified in [22, 23, 28].

## 2 Overparametrized Neural Networks

Consider some feature and label spaces denoted by $\mathsf{X}$ and $\mathsf{Y}$ endowed with $\sigma$-fields $\mathcal{X}$ and $\mathcal{Y}$ respectively. In this paper, we consider a one hidden layer neural network, whose purpose is to classify data from $\mathsf{X}$ with labels in $\mathsf{Y}$. We suppose that the network has $N \in \mathbb{N}^\star$ neurons in the hidden layer whose weights are denoted by $w^{1:N} = \{w^{k,N}\}_{k=1}^N \in (\mathbb{R}^p)^N$. We model the non-linearity by a function $F : \mathbb{R}^p \times \mathsf{X} \to \mathbb{R}$, and consider a loss function $\ell : \mathbb{R} \times \mathsf{Y} \to \mathbb{R}_+$ and a penalty function $V : \mathbb{R}^p \to \mathbb{R}$. Then, the learning problem corresponding to this space of hypothesis consists in minimizing the structural risk

$$\mathscr{R}^N(w^{1:N}) = \int_{\mathsf{X} \times \mathsf{Y}} \ell\left(\frac{1}{N} \sum_{k=1}^N F(w^{k,N}, x), y\right) \mathrm{d}\pi(x, y) + \frac{1}{N} \sum_{k=1}^N V(w^{k,N}), \qquad (1)$$

where $\pi$ is the data distribution on $\mathsf{X} \times \mathsf{Y}$. Note that, in this particular setting, the weights of the second layer are fixed to $(1/N)$. This setting is referred to as "fixed coefficients" in [25, Theorem 1] and is less realistic than the fully-trainable setting. Nevertheless, we believe that this shortcoming can be circumvented upon replacing $F(w^{k,N}, \cdot)$ by $F(u^{k,N}, \cdot)v^{k,N}$ in (1), where $u^{1:N}$ and $v^{1:N}$ are the weights of the hidden and the second layer respectively. However, this raises new theoretical challenges which are left for future work.

Throughout this paper, we consider the following assumptions.

**A1.** *There exist measurable functions $\Phi : \mathsf{X} \to [1, +\infty)$ and $\Psi : \mathsf{Y} \to [1, +\infty)$ such that the following conditions hold.*

*(a)* $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}_+$ *is such that for any* $y \in \mathsf{Y}$, $(\tilde{\mathrm{y}} \mapsto \ell(\tilde{\mathrm{y}}, y))$ *is three-times differentiable and for any* $\mathrm{y} \in \mathbb{R}$ *and* $y \in \mathsf{Y}$ *we have*

$$|\partial_1 \ell(0, y)| \leq \Psi(y) , \qquad \left|\partial_1^2 \ell(\mathrm{y}, y)\right| + \left|\partial_1^3 \ell(\mathrm{y}, y)\right| \leq \Psi(y) ,$$

*where for any* $i \in \{1, 2, 3\}$, $\partial_1^i \ell(\mathrm{y}, y)$ *is the* $i$-*th derivative of* $(\tilde{\mathrm{y}} \mapsto \ell(\tilde{\mathrm{y}}, y))$ *at* $\mathrm{y}$.

*(b)* $F : \mathbb{R}^p \times \mathsf{X} \to \mathbb{R}$ *is such that for any* $x \in \mathsf{X}$, $(\tilde{w} \mapsto F(\tilde{w}, x))$ *is three-times differentiable and for any* $w \in \mathbb{R}^p$ *and* $x \in \mathsf{X}$

$$\|F(w, x)\| + \left\|\mathrm{D}_w^1 F(w, x)\right\| + \left\|\mathrm{D}_w^2 F(w, x)\right\| + \left\|\mathrm{D}_w^3 F(w, x)\right\| \leq \Phi(x) ,$$

*where for any* $i \in \{1, 2, 3\}$, $\mathrm{D}_w^i F(w, x)$ *is the* $i$-*th differential of* $(\tilde{w} \mapsto F(\tilde{w}, x))$ *at* $w$.

*(c)* $V \in \mathrm{C}^3(\mathbb{R}^p, \mathbb{R})$ *satisfies* $\sup_{w \in \mathbb{R}^p}\{\|\mathrm{D}^2 V(w)\| + \|\mathrm{D}^3 V(w)\|\} < +\infty$.

*(d) The data distribution* $\pi$ *satisfies* $\int_{\mathsf{X} \times \mathsf{Y}}\{\Phi^{10}(x) + \Psi^4(y)\}\mathrm{d}\pi(x, y) < \infty$ .

Note that **A1**-(d) is immediately satisfied in the case where $\pi$ is compactly supported, $\mathsf{X}$ and $\mathsf{Y}$ are subsets of $\mathbb{R}^d$ and $\mathbb{R}$ respectively and $\Psi$ and $\Phi$ are bounded on the support of $\pi$. For any $N \in \mathbb{N}^\star$, under **A1**, by the Lebesgue dominated convergence theorem, $\mathscr{R}^N$ given by (1) is well-defined, continuously differentiable with gradient given for any $w^{1:N} \in (\mathbb{R}^p)^N$ by

$$\nabla \mathscr{R}^N(w^{1:N}) = \int_{\mathsf{X} \times \mathsf{Y}} \nabla_w \hat{\mathscr{R}}^N(w^{1:N}, x, y)\mathrm{d}\pi(x, y) ,$$

$$\hat{\mathscr{R}}^N(w^{1:N}, x, y) = \ell\left(\tfrac{1}{N}\textstyle\sum_{k=1}^N F(w^{k,N}, x), y\right) + \tfrac{1}{N}\textstyle\sum_{k=1}^N V(w^{k,N}) ,$$

$$N \nabla_w \hat{\mathscr{R}}^N(w^{1:N}, x, y) = \partial_1 \ell\left(\tfrac{1}{N}\textstyle\sum_{k=1}^N F(w^{k,N}, x), y\right) \nabla_w F^{1:N}(w^{1:N}, x) + \nabla V^{1:N}(w^{1:N}) ,$$

setting $\nabla_w F^{1:N}(w^{1:N}, x) = \{\nabla_w F(w^{k,N}, x)\}_{k=1}^N$, and $\nabla V^{1:N}(w^{1:N}) = \{\nabla V(w^{k,N})\}_{k=1}^N$.

Let $(W_0^k)_{k \in \mathbb{N}^\star}$ be i.i.d. $p$ dimensional random variables with distribution $\mu_0$. Consider the sequence $(W_n^{1:N})_{n \in \mathbb{N}}$ associated with SGD, starting from $W_0^{1:N}$ and defined by the following recursion: for any $n \in \mathbb{N}$ denoting the iteration index

$$W_{n+1}^{1:N} = W_n^{1:N} - \gamma N^\beta(n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha}\nabla \hat{\mathscr{R}}^N(W_n^{1:N}, X_n, Y_n) , \tag{2}$$

where $(X_n, Y_n)_{n \in \mathbb{N}}$ is a sequence of i.i.d. input/label samples distributed according to $\pi$, and $(\gamma N^\beta(n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha})_{n \in \mathbb{N}}$ as a whole denotes a sequence of stepsizes: here, $\beta \in [0, 1]$, $\alpha \in [0, 1)$, and $\gamma_{\alpha,\beta}(N) = \gamma^{1/(1-\alpha)}N^{(\beta-1)/(1-\alpha)}$. Note that in the constant stepsize setting $\alpha = 0$, the recursion (2) consists in using $\gamma N^\beta$ as a stepsize In addition, it also encompasses the case of decreasing stepsizes (as soon as $\alpha > 0$). The term $\gamma_{\alpha,\beta}(N)^{-1}$ in (2) is a scaling parameter which appears naturally in the corresponding continuous-time dynamics, see (4) below. We stress that contrary to previous approaches such as [22, 23, 28], the stepsize appearing in (2) depends on the number of neurons $N$. Our main contribution is to establish that different mean-field limiting behaviors of a continuous counterpart of SGD arise depending on $\beta$.

We will show that the quantity $\gamma_{\alpha,\beta}(N)$ plays the role of a discretization stepsize in the McKean-Vlasov approximation of SGD. The case where $\alpha = 0$ and $\beta = 0$, *i.e.*, the setting considered by [22, 23, 28], corresponds to choosing the stepsize as $\gamma/N$, which decreases with increasing $N$. In the new setting $\alpha = 0$, $\beta = 1$, this corresponds to take a fixed stepsize $\gamma$. This observation further motivates the scaling and the parameter we introduced in (2).

Before stating our result, we present and give an informal derivation of the continuous particle system dynamics we consider to model (2). We first show that (2) can be rewritten as a recursion corresponding to the discretization of a continuous particle system, *i.e.*, a stochastic differential equation (SDE) with coefficients depending on the empirical measure of the particles. Let us denote by $\mathscr{P}(\mathsf{E})$ the set of probability measures on a measurable space $(\mathsf{E}, \mathcal{E})$. Remark that for each particle dynamics $(W_n^{k,N})_{n \in \mathbb{N}}$ the SGD update (2) is a function of the current position and the empirical

measure of the weights. To show this, define the mean-field $h : \mathbb{R}^p \times \mathscr{P}(\mathbb{R}^d) \to \mathbb{R}^p$ and the noise field $\xi : \mathbb{R}^p \times \mathscr{P}(\mathbb{R}^d) \times \mathsf{X} \times \mathsf{Y} \to \mathbb{R}^p$, for any $\mu \in \mathscr{P}(\mathbb{R}^p)$, $w \in \mathbb{R}^p$, $(x, y) \in \mathsf{X} \times \mathsf{Y}$ by

$$h(w, \mu) = - \int_{\mathsf{X} \times \mathsf{Y}} \partial_1 \ell \left( \mu[F(\cdot, x)], y \right) \nabla_w F(w, x) \, \mathrm{d}\pi(x, y) - \nabla V(w) \, ,$$

$$\xi(w, \mu, x, y) = -h(w, \mu) - \partial_1 \ell(\mu[F(\cdot, x)], y) \nabla_w F(w, x) - \nabla V(w) \, .$$

Note that with this notation, $h(w^{k,N}, \nu_n^N) = -N \partial_{w^{k,N}} \mathscr{R}^N(w^{1:N})$ and $\xi(w^{k,N}, \nu_n^N, X_n, Y_n) = N\{-\partial_{w^{k,N}} \hat{\mathscr{R}}^N(w^{1:N}, X_n, Y_n) + \partial_{w^{k,N}} \mathscr{R}^N(w^{1:N})\}$, for any $N \in \mathbb{N}$, $k \in \{1, \dots, N\}$ and $n \in \mathbb{N}$, where $\nu^N$ is the empirical measure of the discrete particle system corresponding to SGD defined by $\nu_n^N = N^{-1} \sum_{k=1}^N \delta_{W_n^{k,N}}$. Then, the recursion (2) can be rewritten as follows:

$$W_{n+1}^{k,N} = W_n^{k,N} + \gamma N^{\beta-1}(n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha} \left\{ h(W_n^{k,N}, \nu_n^N) + \xi(W_n^{k,N}, \nu_n^N, X_n, Y_n) \right\} \, . \quad (3)$$

We now present the continuous model associated with this discrete process. For large $N$ or small $\gamma$ these two processes can be arbitrarily close. For $N \in \mathbb{N}^\star$, consider the particle system diffusion $(\mathbf{W}_t^{1:N})_{t \geq 0} = (\{\mathbf{W}_t^{k,N}\}_{k=1}^N)_{t \geq 0}$ starting from $\mathbf{W}_0^{1:N} = W_0^{1:N}$ defined for any $k \in \{1, \dots, N\}$ by

$$\mathrm{d}\mathbf{W}_t^{k,N} = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N) \mathrm{d}t + \gamma_{\alpha,\beta}(N)^{1/2} \Sigma^{1/2}(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N) \mathrm{d}\mathbf{B}_t^k \right\} \, , \quad (4)$$

where $\{(\mathbf{B}_t^k)_{t \geq 0} \ : \ k \in \mathbb{N}^\star\}$ is a family of independent $p$-dimensional Brownian motions and $\boldsymbol{\nu}_t^N$ is the empirical probability distribution of the particles defined for any $t \geq 0$ by $\boldsymbol{\nu}_t^N = N^{-1} \sum_{k=1}^N \delta_{\mathbf{W}_t^{k,N}}$. In addition in (4), $\Sigma$ is the $p \times p$ matrix given by

$$\Sigma(w, \mu) = \int_{\mathsf{X} \times \mathsf{Y}} \xi(w, \mu, x, y) \xi(w, \mu, x, y)^\top \mathrm{d}\pi(x, y) \, , \qquad \text{for any } w \in \mathbb{R}^p \text{ and } \mu \in \mathscr{P}(\mathbb{R}^p) \, ,$$

which is well-defined under **A**1. In the supplementary material we show that under **A**1, (4) admits a unique strong solution. We now give an informal discussion to justify why (4) can be seen as the continuous-time counterpart of (3). For any $N \in \mathbb{N}^\star$, define $(\tilde{\mathbf{W}}_t^{1:N})_{t \geq 0}$ for any $t \geq 0$ by $\tilde{\mathbf{W}}_t^{1:N} = W_{n_t}^{1:N}$ with $n_t = \lfloor t/\gamma_{\alpha,\beta}(N) \rfloor$ and denote $\tilde{\nu}_t^N$ the empirical measure associated with $\tilde{\mathbf{W}}_t^{1:N}$. In this case, by defining the interval $I_{n,\alpha,\beta}^N = [n\gamma_{\alpha,\beta}(N), (n+1)\gamma_{\alpha,\beta}(N)]$ and using (3) and $\gamma_{\alpha,\beta}(N)^{1-\alpha} = \gamma N^{\beta-1}$, we obtain the following approximation for any $n \in \mathbb{N}$

$$\tilde{\mathbf{W}}_{(n+1)\gamma_{\alpha,\beta}(N)}^{k,N} - \tilde{\mathbf{W}}_{n\gamma_{\alpha,\beta}(N)}^{k,N}$$

$$= \gamma N^{\beta-1}(n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha} \left\{ h(\tilde{\mathbf{W}}_{n\gamma_{\alpha,\beta}(N)}^{k,N}, \nu_n^N) + \xi(\tilde{\mathbf{W}}_{n\gamma_{\alpha,\beta}(N)}^{k,N}, \nu_n^N, X_n, Y_n) \right\}$$

$$\approx \gamma_{\alpha,\beta}(N)(n\gamma_{\alpha,\beta}(N) + 1)^{-\alpha} \left\{ h(\tilde{\mathbf{W}}_{n\gamma_{\alpha,\beta}(N)}^{k,N}, \tilde{\boldsymbol{\nu}}_{n\gamma_{\alpha,\beta}(N)}^N) + \Sigma^{1/2}(\tilde{\mathbf{W}}_{n\gamma_{\alpha,\beta}(N)}^{k,N}, \tilde{\boldsymbol{\nu}}_{n\gamma_{\alpha,\beta}(N)}^N)G \right\}$$

$$\approx \underbrace{\int_{I_{n,\alpha,\beta}^N} (s+1)^{-\alpha} h(\tilde{\mathbf{W}}_s^{k,N}, \tilde{\boldsymbol{\nu}}_s^N) \mathrm{d}s}_{(A)} + \underbrace{\int_{I_{n,\alpha,\beta}^N} \gamma_{\alpha,\beta}^{1/2}(N)(s+1)^{-\alpha} \Sigma^{1/2}(\tilde{\mathbf{W}}_s^{k,N}, \tilde{\boldsymbol{\nu}}_s^N) \mathrm{d}\mathbf{B}_s^k}_{(B)} \, , \quad (5)$$

where $G$ is a $p$-dimensional Gaussian random variable with zero mean and identity covariance matrix. Note that the second line corresponds to (3) and the last to (4). To obtain such proxy, we first remark that for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$, $\xi(w, \mu, X_n, Y_n)$ has zero mean and covariance matrix $\Sigma(w, \mu)$ and assume that the noise term is roughly Gaussian. Second, we use that the covariance of $(B)$ in (5) is equal to $\int_{I_{n,\alpha,\beta}^N} \gamma_{\alpha,\beta}(N)(s+1)^{-2\alpha} \Sigma(\tilde{\mathbf{W}}_s^{k,N}, \tilde{\boldsymbol{\nu}}_s^N) \mathrm{d}s$. To obtain the last line, we use some first-order Taylor expansion of this term and $(A)$ as $\gamma_{\alpha,\beta}(N) \to 0$. Then, (5) corresponds to (4) on $I_{n,\alpha,\beta}^N$. As a result, (4) is the continuous counterpart to (3) and $n$ iterations in (3) correspond to the horizon time $n\gamma_{\alpha,\beta}(N)$ in (4). In the next section, we show that a strong quantitative propagation of chaos holds for (4) *i.e.*, we show that for $N \to +\infty$ the particles become independent and have the same distribution associated with a McKean-Vlasov diffusion. The extension of these results to discrete SGD (3) and the rigorous derivation of (5) can be established using strong functional approximations following [38, Proposition 1]. Due to space constraints, we leave it as future work.

Finally, note that until now we only considered the case where the batch size in SGD is equal to one. For a batch size $M \in \mathbb{N}^\star$, this limitation can be lifted replacing $\pi$ and $\hat{\mathscr{R}}^N$ in (2) by $\pi^{\otimes M}$ and

$$\hat{\mathscr{R}}^{N,M}(w^{1:N}, x, y) = \frac{1}{M} \sum_{i=1}^M \ell \left( \frac{1}{N} \sum_{k=1}^N F(w^{k,N}, x_i), y_i \right) \, ,$$

defined for any $w^{1:N} \in (\mathbb{R}^p)^N$, $x = (x_i)_{i \in \{1,\dots,M\}} \in \mathsf{X}^M$ and $y = (y_i)_{i \in \{1,\dots,M\}} \in \mathsf{Y}^M$. In this case, we obtain that the continuous-time counterpart of (3) is given by (4) upon replacing $\Sigma^{1/2}$ by $\Sigma^{1/2}/M^{1/2}$. This leads to the particle system diffusion $(\mathbf{W}_t^{1:N})_{t \geq 0} = (\{\mathbf{W}_t^{k,N}\}_{k=1}^N)_{t \geq 0}$ starting from $\mathbf{W}_0^{1:N}$ defined for any $k \in \{1,\dots,N\}$ by

$$\mathrm{d}\mathbf{W}_t^{k,N} = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N)\mathrm{d}t + (\gamma_{\alpha,\beta}(N)/M)^{1/2}\Sigma^{1/2}(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N)\mathrm{d}\mathbf{B}_t^k \right\} . \quad (6)$$

In the supplement Section 7, we also present the case of a modified Stochastic Gradient Langevin Dynamics (mSGLD) algorithm [39] which was considered in [23] in the specific case $\beta = 0$. We extend our propagation of chaos results to this setting.

## 3 Mean-Field Approximation and Propagation of Chaos

In this section we identify the mean-field limit of the diffusion (6). More precisely, we show that there exist two regimes depending on how the stepsize scale with the number of hidden units.

Our results are based on the propagation of chaos theory [33–35] and extend the recent works of [22,23, 25,27,28,40]. In what follows, we denote $\mathscr{P}_2(\mathbb{R}^p) = \{\mu \in \mathscr{P}(\mathbb{R}^p) : \int_{\mathbb{R}^p} \|\tilde{w}\|^2 \mathrm{d}\mu(\tilde{w}) < +\infty\}$ and $\mathrm{C}(\mathbb{R}_+, \mathbb{R}^p)$ the set of continuous functions from $\mathbb{R}_+$ to $\mathbb{R}^p$. We also consider the usual metric m on $\mathrm{C}(\mathbb{R}_+, \mathbb{R}^p)$ defined for any $u_1, u_2 \in \mathrm{C}(\mathbb{R}_+, \mathbb{R}^p)$ by $\mathrm{m}(u_1, u_2) = \sum_{n \in \mathbb{N}^\star} 2^{-n} \|u_1 - u_2\|_{\infty,n} /\{1 + \|u_1 - u_2\|_{\infty,n}\}$, where $\|u_1 - u_2\|_{\infty,n} = \sup_{t \in [0,n]} \|u_1(t) - u_2(t)\|$. It is well-known that $(\mathscr{C}, \mathrm{m}) = (\mathrm{C}(\mathbb{R}_+, \mathbb{R}^p), \mathrm{m})$ is a complete separable space. For any metric space $(\mathsf{F}, \mathrm{m}_\mathsf{F})$, with Borel $\sigma$-field $\mathcal{B}(\mathsf{F})$, we define the extended Wasserstein distance of order 2, denoted $\mathcal{W}_2 : \mathscr{P}(\mathsf{F}) \times \mathscr{P}(\mathsf{F}) \to [0, +\infty]$ for any $\mu_1, \mu_2 \in \mathscr{P}(\mathsf{F})$ by $\mathcal{W}_2^2(\mu_1, \mu_2) = \inf_{\Lambda \in \Gamma(\mu_1, \mu_2)} \int_{\mathsf{F} \times \mathsf{F}} \mathrm{m}_\mathsf{F}^2(v_1, v_2)\mathrm{d}\Lambda(v_1, v_2)$, where $\Gamma(\mu_1, \mu_2)$ is the set of transference plans between $\mu_1$ and $\mu_2$, i.e., $\Lambda \in \Gamma(\mu_1, \mu_2)$ if for any $\mathsf{A} \in \mathcal{B}(\mathsf{F})$, $\Lambda(\mathsf{A} \times \mathsf{F}) = \mu_1(\mathsf{A})$ and $\Lambda(\mathsf{F} \times \mathsf{A}) = \mu_2(\mathsf{A})$.

We start by stating our results in the case where $\beta \in [0, 1)$ for which a deterministic mean-field limit is obtained. Consider the mean-field ODE starting from a random variable $\mathbf{W}_0^\star$ given by

$$\mathrm{d}\mathbf{W}_t^\star = (t+1)^{-\alpha}h(\mathbf{W}_t^\star, \boldsymbol{\lambda}_t^\star)\mathrm{d}t , \qquad \text{with } \boldsymbol{\lambda}_t^\star \text{ the distribution of } \mathbf{W}_t^\star . \quad (7)$$

We show in the supplement that this ODE admits a solution on $\mathbb{R}_+$. This mean-field equation (7) is deterministic conditionally to its initialization.

**Theorem 1.** *Assume* **A**1. *Let* $(\mathbf{W}_0^k)_{k \in \mathbb{N}}$ *be a sequence of i.i.d.* $\mathbb{R}^p$-*valued random variables with distribution* $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$ *and set for any* $N \in \mathbb{N}^\star$, $\mathbf{W}_0^{1:N} = \{\mathbf{W}_0^k\}_{k=1}^N$. *Then, for any* $m \in \mathbb{N}^\star$ *and* $T \geq 0$, *there exists* $C_{m,T} \geq 0$ *such that for any* $\alpha \in [0, 1)$, $\beta \in [0, 1)$, $M \in \mathbb{N}^\star$ *and* $N \in \mathbb{N}^\star$ *with* $N \geq m$

$$\mathbb{E}\left[\sup_{t \in [0,T]} \|\mathbf{W}_t^{1:m,N} - \mathbf{W}_t^{1:m,\star}\|^2\right] \leq C_{m,T} \left\{ N^{-(1-\beta)/(1-\alpha)}M^{-1} + N^{-1} \right\} ,$$

*with* $(\mathbf{W}_t^{1:m,N}, \mathbf{W}_t^{1:m,\star}) = \{(\mathbf{W}_t^{k,N}, \mathbf{W}_t^{k,\star})\}_{k=1}^m$, $(\mathbf{W}_t^{1:N})$ *the solution of* (6) *starting from* $\mathbf{W}_0^{1:N}$, *and for any* $k \in \mathbb{N}^\star$, $\mathbf{W}_t^{k,\star}$ *the solution of* (7) *starting from* $\mathbf{W}_0^k$.

In Theorem 1, $m$ is a fixed number of particles. Note that $\{(\mathbf{W}_t^{k,\star})_{t \geq 0} : k \in \mathbb{N}^\star\}$ is i.i.d. with distribution $\boldsymbol{\lambda}^\star$ which is the pushfoward measure of $\mu_0$ by the function $(w_0 \mapsto (w_t)_{t \geq 0})$ which from an initial point $w_0$ gives $(w_t)_{t \geq 0} \in \mathscr{C}$ the solution of (7) on $\mathbb{R}_+$. Theorem 1 shows that the dynamics of the particles become deterministic and independent when $N \to +\infty$. The proofs of Theorem 1 and the following result, Theorem 2, are postponed to Section 9.4.

We now consider the case $\beta = 1$ and derive a similar quantitative theorem as Theorem 1 but with a different dynamics than (7). Consider the mean-field SDE starting from variable $\mathbf{W}_0^\star$ given by

$$\mathrm{d}\mathbf{W}_t^\star = (t+1)^{-\alpha}\left\{ h(\mathbf{W}_t^\star, \boldsymbol{\lambda}_t^\star)\mathrm{d}t + (\gamma^{1/(1-\alpha)}\Sigma(\mathbf{W}_t^\star, \boldsymbol{\lambda}_t^\star)/M)^{1/2}\mathrm{d}\mathbf{B}_t \right\} , \quad (8)$$

where $\boldsymbol{\lambda}_t^\star$ is the distribution of $\mathbf{W}_t^\star$ and $(\mathbf{B}_t)_{t \geq 0}$ is a $p$ dimensional Brownian motion. Note that taking the limit $\gamma \to 0$ or $M \to +\infty$ in (8) we recover (7). We show in the supplement that this SDE admits a solution on $\mathbb{R}_+$. The following theorem is similar to Theorem 1 in the case $\beta = 1$.

**Theorem 2.** *Let $\beta = 1$. Assume **A**1. Let $(\mathbf{W}_0^k)_{k\in\mathbb{N}}$ be a sequence of $\mathbb{R}^p$-valued random variables with distribution $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$ and assume that for any $N \in \mathbb{N}^\star$, $\mathbf{W}_0^{1:N} = \{\mathbf{W}_0^k\}_{k=1}^N$. Then, for any $m \in \mathbb{N}^\star$ and $T \geq 0$, there exists $C_{m,T} \geq 0$ such that for any $\alpha \in [0,1)$, $M \in \mathbb{N}^\star$ and $N \in \mathbb{N}^\star$ with $N \geq m$*

$$\mathbb{E}\left[\sup_{t\in[0,T]} \|\mathbf{W}_t^{1:m,N} - \mathbf{W}_t^{1:m,\star}\|^2\right] \leq C_{m,T} N^{-1} \,,$$

*with $(\mathbf{W}_t^{1:m,N}, \mathbf{W}_t^{1:m,\star}) = \{(\mathbf{W}_t^{k,N}, \mathbf{W}_t^{k,\star})\}_{k=1}^m$, $(\mathbf{W}_t^{1:N})$ the solution of (6) starting from $\mathbf{W}_0^{1:N}$, and for any $k \in \mathbb{N}^\star$, $\mathbf{W}_t^{k,\star}$ the solution of (8) starting from $\mathbf{W}_0^k$ and Brownian motion $(\mathbf{B}_t^k)_{t\geq 0}$.*

The main difference between (7) and (8) is that now this mean-field limit is now longer deterministic up to its initialization but is a SDE driven by a Brownian motion. The stochastic nature of SGD is preserved in this second regime. (7) corresponds to some implicit regularization of (8). In the case where for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$, $\Sigma(w,\mu) = \sigma^2 \mathrm{Id}$ with $\sigma > 0$, it can shown that $(\lambda_t^\star)_{t\geq 0}$ is a gradient flow for an entropic-regularized functional. This relation between our approach and the gradient flow perspective is investigated in the supplement Section 11.

Denote for any $N \in \mathbb{N}^\star$ and $m \in \{1,\ldots,N\}$, $\lambda^{1:m,N}$ the distribution on $\mathscr{C}$ of $\{(\mathbf{W}_t^{k,N})_{t\geq 0}\}_{k=1}^m$. Recall that $\{\mathbf{W}_0^{k,N}\}_{k=1}^N$ are $N$ i.i.d. $\mathbb{R}^p$-valued random variables with distribution $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$. As an immediate consequence of Theorem 1, Theorem 2 and the definition of $\mathcal{W}_2$ for the distance m on $\mathscr{C}$, we have the following propagation of chaos result.

**Corollary 3.** *Assume **A**1. Then for any $\beta \in [0,1]$, $\alpha \in [0,1)$, $M \in \mathbb{N}^\star$ and $m \in \mathbb{N}$ we have $\lim_{N\to+\infty} \mathcal{W}_2(\lambda^{1:m,N}, (\lambda^\star)^{\otimes m}) = 0$ where $\lambda^\star$ is the distribution of $(\mathbf{W}_t^\star)_{t\geq 0}$ solution of (7) if $\beta \in (0,1]$ and (8) if $\beta = 1$ with $\mathbf{W}_0^\star$ distributed according to $\mu_0$.*

Corollary 3 has two main consequences: when the number of hidden units is large (i) all the units have the same distribution $\lambda^\star$, and (ii) the units are independent. Note also that this corollary is valid for the whole trajectory and not only for a fixed time horizon.

Finally, we derive similar results to Corollary 3 for the sequence of the empirical measures. Let $(\boldsymbol{\nu}^N)_{N\in\mathbb{N}^\star}$ be the sequence of empirical measures associated with (6) and given by $\boldsymbol{\nu}^N = N^{-1}\sum_{k=1}^N \delta_{(\mathbf{W}_t^{k,N})_{t\geq 0}}$. Note that for any $N \in \mathbb{N}$, $\boldsymbol{\nu}^N$ is a random probability measure on $\mathscr{P}(\mathscr{C})$. Denote for any $N \in \mathbb{N}^\star$, $\boldsymbol{\Upsilon}^N$ its distribution which then belongs to $\mathscr{P}(\mathscr{P}(\mathscr{C}))$. Since the convergence with respect to the $\mathcal{W}_2$ distance implies the weak convergence, using Corollary 3 and the Tanaka-Sznitman theorem [33, Proposition 2.2], we get that $(\boldsymbol{\Upsilon}^N)_{N\in\mathbb{N}^\star}$ weakly converges towards $\delta_{\lambda^\star}$. In fact, we prove the following stronger proposition whose proof is postponed to Section 9.3.

**Proposition 4.** *Assume **A**1. Then, for any $\beta \in [0,1)$, $\alpha \in [0,1)$ and $M \in \mathbb{N}^\star$ we have $\lim_{N\to+\infty} \mathcal{W}_2(\boldsymbol{\Upsilon}^N, \delta_{\lambda^\star}) = 0$, where $\lambda^\star$ is the distribution of $(\mathbf{W}_t^\star)_{t\geq 0}$ solution of (7) if $\beta \in (0,1]$ and (8) if $\beta = 1$ with $\mathbf{W}_0^\star$ distributed according to $\mu_0$.*

*Proof of Proposition 4.* We consider only the case $\beta = 1$, the proof for $\beta \in [0,1)$ following the same lines. Let $M \in \mathbb{N}^\star$. We have for any $N \in \mathbb{N}^\star$ using Proposition 6,

$$\mathcal{W}_2(\boldsymbol{\Upsilon}^N, \delta_{\lambda^\star})^2 \leq \mathbb{E}[\mathcal{W}_2(\boldsymbol{\nu}^N, \lambda^\star)^2] \leq N^{-1}\sum_{k=1}^N \mathbb{E}[\mathrm{m}^2((\mathbf{W}_t^{k,N})_{t\geq 0}, (\mathbf{W}_t^{k,\star})_{t\geq 0})] \,. \qquad (9)$$

Let $\varepsilon > 0$ and $n_0 \in \mathbb{N}^\star$ such that $\sum_{n=n_0+1}^{+\infty} 2^{-n} \leq \varepsilon$. Combining (40), Theorem 2 and the Cauchy-Schwarz inequality we get that for any $N \in \mathbb{N}^\star$

$$\mathcal{W}_2(\boldsymbol{\Upsilon}^N, \delta_{\lambda^\star})^2 \leq 2\varepsilon^2 + \frac{2n_0}{N}\sum_{k=1}^N\sum_{n=1}^{n_0} \mathbb{E}\left[\sup_{t\in[0,n]} \|\mathbf{W}_t^{k,N} - \mathbf{W}_t^{k,\star}\|^2\right] \leq 2\varepsilon^2 + 2n_0 N^{-1}\sum_{n=0}^{n_0} C_{1,n} \,.$$

Therefore, for any $\varepsilon > 0$ there exists $N_0 \in \mathbb{N}^\star$ such that for any $N \geq N_0$, $\mathcal{W}_2(\boldsymbol{\Upsilon}^N, \delta_{\lambda^\star}) \leq \varepsilon$. $\qquad \square$

**Relation to existing results.** To the authors knowledge, only the case $\beta = 0$ has been considered in the current literature. More precisely, Theorem 1 is a functional and quantitative extension of the results established in [22–24, 28, 40]. First, in [22, Theorem 1.6], it is shown that $(\lambda^{1:m,N})_{N\in\mathbb{N}^\star, N\geq\ell}$ weakly converges towards $(\lambda^\star)^{\otimes m}$. [23, Theorem 3] shows weak convergence of SGD to (7) with high probability in the case $V = 0$ and the quadratic loss $\ell(\mathrm{y}_1, \mathrm{y}_2) = (\mathrm{y}_1 - \mathrm{y}_2)^2$. [40, Theorem 1.5]

establishes a central limit theorem for $(\boldsymbol{\nu}^N)_{N \in \mathbb{N}^\star}$ with rate $N^{-1/2}$ which is in accordance with the convergence rate identified in Theorem 1. Finally, [28, Theorem 2.6] and [24, Proposition 3.2] imply the convergence of $\boldsymbol{\nu}^N$ almost surely under the setting $\Sigma = 0$ in (6) which corresponds to the continuous gradient flow dynamics associated with $\hat{\mathscr{R}}^{N,M}$. We conclude this part by mentioning that similar results are derived for mSGLD in the supplement Section 7 which extend the ones obtained in [25].

Having established the convergence of (6) to (8), we are interested in the long-time behaviour of $(\mathbf{W}_t^\star)_{t \geq 0}$ in the case $\alpha = 0$. To address this problem, the first step is to show that this SDE admits at least one stationary distribution, *i.e.* a probability measure $\mu^\star$ such that if $\mathbf{W}_0^\star$ has distribution $\mu^\star$, then for any $t \geq 0$, $\mathbf{W}_t^\star$ has distribution $\mu^\star$. If $V$ is strongly convex, we are able to answer positively to this question in the case $p = 1$. The proof of this result is postponed to Section 10.

**Proposition 5.** *Assume* **A**1*, $\alpha = 0$ and $p = 1$. In addition, assume that there exist $\eta, \bar{\sigma} > 0$ such that for any $w \in \mathbb{R}$ and $\mu \in \mathscr{P}(\mathbb{R})$, $\Sigma(w, \mu) \geq \bar{\sigma}^2$ and $V$ is $\eta$-strongly convex. Let $H : \mathscr{P}_2(\mathbb{R}) \to \mathscr{P}_2(\mathbb{R})$ defined for any $\mu \in \mathscr{P}_2(\mathbb{R})$ and $w \in \mathbb{R}$ by*

$$(\mathrm{d}H(\mu)/\mathrm{dLeb})(w) \propto \bar{\Sigma}^{-1}(w, \mu) \exp\left[-2 \int_0^w h(\tilde{w}, \mu)/\bar{\Sigma}(\tilde{w}, \mu) \mathrm{d}\tilde{w}\right] ,$$

*where $\bar{\Sigma}(w, \mu) = \gamma^{1/(1-\alpha)} \Sigma(w, \mu)/M$ and* Leb *is the Lebesgue measure on $\mathbb{R}$. Then* $\mathsf{S} = \{\mu \in \mathscr{P}_2(\mathbb{R}) : H(\mu) = \mu\} \neq \emptyset$ *and for any $\mu \in \mathsf{S}$, $\mu$ is invariant for* (8).

## 4  Experiments

We now empirically illustrate the results derived in the previous section. More precisely, we focus on the classification task for two datasets: MNIST [41] and CIFAR-10 [42]. In all of our experiments we consider a fully-connected neural network with one hidden layer and ReLU activation function. We consider the cross-entropy loss in order to train the neural network using SGD as described in Section 2. All along this section we fix a time horizon $T \geq 0$ and sample $W_{n_T}^{1:N}$ defined by (3) with $n_T = \lfloor T/\gamma_{\alpha,\beta}(N) \rfloor$ for $\gamma_{\alpha,\beta}(N) = \gamma^{1/(1-\alpha)} N^{(\beta-1)/(1-\alpha)}$ and taking a batch of size $M \in \mathbb{N}^\star$. We aim at illustrating the results of Section 3 taking $N \to +\infty$ and different sets of values for the parameters $\alpha, \beta, M, \gamma$ in (6). Indeed, recall that as observed in (5), $W_{n_T}^{1:N}$ is an approximation of $\mathbf{W}_T^{1:N}$. See Section 12 for a detailed description of our experimental setting. If not specified, we set $\alpha = 0$, $M = 100$, $T = 100$, $\gamma = 1$.
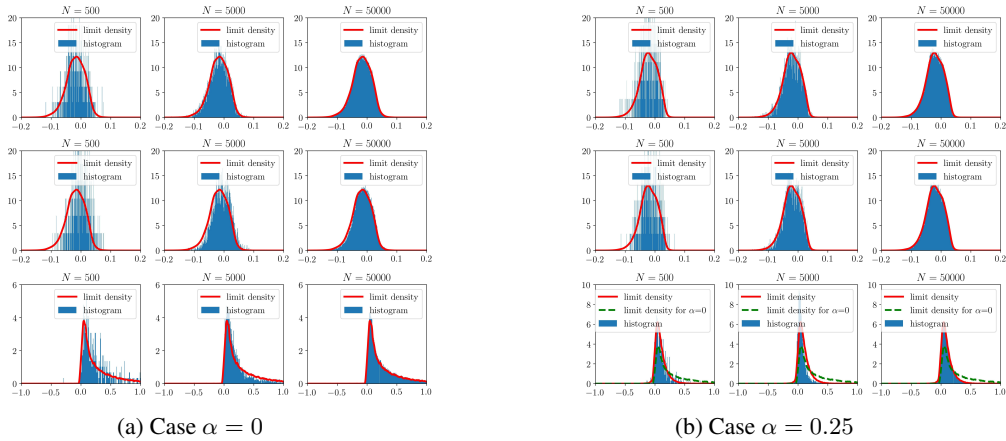


(a) Case $\alpha = 0$                    (b) Case $\alpha = 0.25$

Figure 1: Convergence of the empirical distribution of the weights as $N \to +\infty$.

**Convergence of the empirical measure.**   First we assess the convergence of the empirical distribution of the weights of the hidden layer to a limit distribution when $N \to +\infty$. We focus on the MNIST classification task. Note that in this case $p = 28 \times 28 = 784$. In Figure 1, we observe the behavior of the histograms of the weights $W_{n_T}^{1:N}$ of the hidden layer along the coordinate $(1, 1)$ as $N \to +\infty$. We experimentally observe the existence of two different regimes, one for $\beta < 1$ and the other one for $\beta = 1$. In Figure 1, the first line corresponds to the evolution of the histogram in the case where $\beta = 0.5$. The second and the third lines correspond to the same experiment with

$\beta = 0.75$ and $\beta = 1$, respectively. Note that in both cases the histograms converge to a limit. This limit histogram exhibits two regimes depending if $\beta < 1$ or $\beta = 1$.

**Existence of two regimes.** Now we assess the stochastic nature of the second regime we obtain in the case $\beta = 1$ in contrast to the regime for $\beta < 1$ which is deterministic. In order to highlight this situation, all the weights of the neural network are initialized with a fixed value, *i.e.*, for any $N \in \mathbb{N}^\star$ and $k \in \{1, \dots, N\}$, $W_0^{k,N} = w_0 \in \mathbb{R}^p$. Then, the neural network is trained on the MNIST dataset for $N = 10^6$ and $\beta = 0.75$ or $\beta = 1$. Figure 2 represents 7 samples of the first component of $W_{n_T}^{1:N}$ obtained



Figure 2: Deterministic versus stochastic behavior depending on the value $\beta$.

with independent runs of SGD. We can observe that for $\beta = 0.75$ all the samples converge to the same value which agrees with (7) while in the case where $\beta = 1$ they exhibit different values, which is in accordance with (8).
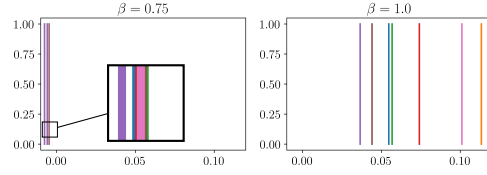
**From stochastic to deterministic.** We illustrate that when $\gamma \to 0$ the dynamics identified in (8) tends to the one identified in (7). We fix $\beta = 1$ and $N = 10000$ and focus on the MNIST classification task. In Figure 3 we show the histogram of the weights $W_{n_T}^{1:N}$ along the coordinate $(1,1)$ for different values of $\gamma$. As expected, see (8) and the follow-
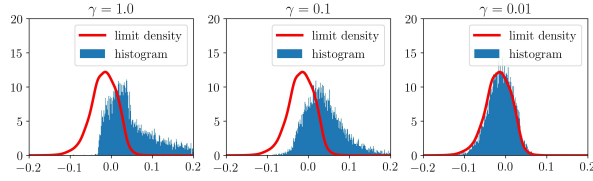


Figure 3: Convergence to the deterministic regime as $\gamma \to 0$.

ing remark, when $\gamma \to 0$ we recover the limit histogram with $\beta < 1$. In Figure 8 we also study the convergence of the empirical measure when $M \to +\infty$ in the case where $\beta < 1$.

**Long-time behavior.** Finally, we illustrate the interest of taking $\beta = 1$ in our setting by considering the more challenging classification task on the CIFAR-10 dataset. We consider the following set of parameters $\alpha = 0$, $M = 100$, $T = 10000$, $\gamma = 0.1$. We emphasize that this experiment aims at comparing the performance of the setting $\beta < 1$ and the one with $\beta = 1$ and that we are not trying to reach state-of-the-art results. In Table 1 we present training and test accuracies for the classification task at hand. To build the classification estimator we average the weights along their trajectory, *i.e.*, we perform averaging and consider the average estimator $\bar{W}_{n_T}^{1:N} = (n_T - n_0 + 1)^{-1} \sum_{n=n_0}^{n_T} W_n^{1:N}$, where $n_0 = 1000$. Using $\beta = 1$ roughly increases the test accuracy by $1\%$, while the training accuracy is not $100\%$. This empirically seems to demonstrate that using a smaller value of $\beta$ tends to overfit the data, whereas using $\beta = 1$ has a regularizing effect.

Table 1: Training and Test accuracies for different settings on the CIFAR-10 dataset, with $\alpha = 0$, $M = 100$ and $\gamma = 0.1$ for $T = 10000$

| Values of $N$ and $\beta$ | $N = 5000$ $\beta = 0.75$ | $N = 5000$ $\beta = 1.0$ | $N = 10000$ $\beta = 0.75$ | $N = 10000$ $\beta = 1.0$ | $N = 50000$ $\beta = 0.75$ | $N = 50000$ $\beta = 1.0$ |
|---|---|---|---|---|---|---|
| Train acc. | **100%** | 97.2% | **100%** | 97.2% | **100%** | 99% |
| Test acc. | 55.5% | **56.5%** | 56.0% | **56.5%** | 56.7% | **57.7%** |

## 5    Conclusion

We show in this paper that taking a stepsize in SGD depending on the number of hidden units leads to particle systems with two possible mean-field behaviours. The first was already identified in [22, 23, 28] and corresponds to a deterministic mean-field ODE. The second is new and corresponds to a McKean-Vlasov diffusion. Our numerical experiments on two real datasets support our findings. In a future work, we intend to follow the same approach for deep neural networks, *i.e.*, with a growing number of hidden layers.

# References

[1] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.

[2] A. Krizhevsky, I. Sutskever, and G. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, pp. 1097–1105, 2012.

[3] C. Manning and H. Schütze, *Foundations of statistical natural language processing*. MIT press, 1999.

[4] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," *arXiv preprint arXiv:1611.03530*, 2016.

[5] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.

[6] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein, "Visualizing the loss landscape of neural nets," in *Advances in Neural Information Processing Systems*, pp. 6389–6399, 2018.

[7] A. Ballard, R. Das, S. Martiniani, D. Mehta, L. Sagun, J. Stevenson, and D. Wales, "Energy landscapes for machine learning," *Physical Chemistry Chemical Physics*, vol. 19, no. 20, pp. 12585–12603, 2017.

[8] M. Soltanolkotabi, A. Javanmard, and J. Lee, "Theoretical insights into the optimization landscape of over-parameterized shallow neural networks," *IEEE Trans. Information Theory*, vol. 65, no. 2, pp. 742–769, 2019.

[9] K. Fukumizu and S. Amari, "Local minima and plateaus in hierarchical structures of multilayer perceptrons," *Neural networks*, vol. 13, no. 3, pp. 317–327, 2000.

[10] A. J. Bray and D. Dean, "Statistics of critical points of gaussian fields on large-dimensional spaces," *Physical review letters*, vol. 98, no. 15, p. 150201, 2007.

[11] R. Pascanu, Y. N. Dauphin, S. Ganguli, and Y. Bengio, "On the saddle point problem for non-convex optimization," *arXiv preprint arXiv:1405.4604*, 2014.

[12] J. Pennington and Y. Bahri, "Geometry of neural network loss surfaces via random matrix theory," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 2798–2806, JMLR. org, 2017.

[13] K. Kawaguchi, "Deep learning without poor local minima," in *Advances in neural information processing systems*, pp. 586–594, 2016.

[14] D. Freeman and J. Bruna, "Topology and geometry of half-rectified network optimization," *arXiv preprint arXiv:1611.01540*, 2016.

[15] L. Venturi, A. Bandeira, and J. Bruna, "Neural networks with finite intrinsic dimension have no spurious valleys," *CoRR*, vol. abs/1802.06384, 2018.

[16] F. Bach, "Breaking the curse of dimensionality with convex neural networks," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 629–681, 2017.

[17] A. Choromanska, M. Henaff, M. Mathieu, G. Ben Arous, and Y. LeCun, "The loss surfaces of multilayer networks," in *Artificial intelligence and statistics*, pp. 192–204, 2015.

[18] L. Venturi, A. S. Bandeira, and J. Bruna, "Spurious valleys in one-hidden-layer neural network optimization landscapes," *Journal of Machine Learning Research*, vol. 20, no. 133, pp. 1–34, 2019.

[19] R. Kuditipudi, X. Wang, H. Lee, Y. Zhang, Z. Li, W. Hu, R. Ge, and S. Arora, "Explaining landscape connectivity of low-cost solutions for multilayer nets," in *Advances in Neural Information Processing Systems 32* (H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, and R. Garnett, eds.), pp. 14601–14610, Curran Associates, Inc., 2019.

[20] Z. Allen-Zhu, Y. Li, and Z. Song, "On the convergence rate of training recurrent neural networks," in *Advances in Neural Information Processing Systems 32* (H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, and R. Garnett, eds.), pp. 6676–6688, Curran Associates, Inc., 2019.

[21] S. S. Mannelli, G. Biroli, C. Cammarota, F. Krzakala, and L. Zdeborová, "Who is afraid of big bad minima? analysis of gradient-flow in spiked matrix-tensor models," in *Advances in Neural Information Processing Systems 32* (H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, and R. Garnett, eds.), pp. 8679–8689, Curran Associates, Inc., 2019.

[22] J. Sirignano and K. Spiliopoulos, "Mean field analysis of neural networks," *arXiv preprint arXiv:1805.01053*, 2018.

[23] S. Mei, A. Montanari, and P. Nguyen, "A mean field view of the landscape of two-layer neural networks," *Proceedings of the National Academy of Sciences*, vol. 115, no. 33, pp. E7665–E7671, 2018.

[24] G. M. Rotskoff and E. Vanden-Eijnden, "Trainability and accuracy of neural networks: An interacting particle system approach," 2018.

[25] S. Mei, T. Misiakiewicz, and A. Montanari, "Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit," *arXiv preprint arXiv:1902.06015*, 2019.

[26] A. Javanmard, M. Mondelli, and A. Montanari, "Analysis of a two-layer neural network via displacement convexity," *arXiv preprint arXiv:1901.01375*, 2019.

[27] L. Chizat, "Sparse optimization on measures with over-parameterized gradient descent," *arXiv preprint arXiv:1907.10300*, 2019.

[28] L. Chizat and F. Bach, "On the global convergence of gradient descent for over-parameterized models using optimal transport," in *Advances in neural information processing systems*, pp. 3036–3046, 2018.

[29] J. Jabir, D. vSivska, and L. Szpruch, "Mean-field neural odes via relaxed optimal control," *arXiv preprint arXiv:1912.05475*, 2019.

[30] L. Ambrosio, N. Gigli, and G. Savaré, *Gradient flows in metric spaces and in the space of probability measures*. Lectures in Mathematics ETH Zürich, Birkhäuser Verlag, Basel, second ed., 2008.

[31] M. Erbar, "The heat equation on manifolds as a gradient flow in the wasserstein space," in *Annales de l'institut Henri Poincaré (B)*, vol. 46, pp. 1–23, 2010.

[32] L. Ambrosio, G. Savaré, and L. Zambotti, "Existence and stability for Fokker–Planck equations with log-concave reference measure," *Probability Theory and Related Fields*, vol. 145, no. 3, pp. 517–564, 2009.

[33] A. Sznitman, "Topics in propagation of chaos," in *Ecole d'été de probabilités de Saint-Flour XIX—1989*, pp. 165–251, Springer, 1991.

[34] A. Gottlieb, "Markov transitions and the propagation of chaos," *arXiv preprint math/0001076*, 2000.

[35] S. Méléard and S. Roelly-Coppoletta, "Systèmes de particules et mesures-martingales: un théorème de propagation du chaos," *Séminaire de probabilités de Strasbourg*, vol. 22, pp. 438–448, 1988.

[36] R. Jordan, D. Kinderlehrer, and F. Otto, "The variational formulation of the fokker–planck equation," *SIAM journal on mathematical analysis*, vol. 29, no. 1, pp. 1–17, 1998.

[37] H. P. McKean, Jr., "Propagation of chaos for a class of non-linear parabolic equations," in *Stochastic Differential Equations (Lecture Series in Differential Equations, Session 7, Catholic Univ., 1967)*, pp. 41–57, Air Force Office Sci. Res., Arlington, Va., 1967.

[38] X. Fontaine, V. D. Bortoli, and A. Durmus, "Continuous and discrete-time analysis of stochastic gradient descent for convex and non-convex functions," 2020.

[39] M. Welling and Y. Teh, "Bayesian learning via stochastic gradient langevin dynamics," in *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688, 2011.

[40] J. Sirignano and K. Spiliopoulos, "Mean field analysis of neural networks: A central limit theorem," *Stochastic Processes and their Applications*, vol. 130, no. 3, pp. 1820–1852, 2020.

[41] Y. LeCun and C. Cortes, "MNIST handwritten digit database," 2010.

[42] A. Krizhevsky, G. Hinton, *et al.*, "Learning multiple layers of features from tiny images," 2009.

[43] C. Villani, *Optimal transport*, vol. 338 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 2009. Old and new.

[44] A. Kechris, *Classical Descriptive Set Theory*. Graduate Texts in Mathematics, Springer New York, 2012.

[45] D. W. Stroock and S. R. S. Varadhan, *Multidimensional diffusion processes*. Classics in Mathematics, Springer-Verlag, Berlin, 2006. Reprint of the 1997 edition.

[46] I. Karatzas and S. E. Shreve, *Brownian motion and stochastic calculus*, vol. 113 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, second ed., 1991.

[47] L. C. G. Rogers and D. Williams, *Diffusions, Markov processes, and martingales. Vol. 2*. Cambridge Mathematical Library, Cambridge University Press, Cambridge, 2000. Itô calculus, Reprint of the second (1994) edition.

[48] Y. Nesterov, *Introductory lectures on convex optimization*, vol. 87 of *Applied Optimization*. Kluwer Academic Publishers, Boston, MA, 2004. A basic course.

[49] L. Ambrosio and N. Gigli, "A user's guide to optimal transport," in *Modelling and optimisation of flows on networks*, vol. 2062 of *Lecture Notes in Math.*, pp. 1–155, Springer, Heidelberg, 2013.

[50] F. F. Bonsall and K. Vedak, *Lectures on some fixed point theorems of functional analysis*. No. 26, Tata Institute of Fundamental Research Bombay, 1962.

[51] J. Kent, "Time-reversible diffusions," *Advances in Applied Probability*, vol. 10, pp. 819–835, 12 1978.

# 6 Preliminaries

## 6.1 Notation

Let $(\mathsf{E}, d_E)$ and $(\mathsf{F}, d_F)$ be two metric spaces. $\mathrm{C}(\mathsf{E}, \mathsf{F})$ stands for the set of continuous $\mathsf{F}$-valued functions. If $\mathsf{F} = \mathbb{R}$, then we simply note $\mathrm{C}(\mathsf{E})$.

We say that $f : \mathsf{E} \to \mathbb{R}^p$ is $L$-Lipschitz if there exists $L \geq 0$ such that for any $x, y \in \mathsf{E}$, $\|f(x) - f(y)\| \leq L d_E(x, y)$. Let $\mathrm{C}_{\mathrm{b}}(\mathsf{E}, \mathbb{R}^p)$ (respectively $\mathrm{C}_{\mathrm{c}}(\mathsf{E}, \mathbb{R}^p)$) be the set of bounded continuous functions from $\mathsf{E}$ to $\mathbb{R}^p$ (respectively the set of compactly supported functions from $\mathsf{E}$ to $\mathbb{R}^p$). If $p = 1$, we simply note $\mathrm{C}_{\mathrm{b}}(\mathsf{E})$ (respectively $\mathrm{C}_{\mathrm{c}}(\mathsf{E})$).

For $\mathsf{U}$ an open set of $\mathbb{R}^d$, $n \in \mathbb{N}^\star$ and define $\mathrm{C}^n(\mathsf{U}, \mathbb{R}^p)$ the set of the $n$-differentiable $\mathbb{R}^p$-valued functions over $\mathsf{U}$. If $p = 1$ then we simply note $\mathrm{C}^n(\mathsf{U})$. Let $f \in \mathrm{C}^1(\mathsf{U})$ we denote by $\nabla f$ its gradient. More generally, if $f \in \mathrm{C}^n(\mathsf{U}, \mathbb{R}^p)$ with $n, p \in \mathbb{N}^\star$, we denote by $\mathrm{D}^k f(x)$ the $k$-th differential of $f$. We also denote for any $i \in \{1, \ldots, d\}$ and $\ell \in \{1, \ldots, k\}$, $\partial_i^\ell f$ the $i$-th partial derivative of $f$ of order $\ell$. If $f \in \mathrm{C}^2(\mathbb{R}^d, \mathbb{R})$, we denote by $\Delta f$ its Laplacian. $\mathrm{C}_{\mathrm{c}}^n(\mathsf{U}, \mathbb{R}^p)$ is the subset of $\mathrm{C}^n(\mathsf{U}, \mathbb{R}^p)$ such that for any $f \in \mathrm{C}_{\mathrm{c}}^n(\mathsf{U}, \mathbb{R}^p)$ and $\ell \in \{0, \ldots, n\}$, $\mathrm{D}^\ell f$ has compact support.

Consider $(\mathsf{F}, d)$ a metric space. Let $\mathscr{P}(\mathsf{F})$ be the space of probability measures over $\mathsf{F}$ equipped with its Borel $\sigma$-field $\mathcal{B}(\mathsf{F})$. For any $\mu \in \mathscr{P}(\mathsf{F})$ and $f : \mathsf{F} \to \mathbb{R}$, we say that $f$ is $\mu$-integrable if $\int_{\mathsf{F}} |f(x)| \mathrm{d}\mu(x) < +\infty$. In this case, we set $\mu[f] = \int_{\mathsf{F}} f(x) \mathrm{d}\mu(x)$. Let $\mu_0 \in \mathscr{P}(\mathsf{F})$. For any $r \geq 1$, define $\mathscr{P}_r(\mathsf{F}) = \{\mu \in \mathscr{P}(\mathsf{F}) : \int_{\mathbb{R}^p} d(\mu_0, \mu)^r \mathrm{d}\mu(x) < +\infty\}$. If not specified, we consider a filtered probability space $(\Omega, \mathcal{F}, \mathbb{P}, (\mathcal{F}_t)_{t \geq 0})$ satisfying the usual conditions and any random variables is defined on this probability space. Let $f : (\mathsf{E}, \mathcal{E}) \to (\mathsf{G}, \mathcal{G})$ be a measurable function. Then for any measure $\mu$ on $\mathcal{E}$ we define its pushforward measure by $f$, $f_\# \mu$, for any $\mathsf{A} \in \mathcal{G}$ by $f_\# \mu(\mathsf{A}) = \mu(f^{-1}(\mathsf{A}))$.

The set of $m \times n$ real matrices is denoted by $\mathbb{R}^{m \times n}$. The set of symmetric real matrices of size $p$ is denoted $\mathbb{S}_p(\mathbb{R})$.

## 6.2 Wasserstein distances

Let $(\mathsf{F}, d)$ be a metric space. Let $\mu_1, \mu_2 \in \mathscr{P}(\mathsf{F})$, where $\mathsf{F}$ is equipped with its Borel $\sigma$-field $\mathcal{B}(\mathsf{F})$. A probability measure $\zeta$ over $\mathcal{B}(\mathsf{F})^{\otimes 2}$ is said to be a transference plan between $\mu_1$ and $\mu_2$ if for any $\mathsf{A} \in \mathcal{B}(\mathsf{F})$, $\zeta(\mathsf{A} \times \mathsf{F}) = \mu_1(\mathsf{A})$ and $\zeta(\mathsf{F} \times \mathsf{A}) = \mu_2(\mathsf{A})$. We denote by $\Lambda(\mu_1, \mu_2)$ the set of all transference plans between $\mu_1$ and $\mu_2$. If $\mu_1, \mu_2 \in \mathscr{P}_r(\mathbb{R}^p)$, we define the Wasserstein distance $\mathscr{W}_r(\mu_1, \mu_2)$ of order $r$ between $\mu_1$ and $\mu_2$ by

$$\mathscr{W}_r^r(\mu_1, \mu_2) = \inf_{\zeta \in \Lambda(\mu_1, \mu_2)} \left\{ \int_{\mathsf{F} \times \mathsf{F}} d(x, y)^r \mathrm{d}\zeta(x, y) \right\} . \tag{10}$$

Note that $\mathcal{W}_r$ is a distance on $\mathscr{P}_r(\mathsf{F})$ by [43, Theorem 6.18]. In addition $(\mathscr{P}_r(\mathbb{R}^p), \mathcal{W}_r)$ is a complete separable metric space. For any $\mu_1, \mu_2 \in \mathscr{P}_p(\mathsf{F})$ we say that a couple of random variables $(X, Y)$ is an optimal coupling of $(\mu_1, \mu_2)$ for $\mathcal{W}_p$ if it has distribution $\xi$ where $\xi$ is an optimal transference plan between $\mu_1$ and $\mu_2$.

For any $T \geq 0$, the space $\mathscr{C}_{2,T}^p = \mathrm{C}([0,T], \mathscr{P}_2(\mathbb{R}^p))$ is a complete separable metric space [44, Theorem 4.19] with the metric $\mathcal{W}_{2,T}$ given for any $(\nu_t)_{t \in [0,T]}$ and $(\mu_t)_{t \in [0,T]}$ by

$$\mathcal{W}_{2,T}((\nu_t)_{t \in [0,T]}, (\mu_t)_{t \in [0,T]}) = \sup_{t \in [0,T]} \mathcal{W}_2(\nu_t, \mu_t) .$$

In the case where the measures we consider can be written as sums of Dirac we have the following proposition.

**Proposition 6.** *Let $r \geq 1$, $N \in \mathbb{N}^\star$, $\{\alpha_k\}_{k=1}^N \in [0,1]^N$ with $\sum_{k=1}^N \alpha_k = 1$, $\{\mu_{k,a}\}_{k=1}^N \in \mathscr{P}(\mathsf{F})^N$ and $\{\mu_{k,b}\}_{k=1}^N \in \mathscr{P}(\mathsf{F})^N$. Then, setting $\nu_i = \sum_{k=1}^N \alpha_k \mu_{k,i}$ with $i \in \{a, b\}$, we have*

$$\mathcal{W}_r^r(\nu_a, \nu_b)w \leq \sum_{k=1}^N \mathcal{W}_r^r(\mu_{k,a}, \mu_{k,b}) .$$

*Proof.* Consider $\zeta = \sum_{k=1}^N \alpha_k \zeta_k \in \Lambda(\nu_a, \nu_b)$ with $\zeta_k$ the optimal transference plan between $\mu_{k,a}$ and $\mu_{k,b}$. Then, we have

$$\mathcal{W}_r^r(\nu_a, \nu_b) \leq \int_{\mathbb{R}^p \times \mathbb{R}^p} d(x, y)^r \mathrm{d}\zeta(x, y) \leq N^{-1} \sum_{k=1}^N \mathcal{W}_r^r(\mu_{k,a}, \mu_{k,b}) .$$

$\square$

As a special case of Proposition 6, we obtain that for any $r \geq 1$, $\{w_{k,a}\}_{k=1}^N \in \mathsf{F}^N$ and $\{w_{k,a}\}_{k=1}^N \in \mathsf{F}^N$,

$$\mathcal{W}_r(N^{-1} \sum_{k=1}^N \delta_{w_{k,a}}, N^{-1} \sum_{k=1}^N \delta_{w_{k,b}}) \leq N^{-1} \sum_{k=1}^N d(w_{k,a}, w_{k,b})^r .$$

As another special case of Proposition 6, we obtain that for any $\mu \in \mathscr{P}_r(\mathsf{F})$ and $\{w_k\}_{k=1}^N \in \mathsf{F}^N$

$$\mathcal{W}_r(N^{-1} \sum_{k=1}^N \delta_{w_k}, N^{-1}, \mu) \leq N^{-1} \sum_{k=1}^N \mathcal{W}_r(w_k, \mu)^r .$$

# 7 A mean-field modification of Stochastic Gradient Langevin Dynamics

## 7.1 Presentation of the modified SGLD and its continuous counterpart

We start by introducing a modified Stochastic Gradient Langevin Dynamics (mSGLD) [39]. In the mean-field regime, this setting was studied in the case $\beta = 0$ in [23]. We recall that the mean-field $h : \mathbb{R}^p \times \mathscr{P}(\mathbb{R}^d) \to \mathbb{R}^p$ and $\xi : \mathbb{R}^p \times \mathscr{P}(\mathbb{R}^d) \times \mathsf{X} \times \mathsf{Y} \to \mathbb{R}^p$ are given for any $\mu \in \mathscr{P}(\mathbb{R}^p)$, $w \in \mathbb{R}^p$, $(x, y) \in \mathsf{X} \times \mathsf{Y}$ by

$$h(w, \mu) = -\int_{\mathsf{X} \times \mathsf{Y}} \partial_1 \ell\left(\mu[F(\cdot, x)], y\right) \nabla_w F(w, x) \, \mathrm{d}\pi(x, y) - \nabla V(w) ,$$

$$\xi(w, \mu, x, y) = -h(w, \mu) - \partial_1 \ell(\mu[F(\cdot, x)], y) \nabla_w F(w, x) - \nabla V(w) .$$

Let $(W_0^k)_{k \in \mathbb{N}^\star}$ be i.i.d. $p$ dimensional random variables with distribution $\mu_0$ and $\{Z_k^n : k, n \in \mathbb{N}^\star\}$ be i.i.d. $p$ dimensional independent Gaussian random variables with zero mean and identity covariance matrix. Consider the sequence $(W_n^{1:N})_{n \in \mathbb{N}}$ associated with mSGLD starting from $W_0^{1:N}$ and defined by the following recursion: for any $n \in \mathbb{N}$, $k \in \{1, \ldots, N\}$,

$$W_{n+1}^{k,N} = W_n^{k,N} + \gamma N^{\beta-1}(n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha} \left\{ h(W_n^{k,N}, \nu_n^N) + \xi(W_n^{k,N}, \nu_n^N, X_n, Y_n) \right\}$$

$$+ \left[2\eta\gamma N^{\beta-1}(n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha}\right]^{1/2} Z_{k,n} \ , \quad (11)$$

where $\eta \geq 0$, $\beta \in [0,1]$, $\alpha \in [0,1)$, $\gamma > 0$, $(X_n, Y_n)_{n\in\mathbb{N}}$ is a sequence of i.i.d. input/label samples distributed according to $\pi$ and $\gamma_{\alpha,\beta}(N) = \gamma^{1/(1-\alpha)} N^{(\beta-1)/(1-\alpha)}$. Note that in the cas $\eta = 0$, we obtain (3). In addition, (11) does not exactly correspond to the usual implementation of mSGLD as introduced in [39]. Indeed, to recover this algorithm, we should replace $[2\eta\gamma N^{\beta-1}(n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha}]^{1/2} Z_{k,n}$ by $[2\eta\gamma N^{\beta}(n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha}]^{1/2} Z_{k,n}$ in (11). The scheme presented in (11) amounts to consider a temperature which scales as $\gamma N^{\beta-1}$ with the number of particles. As emphasized before, this scheme was also considered in [23].

We now present the continuous model associated with this discrete process in the limit $\gamma \to 0$ or $N \to +\infty$. For $N \in \mathbb{N}^\star$, consider the particle system diffusion $(\mathbf{W}_t^{1:N})_{t\geq 0} = (\{\mathbf{W}_t^{k,N}\}_{k=1}^N)_{t\geq 0}$ starting from $\mathbf{W}_0^{1:N}$ defined for any $k \in \{1, \ldots, N\}$ by

$$d\mathbf{W}_t^{k,N} = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N)dt + \gamma_{\alpha,\beta}(N)^{1/2}\Sigma^{1/2}(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N)d\mathbf{B}_t^k + \sqrt{2\eta}d\tilde{\mathbf{B}}_t^k \right\} \ , \tag{12}$$

where $\{(\mathbf{B}_t^k)_{t\geq 0} : k \in \mathbb{N}^\star\}$ and $\{(\tilde{\mathbf{B}}_t^k)_{t\geq 0} : k \in \mathbb{N}^\star\}$ are two independent families of independent $p$ dimensional Brownian motions and $\boldsymbol{\nu}_t^N$ is the empirical probability distribution of the particles defined for any $t \geq 0$ by $\boldsymbol{\nu}_t^N = N^{-1}\sum_{k=1}^N \delta_{\mathbf{W}_t^{k,N}}$. Similarly to Section 2, (12) is the continuous counterpart of (11). Let $M \in \mathbb{N}^\star$. Similarly to (6), we consider the following particle system diffusion $(\mathbf{W}_t^{1:N})_{t\geq 0} = (\{\mathbf{W}_t^{k,N}\}_{k=1}^N)_{t\geq 0}$ starting from $\mathbf{W}_0^{1:N}$ defined for any $k \in \{1, \ldots, N\}$ by

$$d\mathbf{W}_t^{k,N} = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N)dt + (\gamma_{\alpha,\beta}(N)/M)^{1/2}\Sigma^{1/2}(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N)d\mathbf{B}_t^k + \sqrt{2\eta}d\tilde{\mathbf{B}}_t^k \right\} \ . \tag{13}$$

### 7.2 Mean field approximation and propagation of chaos for mSGLD

The following theorems are the extensions of Theorem 1 and Theorem 2 to (12) for any $\eta \geq 0$. Note that in the case $\eta = 0$, Theorem 7 boils down to Theorem 1 and Theorem 8 to Theorem 2.

We start by stating our results in the case $\beta \in [0,1)$. Consider the mean-field SDE starting from a random variable $\mathbf{W}_0^\star$ given by

$$d\mathbf{W}_t^\star = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^\star, \boldsymbol{\lambda}_t^\star)dt + \sqrt{2\eta}\tilde{\mathbf{B}}_t \right\} \ , \qquad \text{with } \boldsymbol{\lambda}_t^\star \text{ the distribution of } \mathbf{W}_t^\star \ . \tag{14}$$

**Theorem 7.** *Assume* **A**1. *Let* $(\mathbf{W}_0^k)_{k\in\mathbb{N}}$ *be a sequence of i.i.d.* $\mathbb{R}^p$*-valued random variables with distribution* $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$ *and set for any* $N \in \mathbb{N}^\star$, $\mathbf{W}_0^{1:N} = \{\mathbf{W}_0^k\}_{k=1}^N$. *Then, for any* $m \in \mathbb{N}^\star$ *and* $T \geq 0$, *there exists* $C_{m,T} \geq 0$ *such that for any* $\alpha \in [0,1)$, $\beta \in [0,1)$, $M \in \mathbb{N}^\star$ *and* $N \in \mathbb{N}^\star$

$$\mathbb{E}\left[\sup_{t\in[0,T]} \|\mathbf{W}_t^{1:m,N} - \mathbf{W}_t^{1:m,\star}\|^2\right] \leq C_{m,T}\left\{N^{-(1-\beta)/(1-\alpha)}M^{-1} + N^{-1}\right\} \ ,$$

*with* $(\mathbf{W}_t^{1:m,N}, \mathbf{W}_t^{1:m,\star}) = \{(\mathbf{W}_t^{k,N}, \mathbf{W}_t^{k,\star})\}_{k=1}^m$, $(\mathbf{W}_t^{1:N})$ *is the solution of* (13) *starting from* $\mathbf{W}_0^{1:N}$, *and for any* $k \in \{1, \ldots, N\}$, $\mathbf{W}_t^{k,\star}$ *is the solution of* (14) *starting from* $\mathbf{W}_0^k$ *and Brownian motion* $(\tilde{\mathbf{B}}_t^k)_{t\geq 0}$.

*Proof.* The proof is postponed to Section 9.4 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ □

Consider now the mean-field SDE starting from a random variable $\mathbf{W}_0^\star$ given by

$$d\mathbf{W}_t^\star = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^\star, \boldsymbol{\lambda}_t^\star)dt + (\gamma^{1/(1-\alpha)}\Sigma(\mathbf{W}_t^\star, \boldsymbol{\lambda}_t^\star)/M)^{1/2}d\mathbf{B}_t + \sqrt{2\eta}d\tilde{\mathbf{B}}_t \right\} \ , \tag{15}$$

where $\boldsymbol{\lambda}_t^\star$ is the distribution of $\mathbf{W}_t^\star$ and $(\mathbf{B}_t)_{t\geq 0}$ and $(\tilde{\mathbf{B}}_t)_{t\geq 0}$ are independent $p$ dimensional Brownian motions.

**Theorem 8.** *Let* $\beta = 1$. *Assume* **A**1. *Let* $(\mathbf{W}_0^k)_{k\in\mathbb{N}}$ *be a sequence of* $\mathbb{R}^p$*-valued random variables with distribution* $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$ *and assume that for any* $N \in \mathbb{N}^\star$, $\mathbf{W}_0^{1:N} = \{\mathbf{W}_0^k\}_{k=1}^N$. *Then, for any* $m \in \mathbb{N}^\star$ *and* $T \geq 0$, *there exists* $C_{m,T} \geq 0$ *such that for any* $\alpha \in [0,1)$, $M \in \mathbb{N}^\star$ *and* $N \in \mathbb{N}^\star$ *we have*

$$\mathbb{E}\left[\sup_{t\in[0,T]} \|\mathbf{W}_t^{1:m,N} - \mathbf{W}_t^{1:m,\star}\|^2\right] \leq C_{m,T}N^{-1} \ ,$$

13

with $(\mathbf{W}_t^{1:m,N}, \mathbf{W}_t^{1:m,\star}) = \{(\mathbf{W}_t^{k,N}, \mathbf{W}_t^{k,\star})\}_{k=1}^m$, $(\mathbf{W}_t^{1:N})$ *is the solution of* (13) *starting from* $\mathbf{W}_0^{1:N}$, *and for any* $k \in \{1, \ldots, N\}$, $\mathbf{W}_t^{k,\star}$ *is the solution of* (15) *starting from* $\mathbf{W}_0^k$ *and Brownian motions* $(\mathbf{B}_t^k)_{t \geq 0}$ *and* $(\tilde{\mathbf{B}}_t^k)_{t \geq 0}$.

*Proof.* The proof is postponed to Section 9.4 $\qquad\square$

## 8 Technical results

In this section, we derive technical results needed to establish Theorem 1, Theorem 2, Theorem 7 and Theorem 8. In particular, we are interested in the regularity properties of the mean field $h$ and the diffusion matrix $\Sigma$ under **A**1. We recall that in this setting, for any $w \in \mathbb{R}^p$, $\mu \in \mathscr{P}(\mathbb{R}^p)$, $(x, y) \in \mathsf{X} \times \mathsf{Y}$, we have

$$h(w, \mu) = \tilde{h}(w, \mu) - \nabla V(w) ,$$

$$\text{with} \quad \tilde{h}(w, \mu) = -\int_{\mathsf{X} \times \mathsf{Y}} \partial_1 \ell \left( \int_{\mathbb{R}^p} F(\zeta, x) \, \mathrm{d}\mu(\zeta), y \right) \nabla_w F(w, x) \, \mathrm{d}\pi(x, y) ,$$

$$\xi(w, \mu, x, y) = -\tilde{h}(w, \mu) - \partial_1 \ell \left( \int_{\mathbb{R}^p} F(\zeta, x) \, \mathrm{d}\mu(\zeta), y \right) \nabla_w F(w, x) ,$$

$$\Sigma(w, \mu) = \int_{\mathsf{X} \times \mathsf{Y}} \{\xi \xi^\top\}(w, \mu, x, y) \mathrm{d}\pi(x, y) , \qquad \mathrm{S}(w, \mu) = \Sigma^{1/2}(w, \mu) . \tag{16}$$

Note that by **A**1-(a), we obtain the following estimate used in the proof of the results of this Section: for any $\mathrm{y}, y \in \mathbb{R}$

$$|\partial_1 \ell(\mathrm{y}, y)| \leq |\partial_1 \ell(0, y)| + \Psi(y) |\mathrm{y}| \leq 2\Psi(y) \max(1, |\mathrm{y}|) . \tag{17}$$

In addition, note that under **A**1-(c), there exists $\mathrm{K} \geq 0$ such that for any $w \in \mathbb{R}^p$

$$\left\| \nabla^2 V(w) \right\| + \left\| \mathrm{D}^3 V(w) \right\| \leq \mathrm{K} , \qquad \|\nabla V(w)\| \leq \mathrm{K}(1 + \|w\|) . \tag{18}$$

Let $G : \mathbb{R}^p \times \mathsf{X} \times \mathsf{Y} \to \mathbb{R}$ given for any $(x, y) \in \mathsf{X} \times \mathsf{Y}$ and $w \in \mathbb{R}^p$ by

$$G(w, x, y) = \{\Phi^4(x) + \Psi^2(y)\} F(w, x) . \tag{19}$$

We now state our main regularity/boundedness proposition.

**Proposition 9.** *Assume* **A**1. *Then, there exists* $\mathrm{L} \geq 0$ *such that the following hold.*

*(a) For any* $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$ *and* $w_1, w_2 \in \mathbb{R}^p$ *we have*

$$\|h(w_1, \mu_1) - h(w_2, \mu_2)\|$$
$$\leq \mathrm{L} \left\{ \|w_1 - w_2\| + \left( \int_{\mathsf{X} \times \mathsf{Y}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 \, \mathrm{d}\pi(x, y) \right)^{1/2} \right\} . \tag{20}$$

*In addition, we have for any* $\mu \in \mathscr{P}(\mathbb{R}^p)$ *and* $w \in \mathbb{R}^p$, $\|h(w, \mu)\| \leq \mathrm{L}(1 + \|w\|)$ *and* $\|\bar{h}(w, \mu)\| \leq \mathrm{L}$.

*(b) For any* $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $w_1, w_2 \in \mathbb{R}^p$ *and* $i, j \in \{1, \ldots, p\}$ *we have*

$$|\mathrm{S}_{i,j}(w_1, \mu_1) - \mathrm{S}_{i,j}(w_2, \mu_2)|$$
$$\leq \mathrm{L} \left\{ \|w_1 - w_2\| + \left( \int_{\mathsf{X} \times \mathsf{Y}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 \, \mathrm{d}\pi(x, y) \right)^{1/2} \right\} . \tag{21}$$

*In addition, we have for any* $\mu \in \mathscr{P}(\mathbb{R}^p)$, $w \in \mathbb{R}^p$ *and* $i, j \in \{1, \ldots, p\}$, $|\mathrm{S}_{i,j}(w, \mu)| \leq \mathrm{L}$.

*(c) For any* $\mu \in \mathscr{P}(\mathbb{R}^p)$ *and* $w \in \mathbb{R}^p$, $\int_{\mathsf{X} \times \mathsf{Y}} \|\xi(w, \mu, x, y)\|^2 \, \mathrm{d}\pi(x, y) \leq p^2 \mathrm{L}^2$.

*Proof.* (a) First, we show that (20) holds. Note that by the triangle inequality and (16), we only need to consider $h \leftarrow \tilde{h}$ and $h \leftarrow V$. The case $h \leftarrow V$ is straightforward using (18). We now deal with the first case. For any $w_1, w_2 \in \mathbb{R}^p$ and $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, consider the decomposition,

$$\|\tilde{h}(w_1, \mu_1) - \tilde{h}(w_2, \mu_2)\| \leq \|\tilde{h}(w_1, \mu_1) - \tilde{h}(w_2, \mu_1)\| + \|\tilde{h}(w_2, \mu_1) - \tilde{h}(w_2, \mu_2)\| .$$

In what follows, we bound separately the two terms in the right-hand side. Using **A**1-(a), **A**1-(b), (16) and (17) we have for any $w_1, w_2 \in \mathbb{R}^p$ and $\mu_1 \in \mathscr{P}(\mathbb{R}^p)$

$$\begin{aligned}
\|\tilde{h}(w_1, \mu_1) - \tilde{h}(w_2, \mu_1)\| &\leq \left\| \int_{\mathsf{X} \times \mathsf{Y}} \partial_1 \ell(\mu_1[F(\cdot, x)], y) \nabla_w F(w_1, x) \, \mathrm{d}\pi(x, y) \right. \\
&\qquad \left. - \int_{\mathsf{X} \times \mathsf{Y}} \partial_1 \ell(\mu_1[F(\cdot, x)], y) \nabla_w F(w_2, x) \, \mathrm{d}\pi(x, y) \right\| \\
&\leq \int_{\mathsf{X} \times \mathsf{Y}} |\partial_1 \ell(\mu_1[F(\cdot, x)], y)| \, \Phi(x) \mathrm{d}\pi(x, y) \, \|w_1 - w_2\| \\
&\leq \int_{\mathsf{X} \times \mathsf{Y}} \Psi(y) \Phi(x) \left(1 + |\mu_1[F(\cdot, x)]|\right) \mathrm{d}\pi(x, y) \, \|w_1 - w_2\| \\
&\leq 2 \int_{\mathsf{X} \times \mathsf{Y}} \Psi(y) \Phi^2(x) \mathrm{d}\pi(x, y) \, \|w_1 - w_2\| .
\end{aligned} \tag{22}$$

Using **A**1-(a), **A**1-(b), (16) and the Cauchy-Schwarz inequality, we also have for any $w_1 \in \mathbb{R}^p$ and $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$

$$\begin{aligned}
\|\tilde{h}(\mu_1, w_1) &- \tilde{h}(\mu_2, w_1)\| \\
&\leq \left\| \int_{\mathsf{X} \times \mathsf{Y}} \{\partial_1 \ell(\mu_1[F(\cdot, x)], y) \nabla_w F(w_1, x) - \partial_1 \ell(\mu_2[F(\cdot, x)], y) \nabla_w F(w_1, x)\} \, \mathrm{d}\pi(x, y) \right\| \\
&\leq \int_{\mathsf{X} \times \mathsf{Y}} |\partial_1 \ell(\mu_1[F(\cdot, x)], y) - \partial_1 \ell(\mu_2[F(\cdot, x)], y)| \, \|\nabla_w F(w_1, x)\| \, \mathrm{d}\pi(x, y) \\
&\leq \int_{\mathsf{X} \times \mathsf{Y}} \Psi(y) \, \|\mu_1[F(\cdot, x)] - \mu_2[F(\cdot, x)]\| \, \Phi(x) \mathrm{d}\pi(x, y) \\
&\leq \left( \int_{\mathsf{X} \times \mathsf{Y}} \Psi^2(y) \Phi^2(x) \mathrm{d}\pi(x, y) \right)^{1/2} \left( \int_{\mathsf{X}} \|\mu_1[F(\cdot, x)] - \mu_2[F(\cdot, x)]\|^2 \, \mathrm{d}\pi(x) \right)^{1/2} . \quad (23)
\end{aligned}$$

Combining (19), (22), (23), the fact that for any $a, b \geq 0$, $2ab \leq a^2 + b^2$ and **A**1-(d), we obtain that there exists $\mathsf{L}_1 \geq 0$ such that for any $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$ and $w_1, w_2 \in \mathbb{R}^p$ we have

$$\begin{aligned}
\|\tilde{h}(w_1, \mu_1) &- \tilde{h}(w_2, \mu_2)\| \\
&\leq \mathsf{L}_1 \left\{ \|w_1 - w_2\| + \left( \int_{\mathsf{X} \times \mathsf{Y}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 \, \mathrm{d}\pi(x, y) \right)^{1/2} \right\} .
\end{aligned}$$

In addition, using **A**1-(b) and (17), we have for any $w \in \mathbb{R}^p$, $\mu \in \mathscr{P}(\mathbb{R}^p)$, $x \in \mathsf{X}$ and $y \in \mathsf{Y}$

$$|\partial_1 \ell(\mu[F(\cdot, x)], y)| \, \|\nabla_w F(w, x)\| \leq \Psi(y) \Phi(x)(1 + \Phi(x)) \leq 2\Psi(y) \Phi^2(x) . \tag{24}$$

Therefore, combining this result and (16), we get that for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$

$$\|\tilde{h}(w, \mu)\| \leq \int_{\mathsf{X} \times \mathsf{Y}} 2\Psi(y) \Phi^2(x) \mathrm{d}\pi(x, y) .$$

Using the fact that for any $a, b \geq 0$, $2ab \leq a^2 + b^2$ and **A**1-(d), there exists $\mathsf{L}_2 \geq 0$ such that for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$,

$$\|\tilde{h}(w, \mu)\| \leq \mathsf{L}_2 \tag{25}$$

(b) Second, we first show that there exists $\mathsf{L}_3 \geq 0$ such that for any $\mu \in \mathscr{P}(\mathbb{R}^p)$, $w \in \mathbb{R}^p$ and $i, j \in \{1, \ldots, p\}$, $|\mathsf{S}_{i,j}(w, \mu)| \leq \mathsf{L}$. Let $i, j \in \{1, \ldots, p\}$. We have for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$

$$|\mathsf{S}_{i,j}(w, \mu)| \leq \|\mathsf{S}(w, \mu)\| \leq \mathrm{Tr}^{1/2} \left(\Sigma(w, \mu)\right) . \tag{26}$$

15

Similarly to (24), using (16), (25), the fact that for any $a, b \geq 0$, $(a+b)^2 \leq 2(a^2 + b^2)$ and the Cauchy-Schwarz inequality, we get for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$

$$\mathrm{Tr}\left(\Sigma(w,\mu)\right) \leq \int_{\mathsf{X} \times \mathsf{Y}} \|\xi(w,\mu,x,y)\|^2 \, \mathrm{d}\pi(x,y) \leq 2 \int_{\mathsf{X} \times \mathsf{Y}} \{\mathrm{L}_2^2 + 2\Psi^2(y)\Phi^4(x)\} \mathrm{d}\pi(x,y) \,. \quad (27)$$

Combining (26), (27) and **A**1-(d), there exists $\mathrm{L}_3 \geq 0$ such that for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$, $\max_{1 \leq i,j \leq p} |\mathrm{S}_{i,j}(w,\mu)| \leq \mathrm{L}_3$.

We now show that (21) holds. For any $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$ define $\varphi_\Sigma : [0,1] \to \mathbb{S}_p(\mathbb{R})$ for any $t \in [0,1]$ by

$$\varphi_\Sigma(t) = \Sigma(tw_1 + (1-t)w_2, t\mu_1 + (1-t)\mu_2) \,. \quad (28)$$

For ease of notation, the dependency of $\varphi_\Sigma$ with respect to $w_1, w_2 \in \mathbb{R}^p$ and $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$ is omitted. In what follows, we show that for any $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $\varphi_\Sigma \in \mathrm{C}^2([0,1], \mathbb{S}_p(\mathbb{R}))$ and that there exists $\mathrm{L}_4 \geq 0$ such that for any $t \in [0,1]$

$$\|\varphi_\Sigma''(t)\| \leq \mathrm{L}_4 \left\{ \|w_1 - w_2\| + \left( \int_{\mathsf{X} \times \mathsf{Y}} \|\mu_1[G(\cdot,x,y)] - \mu_2[G(\cdot,x,y)]\|^2 \, \mathrm{d}\pi(x,y) \right)^{1/2} \right\}^2 \,,$$

which will conclude the proof of (21) upon using a straightforward adaptation of [45, Lemma 3.2.3, Theorem 5.2.3]. We conclude the proof of Proposition 9 upon letting $\mathrm{L} = \max(\mathrm{L}_1, \mathrm{L}_2, \mathrm{L}_3, \mathrm{L}_4)$.

For any $t \in [0,1]$, let $\mu_t = \mu_1 + t(\mu_2 - \mu_1) \in \mathscr{P}(\mathbb{R}^p)$ and $w_t = w_1 + t(w_2 - w_1) \in \mathbb{R}^p$ and for any $(x,y) \in \mathsf{X} \times \mathsf{Y}$ define

$$\begin{aligned} \mathrm{f}(t,x,y) &= \partial_1 \ell(\mu_t[F(\cdot,x)], y) \nabla_w F(w_t, x) \,, \\ \tilde{\mathrm{f}}(t,x,y) &= \xi(w_t, \mu_t, x, y) = \int_{\mathsf{X} \times \mathsf{Y}} \mathrm{f}(t,x,y) \mathrm{d}\pi(x,y) - \mathrm{f}(t,x,y) \,. \end{aligned} \quad (29)$$

The rest of the proof consists in showing that $\varphi_\Sigma$ is twice differentiable with dominated derivatives using the Lebesgue convergence theorem.

By (16), (24) and (25), we get that for any $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $(x,y) \in \mathsf{X} \times \mathsf{Y}$ and $t \in [0,1]$

$$\|\mathrm{f}(t,x,y)\| \leq 2\Psi(y)\Phi^2(x) \,, \qquad \|\tilde{\mathrm{f}}(t,x,y)\| \leq \mathrm{L}_2 + 2\Psi(y)\Phi^2(x) \,. \quad (30)$$

Using (29), **A**1-(a) and **A**1-(b), we have that for any $(x,y) \in \mathsf{X} \times \mathsf{Y}$, $\mathrm{f}(\cdot,x,y) \in \mathrm{C}^1([0,1], \mathbb{R}^p)$ and for any $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $(x,y) \in \mathsf{X} \times \mathsf{Y}$ and $t \in [0,1]$

$$\begin{aligned} \partial_1 \mathrm{f}(t,x,y) &= \partial_1^2 \ell(\mu_t[F(\cdot,x)], y) \nabla_w F(w_t, x) \left(\mu_2[F(\cdot,x)] - \mu_1[F(\cdot,x)]\right) \\ &\quad + \partial_1 \ell(\mu_t[F(\cdot,x)], y) \nabla_w^2 F(w_t, x)(w_2 - w_1) \,. \quad (31) \end{aligned}$$

Using **A**1-(a), **A**1-(b), (19) and (17), we get that for any $(x,y) \in \mathsf{X} \times \mathsf{Y}$ and $t \in [0,1]$

$$\|\partial_1 \mathrm{f}(t,x,y)\| \leq 3\Psi(y)\Phi^2(x) \left(\|w_2 - w_1\| + \|\mu_1[F(\cdot,x)] - \mu_2[F(\cdot,x)]\|\right) \,, \quad (32)$$

Similarly, using (31), **A**1-(a) and **A**1-(b), we have that for any $(x,y) \in \mathsf{X} \times \mathsf{Y}$, $\mathrm{f}(\cdot,x,y) \in \mathrm{C}^2([0,1], \mathbb{R}^p)$ and for any $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $(x,y) \in \mathsf{X} \times \mathsf{Y}$ and $t \in [0,1]$

$$\begin{aligned} \partial_1^2 \mathrm{f}(t,x,y) &= \partial_1^3 \ell(\mu_t[F(\cdot,x)], y) \nabla_w F(w_t, x) \left(\mu_2[F(\cdot,x)] - \mu_1[F(\cdot,x)]\right)^2 \\ &\quad + 2\partial_1^2 \ell(\mu_t[F(\cdot,x)], y) \nabla_w^2 F(w_t, x)(w_2 - w_1) \left(\mu_2[F(\cdot,x)] - \mu_1[F(\cdot,x)]\right) \\ &\quad\quad + \partial_1 \ell(\mu_t[F(\cdot,x)], y) \mathrm{D}_w^3 F(w_t, x)(w_2 - w_1)^{\otimes 2} \,. \end{aligned}$$

Using **A**1-(a), **A**1-(b) and (17) and that for any $a, b \geq 0$, $2ab \leq a^2 + b^2$, we get that for any $(x,y) \in \mathsf{X} \times \mathsf{Y}$ and $t \in [0,1]$

$$\left\|\partial_1^2 \mathrm{f}(t,x,y)\right\| \leq 5\Psi(y)\Phi^2(x) \left(\|w_2 - w_1\|^2 + \|\mu_1[F(\cdot,x)] - \mu_2[F(\cdot,x)]\|^2\right) \,. \quad (33)$$

Combining (29), (32), (33), **A**1-(d) and the dominated convergence theorem, we get that for any $(x,y) \in \mathsf{X} \times \mathsf{Y}$, $\tilde{\mathrm{f}}(\cdot,x,y) \in \mathrm{C}^2([0,1], \mathbb{R}^p)$. In addition, using (29), (30), (32), (33), the Cauchy-Schwarz inequality and the fact that for any $a, b \geq 0$, $2ab \leq a^2 + b^2$, there exists $C \geq 0$, such that for any $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $(x,y) \in \mathsf{X} \times \mathsf{Y}$ and $t \in [0,1]$

$$\|\tilde{\mathrm{f}}(t,x,y)\| \leq C \left(\Phi^4(x) + \Psi^2(y)\right) \,,$$

16

$$\|\partial_1 \tilde{f}(t, x, y)\| \leq C \left( \Phi^4(x) + \Psi^2(y) \right) \chi(w_1, w_2, \mu_1, \mu_2, x) \,,$$
$$\|\partial_1^2 \tilde{f}(t, x, y)\| \leq C \left( \Phi^4(x) + \Psi^2(y) \right) \chi^2(w_1, w_2, \mu_1, \mu_2, x) \,, \tag{34}$$

where

$$\chi(w_1, w_2, \mu_1, \mu_2, x) = \|w_1 - w_2\|$$
$$+ \|\mu_1[F(\cdot, x)] - \mu_2[F(\cdot, x)]\| + \left( \int_{\mathsf{X} \times \mathsf{Y}} \|\mu_1[G(\cdot, \tilde{x}, \tilde{y})] - \mu_2[G(\cdot, \tilde{x}, \tilde{y})]\|^2 \, \mathrm{d}\pi(\tilde{x}, \tilde{y}) \right)^{1/2} .$$

Using (28) and (16), we have that for any $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $t \in [0, 1]$

$$\varphi_\Sigma(t) = \int_{\mathsf{X} \times \mathsf{Y}} \tilde{f}(t, x, y) \tilde{f}(t, x, y)^\top \mathrm{d}\pi(x, y) \,.$$

Combining this result, (34) and **A**1-(d) we get that for any $w_1, w_2 \in \mathbb{R}^p$ and $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $\varphi_\Sigma \in \mathrm{C}^2([0, 1], \mathbb{S}_p(\mathbb{R}))$ and, using the Cauchy-Schwarz inequality, there exist $C_1, C_2 \geq 0$ such that for any $w_1, w_2 \in \mathbb{R}^p$ and $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $t \in [0, 1]$ and $\mathrm{u} \in \mathbb{R}^p$ with $\|\mathrm{u}\| = 1$, we have

$$\langle \mathrm{u}, \varphi_\Sigma''(t) \mathrm{u} \rangle = \int_{\mathsf{X} \times \mathsf{Y}} \partial_1^2 \left( \langle \mathrm{u}, \tilde{f}(t, x, y) \rangle^2 \right) \mathrm{d}\pi(x, y)$$

$$\leq 2 \int_{\mathsf{X} \times \mathsf{Y}} \|\partial_1 \tilde{f}(t, x, y)\|^2 \mathrm{d}\pi(x, y) + 2 \int_{\mathsf{X} \times \mathsf{Y}} \|\partial_1^2 \tilde{f}(t, x, y)\| \|\tilde{f}(t, x, y)\| \mathrm{d}\pi(x, y)$$

$$\leq C_1 \int_{\mathsf{X} \times \mathsf{Y}} \left( \Phi^8(x) + \Psi^4(y) \right) \chi^2(w_1, w_2, x, y) \mathrm{d}\pi(x, y)$$

$$\leq C_2 \left\{ \|w_1 - w_2\| + \left( \int_{\mathsf{X}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 \, \mathrm{d}\pi(x, y) \right)^{1/2} \right\}^2 \,,$$

Therefore, we get that for any $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$, $t \in [0, 1]$

$$\|\varphi_\Sigma''(t)\| = \sup_{\mathrm{u} \in \mathbb{R}^p, \|\mathrm{u}\| = 1} \langle \mathrm{u}, \varphi_\Sigma''(t) \mathrm{u} \rangle$$

$$\leq C \left\{ \|w_1 - w_2\| + \left( \int_{\mathsf{X}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 \, \mathrm{d}\pi(x, y) \right)^{1/2} \right\}^2 \,.$$

Combining this result and a straightforward adaptation of [45, Lemma 3.2.3, Theorem 5.2.3] we obtain that for any $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R}^p)$

$$|\mathsf{S}_{i,j}(w_1, \mu_1) - \mathsf{S}_{i,j}(w_2, \mu_2)| \leq \mathsf{L}_4 \left\{ \|w_1 - w_2\| + \left( \int_{\mathsf{X}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 \, \mathrm{d}\pi(x, y) \right)^{1/2} \right\} ,$$

with $\mathsf{L}_4 = \sqrt{2C} p$.

(c) Using (16), we have for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$

$$\int_{\mathsf{X} \times \mathsf{Y}} \|\xi(w, \mu, x, y)\|^2 \, \mathrm{d}\pi(x, y) = \int_{\mathsf{X} \times \mathsf{Y}} \mathrm{Tr} \left( \xi \xi^\top (w, \mu, x, y) \right) \mathrm{d}\pi(x, y) = \sum_{i,j=1}^{p} |\mathsf{S}_{i,j}(w, \mu)|^2 \leq p^2 \mathsf{L}^2 \,.$$

$\square$

# 9 Quantitative propagation of chaos

## 9.1 Existence of strong solutions to the particle SDE

In this section, for two functions $A, B : \bigcup_{N \in \mathbb{N}^\star} \left\{ \{1, \ldots, N\} \times \mathbb{R}_+ \times (\mathbb{R}^p)^2 \times (\mathscr{P}_2(\mathbb{R}^p))^2 \right\} \to \mathbb{R}$, the notation $A_N(k, t, w_1, w_2, \mu_1, \mu_2) \lesssim B_N(k, t, w_1, w_2, \mu_1, \mu_2)$ stands for the statement that there

exists $C \geq 0$ such that for any $N \in \mathbb{N}^\star$, $k \in \{1, \ldots, N\}$, $t \in \mathbb{R}_+$, $w_1, w_2 \in \mathbb{R}^p$, $\mu_1, \mu_2 \in \mathscr{P}_2(\mathbb{R}^p)$, $A_N(k, t, w_1, w_2, \mu_1, \mu_2) \leq CB_N(k, t, w_1, w_2, \mu_1, \mu_2)$, where $A_N$ and $B_N$ are the restrictions of $A$ and $B$ to $\{1, \ldots, N\} \times \mathbb{R}_+ \times (\mathbb{R}^p)^2 \times (\mathscr{P}_2(\mathbb{R}^p))^2$.

We consider for $N \in \mathbb{N}^\star$, $p$ dimensional particle system $(\mathbf{W}_t^{1:N})_{t \geq 0}$ associated with the SDE: for any $k \in \{1, \ldots, N\}$

$$\mathrm{d}\mathbf{W}_t^{k,N} = b_N(t, \mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N)\mathrm{d}t + \sigma_N(t, \mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N)\mathrm{d}\mathbf{B}_t^k , \qquad \boldsymbol{\nu}_t^N = (1/N)\sum_{k=1}^N \delta_{\mathbf{W}_t^{k,N}} , \quad (35)$$

where $(\mathbf{B}_t^k)_{k \in \mathbb{N}^\star}$ are independent $r$-dimensional Brownian motions and where $(b_N)_{N \in \mathbb{N}^\star}$ and $(\sigma_N)_{N \in \mathbb{N}^\star}$ are family of measurable functions such that for any $N \in \mathbb{N}^\star$, $b_N : \mathbb{R}_+ \times \mathbb{R}^p \times \mathscr{P}_2(\mathbb{R}^p) \to \mathbb{R}^p$ and $\sigma_N : \mathbb{R}_+ \times \mathbb{R}^p \times \mathscr{P}_2(\mathbb{R}^p) \to \mathbb{R}^{p \times r}$. We make the following assumption ensuring the existence and uniqueness of solutions of (35) for any $N \in \mathbb{N}^\star$. Consider in the sequel a measurable space $(\mathsf{Z}, \mathcal{Z})$ and a probability measure $\pi_\mathsf{Z}$ on this space.

**B1.** *There exist a measurable function* $\mathrm{g} : \mathbb{R}^p \times \mathsf{Z} \to \mathbb{R}$, $\mathtt{M}_1 \geq 0$ *and* $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$ *such that for any* $N \in \mathbb{N}^\star$, *the following hold.*

*(a) For any* $w_1, w_2 \in \mathbb{R}^p$ *and* $z \in \mathsf{Z}$ *we have*

$$\|\mathrm{g}(w_1, z) - \mathrm{g}(w_2, z)\| \leq \zeta(z) \|w_1 - w_2\| , \quad \|\mathrm{g}(w_1, z)\| \leq \zeta(z) , \quad \text{with } \int_\mathsf{Z} \zeta^2(z)\mathrm{d}\pi_\mathsf{Z}(z) < +\infty .$$

*(b)* $b_N \in \mathrm{C}(\mathbb{R}_+ \times \mathbb{R}^p \times \mathscr{P}_2(\mathbb{R}^p), \mathbb{R}^p)$ *and* $\sigma_N \in \mathrm{C}(\mathbb{R}_+ \times \mathbb{R}^p \times \mathscr{P}_2(\mathbb{R}^p), \mathbb{R}^{p \times r})$.

*(c) For any* $w_1, w_2 \in \mathbb{R}^p$ *and* $\mu_1, \mu_2 \in \mathscr{P}_2(\mathbb{R}^p)$

$$\sup_{t \geq 0}\{\|b_N(t, w_1, \mu_1) - b_N(t, w_2, \mu_2)\| + \|\sigma_N(t, w_1, \mu_1) - \sigma_N(t, w_2, \mu_2)\|\}$$
$$\leq \mathtt{M}_1 \left\{ \|w_1 - w_2\| + \left( \int_\mathsf{Z} |\mu_1[\mathrm{g}(\cdot, z)] - \mu_2[\mathrm{g}(\cdot, z)]|^2 \, \mathrm{d}\pi_\mathsf{Z}(z) \right)^{1/2} \right\} ,$$
$$\sup_{t \geq 0} \{\|b_N(t, 0, \mu_0)\| + \|\sigma_N(t, 0, \mu_0)\|\} \leq \mathtt{M}_1 .$$

**B2.** *There exist* $\mathtt{M}_2 \geq 0$, $\kappa > 0$, $b \in \mathrm{C}(\mathbb{R}_+ \times \mathbb{R}^p \times \mathscr{P}_2(\mathbb{R}^p), \mathbb{R}^p)$ *and* $\sigma \in \mathrm{C}(\mathbb{R}_+ \times \mathbb{R}^p \times \mathscr{P}_2(\mathbb{R}^p), \mathbb{R}^{p \times r})$ *such that*

$$\sup_{t \geq 0, w \in \mathbb{R}^p, \mu \in \mathscr{P}_2(\mathbb{R}^p)} \{\|b_N(t, w, \mu) - b(t, w, \mu)\| + \|\sigma_N(t, w, \mu) - \sigma(t, w, \mu)\|\} \leq \mathtt{M}_2 N^{-\kappa} .$$

Note that under **B1**, we have the following estimate which will be used in our next result,

$$\|b_N(t, w, \mu)\| + \|\sigma_N(t, w, \mu)\| \lesssim \left[ 1 + \|w\| + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2)\mathrm{d}\mu(\tilde{w}) \right)^{1/2} \right] , \qquad (36)$$

$$\sup_{t \geq 0}\{\|b_N(t, w_1, \mu_1) - b_N(t, w_2, \mu_2)\| + \|\sigma_N(t, w_1, \mu_1) - \sigma_N(t, w_2, \mu_2)\|\}$$
$$\lesssim \|w_1 - w_2\| + \mathcal{W}_2(\mu_1, \mu_2) .$$

**Theorem 10.** *Assume* **B1**. *Then for any* $N \in \mathbb{N}^\star$, *(35) admits a unique strong solution. If in addition, there exists* $m \geq 1$ *such that* $\sup_{N \in \mathbb{N}^\star} \sup_{k \in \{1, \ldots, N\}} \mathbb{E}[\|\mathbf{W}_0^{k,N}\|^{2m}] < +\infty$, *then for any* $T \geq 0$, *there exists* $C \geq 0$ *such that*

$$\sup_{N \in \mathbb{N}^\star} \sup_{k \in \{1, \ldots, N\}} \mathbb{E}\left[ \sup_{t \in [0,T]} \left\| \mathbf{W}_t^{k,N} \right\|^{2m} \right] \leq C .$$

*Proof.* First, we show that for any $N \in \mathbb{N}^\star$, (35) admits a unique strong solution. Let $\tilde{b}_N : \mathbb{R}_+ \times (\mathbb{R}^p)^N \to (\mathbb{R}^p)^N$ and $\tilde{\sigma}_N : \mathbb{R}_+ \times (\mathbb{R}^p)^N \to (\mathbb{R}^{p \times r})^N$ given, setting $\nu^{N,w} = (1/N)\sum_{j=1}^N \delta_{w^{j,N}}$ for any $t \geq 0$ and $w^{1:N} \in (\mathbb{R}^p)^N$, by

$$\tilde{b}_N(t, w^{1:N}) = \left( b_N\left(t, w^{k,N}, \nu^{N,w}\right) \right)_{k \in \{1, \ldots, N\}} , \tilde{\sigma}_N(t, w^{1:N}) = \left( \sigma_N\left(t, w^{k,N}, \nu^{N,w}\right) \right)_{k \in \{1, \ldots, N\}} .$$

Let $w_1^{1:N}, w_2^{1:N} \in (\mathbb{R}^p)^N$. Using **B1**, Proposition 6 and that for any $a, b \geq 0$, $(a+b)^{1/2} \leq a^{1/2} + b^{1/2}$, we have

$$\|b_N(t, w_1^{k,N}, \nu^{N,w_1}) - b_N(t, w_2^{k,N}, \nu^{N,w_2})\| \lesssim \|w_1^{k,N} - w_2^{k,N}\| + \mathcal{W}_2(\nu^{N,w_1}, \nu^{N,w_2})$$
$$\lesssim \|w_1^{k,N} - w_2^{k,N}\| + (N^{-1} \textstyle\sum_{j=1}^N \|w_1^{j,N} - w_2^{j,N}\|^2)^{1/2} \lesssim \|w_1^{1:N} - w_2^{1:N}\| .$$

Similarly, we have $\|\sigma_N(t, w_1^{k,N}, \nu^{N,w_1}) - \sigma_N(t, w_2^{k,N}, \nu^{N,w_2})\| \lesssim \|w^{1:N} - w_2^{1:N}\|$. Therefore, we obtain that for any $N \in \mathbb{N}^\star$, $\tilde{b}_N$ and $\tilde{\sigma}_N$ are Lipschitz-continuous and using [46, Theorem 2.9], we get that there exists a unique strong solution to (35). Let $m \geq 1$ and assume that $\sup_{N \in \mathbb{N}^\star} \sup_{k \in \{1,\ldots,N\}} \mathbb{E}[\|\mathbf{W}_0^{k,N}\|^{2m}] < +\infty$, we now show that for any $T \geq 0$, there exists $C \geq 0$ such that

$$\sup_{t \in [0,T]} \sup_{N \in \mathbb{N}^\star} \sup_{k \in \{1,\ldots,N\}} \mathbb{E}\left[\left\|\mathbf{W}_t^{k,N}\right\|^{2m}\right] \leq C .$$

Let $V_m : \mathbb{R}^p \to \mathbb{R}_+$ given for any $w \in \mathbb{R}^p$ by $V_m(w) = 1 + \|w\|^{2m}$. For any $w \in \mathbb{R}^p$ we have

$$\|\nabla V_m(w)\| = 2m \|w\|^{2m-1} , \qquad \left\|\nabla^2 V_m(w)\right\| \leq 2m(2m-1) \|w\|^{2m-2} .$$

Combining this result with (36), the Cauchy-Schwarz inequality and the fact that for any $a, b \geq 0$ and $n_1, n_2 \in \mathbb{N}$, $a^{n_1} b^{n_2} \leq a^{n_1+n_2} + b^{n_1+n_2}$, we get that

$$|\langle \nabla V_m(w), b_N(t,w,\mu)\rangle| + |\langle \nabla^2 V_m(w), \sigma_N \sigma_N^\top(t,w,\mu)\rangle|$$

$$\lesssim \left[ 1 + \|w\| + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) \mathrm{d}\mu(\tilde{w}) \right)^{1/2} \right] \|\nabla V_m(w)\|$$

$$+ \left[ 1 + \|w\| + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) \mathrm{d}\mu(\tilde{w}) \right)^{1/2} \right]^2 \left\|\nabla^2 V_m(w)\right\|$$

$$\lesssim \left[ 1 + \|w\| + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) \mathrm{d}\mu(\tilde{w}) \right)^{1/2} \right] \|w\|^{2m-1}$$

$$+ \left[ 1 + \|w\|^2 + \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) \mathrm{d}\mu(\tilde{w}) \right] \|w\|^{2m-2}$$

$$\lesssim 1 + \|w\|^{2m} + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) \mathrm{d}\mu(\tilde{w}) \right)^m \lesssim 1 + \|w\|^{2m} + \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^{2m}) \mathrm{d}\mu(\tilde{w}) .$$
$$(37)$$

Now let $\tau_n^N = \inf\{t \geq 0 : \|\mathbf{W}_t^{k,N}\| \geq n \text{ for some } k \in \{1,\ldots,N\}\}$. Using Itô's lemma, (37) and (35), we have

$$\mathbb{E}\left[V_m(\mathbf{W}_{t \wedge \tau_n^N}^{k,N})\right] = \mathbb{E}\left[V_m(\mathbf{W}_{0 \wedge \tau_n^N}^{k,N})\right] + \mathbb{E}\left[\int_0^{t \wedge \tau_n^N} \langle \nabla V_m(\mathbf{W}_s^{k,N}), b_N(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N)\rangle \mathrm{d}s\right]$$

$$+ (1/2)\mathbb{E}\left[\int_0^{t \wedge \tau_n^N} \langle \nabla^2 V_m(\mathbf{W}_s^{k,N}), \sigma_N \sigma_N^\top(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N)\rangle \mathrm{d}s\right]$$

$$\lesssim \mathbb{E}\left[V_m(\mathbf{W}_{0 \wedge \tau_n^N}^{k,N})\right] + \mathbb{E}\left[\int_0^{t \wedge \tau_n^N} \left\{ V_m(\mathbf{W}_s^{k,N}) + (1/N) \sum_{j=1}^N V_m(\mathbf{W}_s^{j,N}) \right\} \mathrm{d}s\right]$$

Using Fatou's lemma, since almost surely $\tau_n^N \to +\infty$ as $n \to +\infty$, we get that

$$\mathbb{E}\left[V_m(\mathbf{W}_t^{k,N}) + (1/N) \sum_{j=1}^N V_m(\mathbf{W}_t^{j,N})\right]$$

$$\lesssim \mathbb{E}\left[V_m(\mathbf{W}_0^{k,N}) + (1/N) \sum_{j=1}^N V_m(\mathbf{W}_0^{j,N})\right] + \int_0^t \mathbb{E}\left[V_m(\mathbf{W}_s^{k,N}) + (1/N) \sum_{j=1}^N V_m(\mathbf{W}_s^{j,N})\right] \mathrm{d}s .$$

Using Grönwall's lemma, we get that for any $T \geq 0$, there exists $C \geq 0$ such that

$$\sup_{t \in [0,T]} \sup_{N \in \mathbb{N}^\star} \sup_{k \in \{1,\ldots,N\}} \mathbb{E}\left[\left\|\mathbf{W}_t^{k,N}\right\|^{2m}\right] \leq C .$$

We now show that there exists $C \geq 0$ such that

$$\sup_{N \in \mathbb{N}^\star} \sup_{k \in \{1,\ldots,N\}} \mathbb{E}\left[\sup_{t \in [0,T]} \left\|\mathbf{W}_t^{k,N}\right\|^{2m}\right] \leq C .$$

Using Jensen's inequality, Burkholder-Davis-Gundy's inequality [47, IV.42], (36) and the fact that for any $(a_j)_{j \in \{1,\ldots,M\}}$ and $r \geq 1$ such that $a_j \geq 0$, $(\sum_{j=1}^M a_j)^r \leq M^{r-1} \sum_{j=1}^M a_j^r$ we get for any $m \in \mathbb{N}^\star$

$$\begin{aligned}
&\mathbb{E}\left[\sup_{t \in [0,T]} \left\|\mathbf{W}_t^{k,N}\right\|^{2m}\right] \\
&\lesssim \mathbb{E}\left[\sup_{t \in [0,T]} \left\|\int_0^t b_N(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N) \mathrm{d}s\right\|^{2m}\right] + \mathbb{E}\left[\sup_{t \in [0,T]} \left\|\int_0^t \sigma_N^{1/2}(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N) \mathrm{d}\mathbf{B}_s\right\|^{2m}\right] \\
&\lesssim \mathbb{E}\left[\int_0^T \left\|b_N(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N)\right\|^{2m} \mathrm{d}s\right] + \mathbb{E}\left[\left(\int_0^T \mathrm{Tr}(\sigma_N \sigma_N^\top(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N)) \mathrm{d}s\right)^m\right] \\
&\lesssim \int_0^T \left\{\mathbb{E}\left[\left\|b_N(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N)\right\|^{2m}\right] + \mathbb{E}\left[\left\|\sigma_N(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N)\right\|^{2m}\right]\right\} \mathrm{d}s \\
&\lesssim \int_0^T \left\{1 + \mathbb{E}\left[\left\|\mathbf{W}_s^{k,N}\right\|^{2m}\right] + \mathbb{E}\left[\int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^{2m}) \mathrm{d}\boldsymbol{\nu}_s^N(\tilde{w})\right]\right\} \mathrm{d}s \\
&\lesssim \int_0^T \left\{1 + \mathbb{E}\left[\left\|\mathbf{W}_s^{k,N}\right\|^{2m}\right] + (1/N) \sum_{j=1}^N \mathbb{E}\left[\left\|\mathbf{W}_s^{j,N}\right\|^{2m}\right]\right\} \mathrm{d}s \\
&\lesssim 1 + \sup_{N \in \mathbb{N}^\star} \sup_{j \in \{1,\ldots,N\}} \sup_{t \in [0,T]} \mathbb{E}\left[\left\|\mathbf{W}_s^{j,N}\right\|^{2m}\right] ,
\end{aligned}$$

which concludes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

## 9.2 Existence of solutions to the mean-field SDE

The following result is based on [33, Theorem 1.1] showing, under **B**1 and **B**2, the existence of strong solutions and pathwise uniqueness for non-homogeneous McKean-Vlasov SDE with non-constant covariance matrix:

$$\mathrm{d}\mathbf{W}_t^\star = b(t, \mathbf{W}_t^\star, \boldsymbol{\lambda}_t^\star)\mathrm{d}t + \sigma(t, \mathbf{W}_t^\star, \boldsymbol{\lambda}_t^\star)\mathrm{d}\mathbf{B}_t , \qquad (38)$$

where $b$ and $\sigma$ are given in **B**2 and where for any $t \geq 0$, $\mathbf{W}_t^\star$ has distribution $\boldsymbol{\lambda}_t^\star \in \mathscr{P}_2(\mathbb{R}^p)$, $(\mathbf{B}_t)_{t \geq 0}$ is a $r$ dimensional Brownian motion and $\mathbf{W}_0^\star$ has distribution $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$.

**Proposition 11.** *Assume* **B**1 *and* **B**2*. Let* $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$*. Then, there exists a* $(\mathcal{F}_t)_{t \geq 0}$*-adapted continuous process* $(\mathbf{W}_t^\star)_{t \geq 0}$ *which is the unique strong solution of* (38) *satisfying for any* $T \geq 0$, $\sup_{t \in [0,T]} \mathbb{E}[\|\mathbf{W}_t^\star\|^2] < +\infty$.

*Proof.* Let $\delta \geq 0$ and $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$. Note that we only need to show that (38) admits a strong solution up to $\bar{\delta} > 0$. First, using [46, Theorem 2.9], note that for any $(\boldsymbol{\mu}_t)_{t \in [0,\delta]} \in \mathscr{C}_{2,\delta}^p$ the SDE,

$$\mathrm{d}\mathbf{W}_t^{\boldsymbol{\mu}} = b(t, \mathbf{W}_t^{\boldsymbol{\mu}}, \boldsymbol{\mu}_t)\mathrm{d}t + \sigma(t, \mathbf{W}_t^{\boldsymbol{\mu}}, \boldsymbol{\mu}_t)\mathrm{d}\mathbf{B}_t ,$$

admits a unique strong solution, since for any $t \in [0,\delta]$ and $w_1, w_2 \in \mathbb{R}^p$

$$\|b(t, w_1, \boldsymbol{\mu}_t) - b(t, w_2, \boldsymbol{\mu}_t)\| + \|\sigma(t, w_1, \boldsymbol{\mu}_t) - \sigma(t, w_2, \boldsymbol{\mu}_t)\| \leq \mathtt{M}_1 \|w_1 - w_2\| . \qquad (39)$$

In addition, $\sup_{t \in [0,\delta]} \mathbb{E}[\|\mathbf{W}_t^{\boldsymbol{\mu}}\|^2] < +\infty$.

In the rest of the proof, the strategy is to adapt the well-known Cauchy-Lipschitz approach using the Picard fixed point theorem. More precisely, we define below for $\delta > 0$ small enough, a contractive mapping $\mathbf{\Phi}_\delta : \mathscr{C}_{2,\delta}^p \to \mathscr{C}_{2,\delta}^p$ such that the unique fixed point $(\boldsymbol{\lambda}_t^\star)_{t \in [0,\delta]}$ is a weak solution of (38). Considering $(\mathbf{W}_t^{\boldsymbol{\lambda}^\star})_{t \in [0,\delta]}$, we obtain the unique strong solution of (38) on $[0,\delta]$.

Let $\delta > 0$. Denote $(\boldsymbol{\lambda}_t^{\boldsymbol{\mu}})_{t \in [0,\delta]} \in \mathscr{P}_2(\mathbb{R}^p)^{[0,\delta]}$ such that for any $t \in [0,\delta]$, $\boldsymbol{\lambda}_t^{\boldsymbol{\mu}}$ is the distribution of $\mathbf{W}_t^{\boldsymbol{\mu}}$ with initial condition $\mathbf{W}_0^\star$ with distribution $\boldsymbol{\lambda}_0^{\boldsymbol{\mu}} = \mu_0$. In addition, using (10), (36), (39), **B**1, **B**2, the Cauchy-Schwarz inequality, the Itô isometry and the fact that for any $a, b \geq 0$, $2ab \leq a^2 + b^2$, there exists $C \geq 0$ such that for any $t, s \in [0,\delta]$ with $t \geq s$,

$$
\mathcal{W}_2(\boldsymbol{\lambda}_t^{\boldsymbol{\mu}}, \boldsymbol{\lambda}_s^{\boldsymbol{\mu}})^2 \leq \mathbb{E}\left[\|\mathbf{W}_t^{\boldsymbol{\mu}} - \mathbf{W}_s^{\boldsymbol{\mu}}\|^2\right]
$$

$$
\leq 2\mathbb{E}\left[\left\|\int_s^t b(u, \mathbf{W}_u^{\boldsymbol{\mu}}, \boldsymbol{\mu}_u)\mathrm{d}u\right\|^2\right] + 2\mathbb{E}\left[\left\|\int_s^t \sigma(u, \mathbf{W}_u^{\boldsymbol{\mu}}, \boldsymbol{\mu}_u)\mathrm{d}\mathbf{B}_u\right\|^2\right]
$$

$$
\leq 2(t-s)\int_s^t \mathbb{E}\left[\|b(u, \mathbf{W}_u^{\boldsymbol{\mu}}, \boldsymbol{\mu}_u)\|^2\right]\mathrm{d}u + 2\int_s^t \mathbb{E}\left[\mathrm{Tr}(\sigma\sigma^\top(u, \mathbf{W}_u^{\boldsymbol{\mu}}, \boldsymbol{\mu}_u))\right]\mathrm{d}u
$$

$$
\leq 4(t-s)\int_s^t \left\{\|b(u, 0, \boldsymbol{\mu}_u)\|^2 + \mathtt{M}_1^2\mathbb{E}\left[\|\mathbf{W}_u^{\boldsymbol{\mu}}\|^2\right]\right\}\mathrm{d}u
$$

$$
+ 4\int_s^t \left\{\|\sigma(u, 0, \boldsymbol{\mu}_u)\|^2 + \mathtt{M}_1^2\mathbb{E}\left[\|\mathbf{W}_u^{\boldsymbol{\mu}}\|^2\right]\right\}\mathrm{d}u
$$

$$
\leq 4(1+\delta)(t-s)\left[\mathtt{M}_1^2 \sup_{t \in [0,\delta]} \mathbb{E}[\|\mathbf{W}_t^{\boldsymbol{\mu}}\|^2] + \sup_{t \in [0,\delta]}\left\{\|b(t, 0, \boldsymbol{\mu}_t)\|^2 + \|\sigma(t, 0, \boldsymbol{\mu}_t)\|^2\right\}\right]
$$

$$
\leq C(t-s)\{1 + \sup_{t \in [0,\delta]} \mathbb{E}[\|\mathbf{W}_t^{\boldsymbol{\mu}}\|^2]\} .
$$

Therefore, $(\boldsymbol{\lambda}_t^{\boldsymbol{\mu}})_{t \in [0,\delta]} \in \mathscr{C}_{2,\delta}^p$. Let $\mathbf{\Phi}_\delta : \mathscr{C}_{2,\delta}^p \to \mathscr{C}_{2,\delta}^p$ given for any $(\boldsymbol{\mu}_t)_{t \in [0,\delta]} \in \mathscr{C}_{2,\delta}^p$ by $\mathbf{\Phi}_\delta((\boldsymbol{\mu}_t)_{t \in [0,\delta]}) = (\boldsymbol{\lambda}_t^{\boldsymbol{\mu}})_{t \in [0,\delta]}$. Let $(\boldsymbol{\mu}_{1,t})_{t \in [0,\delta]}, (\boldsymbol{\mu}_{2,t})_{t \in [0,\delta]} \in \mathscr{C}_{2,\delta}^p$, using (10), (39), **B**1, **B**2, the Cauchy-Schwarz inequality, the Itô isometry, the fact that for any $a, b \geq 0$, $2ab \leq a^2 + b^2$ and Grönwall's inequality we have for any $t \in [0,\delta]$

$$
\mathbb{E}\left[\|\mathbf{W}_t^{\boldsymbol{\mu}_1} - \mathbf{W}_t^{\boldsymbol{\mu}_2}\|^2\right] \leq 2\mathbb{E}\left[\left\|\int_0^t \{b(s, \mathbf{W}_s^{\boldsymbol{\mu}_1}, \boldsymbol{\mu}_{1,s}) - b(s, \mathbf{W}_s^{\boldsymbol{\mu}_2}, \boldsymbol{\mu}_{2,s})\}\mathrm{d}s\right\|^2\right]
$$

$$
+ 2\mathbb{E}\left[\left\|\int_0^t \{\sigma(s, \mathbf{W}_s^{\boldsymbol{\mu}_1}, \boldsymbol{\mu}_{1,s}) - \sigma(s, \mathbf{W}_s^{\boldsymbol{\mu}_2}, \boldsymbol{\mu}_{2,s})\}\mathrm{d}\mathbf{B}_s\right\|^2\right]
$$

$$
\leq 2\delta\int_0^t \mathbb{E}\left[\|b(s, \mathbf{W}_s^{\boldsymbol{\mu}_1}, \boldsymbol{\mu}_{1,s}) - b(s, \mathbf{W}_s^{\boldsymbol{\mu}_2}, \boldsymbol{\mu}_{2,s})\|^2\right]\mathrm{d}s
$$

$$
+ 2\int_0^t \mathbb{E}\left[\|\sigma(s, \mathbf{W}_s^{\boldsymbol{\mu}_1}, \boldsymbol{\mu}_{1,s}) - \sigma(s, \mathbf{W}_s^{\boldsymbol{\mu}_2}, \boldsymbol{\mu}_{2,s})\|^2\right]\mathrm{d}s
$$

$$
\leq 4\mathtt{M}_1^2(1+\delta)\int_0^t \left\{\mathbb{E}\left[\|\mathbf{W}_s^{\boldsymbol{\mu}_1} - \mathbf{W}_s^{\boldsymbol{\mu}_2}\|^2\right] + \int_\mathsf{Z} \zeta^2(z)\mathrm{d}\pi_\mathsf{Z}(z)\mathcal{W}_2^2(\boldsymbol{\mu}_{1,s}, \boldsymbol{\mu}_{2,s})\right\}\mathrm{d}s
$$

$$
\leq 4\mathtt{M}_1^2\delta(1+\delta)\int_\mathsf{Z} \zeta^2(z)\mathrm{d}\pi_\mathsf{Z}(z)\mathcal{W}_{2,\delta}^2(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2) + 4\mathtt{M}_1^2(1+\delta)\int_0^t \mathbb{E}\left[\|\mathbf{W}_s^{\boldsymbol{\mu}_1} - \mathbf{W}_s^{\boldsymbol{\mu}_2}\|^2\right]\mathrm{d}s
$$

$$
\leq 4\mathtt{M}_1^2\delta(1+\delta)\exp\left[4\mathtt{M}_1^2(1+\delta)\delta\int_\mathsf{Z} \zeta^2(z)\mathrm{d}\pi_\mathsf{Z}(z)\right]\mathcal{W}_{2,\delta}^2(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2) .
$$

Using this result, we obtain that for any $(\boldsymbol{\mu}_{1,t})_{t \in [0,\delta]}, (\boldsymbol{\mu}_{2,t})_{t \in [0,\delta]} \in \mathrm{C}([0,\delta], \mathscr{P}_2(\mathbb{R}^p))$,

$$
\mathcal{W}_{2,\delta}^2(\mathbf{\Phi}_\delta(\boldsymbol{\mu}_1), \mathbf{\Phi}_\delta(\boldsymbol{\mu}_2)) \leq \sup_{t \in [0,\delta]} \mathbb{E}\left[\|\mathbf{W}_t^{\boldsymbol{\mu}_1} - \mathbf{W}_t^{\boldsymbol{\mu}_2}\|^2\right]
$$

$$
\leq 4\mathtt{M}_1^2\delta(1+\delta)\exp\left[4\mathtt{M}_1^2(1+\delta)\delta\int_\mathsf{Z} \zeta^2(z)\mathrm{d}\pi_\mathsf{Z}(z)\right]\mathcal{W}_{2,\delta}^2(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2) .
$$

Hence, for $\delta > 0$ small enough, $\boldsymbol{\Phi}_\delta$ is contractive and since $\mathrm{C}([0,\delta], \mathscr{P}_2(\mathbb{R}^p))$ is a complete metric space, we get, using Picard fixed point theorem, that there exists a unique $(\boldsymbol{\lambda}_t^\star)_{t\in[0,\delta]} \in \mathrm{C}([0,\delta], \mathscr{P}_2(\mathbb{R}^p))$ such that, $\boldsymbol{\Phi}_\delta(\boldsymbol{\lambda}^\star) = \boldsymbol{\lambda}^\star$. For this $\boldsymbol{\lambda}^\star$, we have that $(\mathbf{W}_t^{\boldsymbol{\lambda}^\star})_{t\in[0,\delta]}$ is a strong solution to (38). We have shown that (38) admits a strong solution for any initial condition $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$.

We now show that pathwise uniqueness holds for (38). Let $(\mathbf{W}_t^1)_{t\in[0,\delta]}$ and $(\mathbf{W}_t^2)_{t\in[0,\delta]}$ be two strong solutions of (38) such that $\mathbf{W}_0^1 = \mathbf{W}_0^2 = w_0 \in \mathbb{R}^p$. Let, $(\boldsymbol{\mu}_{1,t})_{t\in[0,\delta]}$ and $(\boldsymbol{\mu}_{2,t})_{t\in[0,\delta]}$ such that for any $t \in [0,\delta]$, $\boldsymbol{\mu}_{1,t}$ is the distribution of $\mathbf{W}_t^1$ and $\boldsymbol{\mu}_{2,t}$ the one of $\mathbf{W}_t^2$. Since $\boldsymbol{\Phi}_\delta$ admits a unique fixed point, we get that $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$. Hence, $(\mathbf{W}_t^1)_{t\in[0,\delta]}$ and $(\mathbf{W}_t^2)_{t\in[0,\delta]}$ are strong solutions of (39) with $\boldsymbol{\mu} \leftarrow \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$ and since pathwise uniqueness holds for (39), we get that $(\mathbf{W}_t^1)_{t\in[0,\delta]} = (\mathbf{W}_t^1)_{t\in[0,\delta]}$. $\qquad\square$

## 9.3 Main result

**Theorem 12.** *Assume* **B**1 *and* **B**2. *For any* $N \in \mathbb{N}^\star$, *let* $(\mathbf{W}_t^{1:N})_{t\geq 0}$ *be a strong solution of* (35) *and for any* $N \in \mathbb{N}^\star$ *and* $k \in \{1,\ldots,N\}$, *let* $(\mathbf{W}_t^{k,\star})_{t\geq 0}$ *be a strong solution of* (38) *with Brownian motion* $(\mathbf{B}_t^k)_{t\geq 0}$. *Assume that there exists* $\mu_0 \in \mathscr{P}_2(\mathbb{R}^p)$ *such that for any* $N \in \mathbb{N}^\star$, $\mathbf{W}_0^{1:N} = \mathbf{W}_0^{\star,1:N}$ *has distribution* $\mu_0^{\otimes N}$. *Then for any* $T \geq 0$, $N \in \mathbb{N}^\star$ *and* $k \in \{1,\ldots,N\}$

$$
\mathbb{E}\left[\sup_{t\in[0,T]}\left\|\mathbf{W}_t^{k,N} - \mathbf{W}_t^{k,\star}\right\|^2\right] \leq 32(1+T)^2\left(1+\int_{\mathsf{Z}}\zeta^2(z)\mathrm{d}\pi_{\mathsf{Z}}(z)\right)\left(\mathtt{M}_2^2 N^{-2\kappa} + \mathtt{M}_1^2 N^{-1}\right)
$$
$$
\times \exp\left[16(1+T)^2\left(1+\int_{\mathsf{Z}}\zeta^2(z)\mathrm{d}\pi_{\mathsf{Z}}(z)\right)\mathtt{M}_1^2\right].
$$

*Proof.* Let $T \geq 0$. For any $N \in \mathbb{N}^\star$, $t \geq 0$, let $\boldsymbol{\nu}_t^{\star,N} = (1/N)\sum_{j=1}^N \delta_{\mathbf{W}_s^{\star,j}}$. Using **B**1, **B**2, Itô's isometry, Doob's inequality, Jensen's inequality and the fact that for any $a,b \geq 0$, $(a+b)^2 \leq 2(a^2 + b^2)$, we have for any $N \in \mathbb{N}^\star$ and $k \in \{1,\ldots,N\}$

$$
\mathbb{E}\left[\sup_{t\in[0,T]}\left\|\mathbf{W}_t^{k,N} - \mathbf{W}_t^{k,\star}\right\|^2\right] \leq 2\mathbb{E}\left[\sup_{t\in[0,T]}\left\|\int_0^t \left(b_N(s,\mathbf{W}_s^{k,N},\boldsymbol{\nu}_s^N) - b(s,\mathbf{W}_s^{k,\star},\boldsymbol{\lambda}_s^\star)\right)\mathrm{d}s\right\|^2\right]
$$
$$
+ 2\mathbb{E}\left[\sup_{t\in[0,T]}\left\|\int_0^t \left(\sigma_N(s,\mathbf{W}_s^{k,N},\boldsymbol{\nu}_s^N) - \sigma(s,\mathbf{W}_s^{k,\star},\boldsymbol{\lambda}_s^\star)\right)\mathrm{d}\mathbf{B}_s^{k,N}\right\|^2\right]
$$
$$
\leq 2T\int_0^T \mathbb{E}\left[\left\|b_N(s,\mathbf{W}_s^{k,N},\boldsymbol{\nu}_s^N) - b(s,\mathbf{W}_s^{k,\star},\boldsymbol{\lambda}_s^\star)\right\|^2\right]\mathrm{d}s
$$
$$
+ 2\mathbb{E}\left[\left\|\int_0^T \left(\sigma_N(s,\mathbf{W}_s^{k,N},\boldsymbol{\nu}_s^N) - \sigma(s,\mathbf{W}_s^{k,\star},\boldsymbol{\lambda}_s^\star)\right)\mathrm{d}\mathbf{B}_s^{k,N}\right\|^2\right]
$$
$$
\leq 2(1+T)\int_0^T \left\{\mathbb{E}\left[\left\|b_N(s,\mathbf{W}_s^{k,N},\boldsymbol{\nu}_s^N) - b(s,\mathbf{W}_s^{k,\star},\boldsymbol{\lambda}_s^\star)\right\|^2\right]\right.
$$
$$
\left. + \mathbb{E}\left[\left\|\sigma_N(s,\mathbf{W}_s^{k,N},\boldsymbol{\nu}_s^N) - \sigma(s,\mathbf{W}_s^{k,\star},\boldsymbol{\lambda}_s^\star)\right\|^2\right]\right\}\mathrm{d}s
$$
$$
\leq 8\mathtt{M}_2^2(1+T)^2 N^{-2\kappa} + 4(1+T)\int_0^T \left\{\mathbb{E}\left[\left\|b(s,\mathbf{W}_s^{k,N},\boldsymbol{\nu}_s^N) - b(s,\mathbf{W}_s^{k,\star},\boldsymbol{\lambda}_s^\star)\right\|^2\right]\right.
$$
$$
\left. + \mathbb{E}\left[\left\|\sigma(s,\mathbf{W}_s^{k,N},\boldsymbol{\nu}_s^N) - \sigma(s,\mathbf{W}_s^{k,\star},\boldsymbol{\lambda}_s^\star)\right\|^2\right]\right\}\mathrm{d}s
$$
$$
\leq 8\mathtt{M}_2^2(1+T)^2 N^{-2\kappa} + 8\mathtt{M}_1^2(1+T)
$$
$$
\times \int_0^T \left\{\int_{\mathsf{Z}}\mathbb{E}\left[\left\|\boldsymbol{\nu}_s^N[\mathsf{g}(\cdot,z)] - \boldsymbol{\lambda}_s^\star[\mathsf{g}(\cdot,z)]\right\|^2\right]\mathrm{d}\pi_{\mathsf{Z}}(z) + \mathbb{E}\left[\left\|\mathbf{W}_s^{k,N} - \mathbf{W}_s^{k,\star}\right\|^2\right]\right\}\mathrm{d}s
$$
$$
\leq 8\mathtt{M}_2^2(1+T)^2 N^{-2\kappa} + 16\mathtt{M}_1^2(1+T)
$$

$$\times \int_0^T \left\{ \int_{\mathsf{Z}} \left( \mathbb{E}\left[ \left\| \boldsymbol{\nu}_s^N[\mathrm{g}(\cdot, z)] - \boldsymbol{\nu}_s^{\star,N}[\mathrm{g}(\cdot,z)] \right\|^2 \right] + \mathbb{E}\left[ \left\| \boldsymbol{\nu}_s^{\star,N}[\mathrm{g}(\cdot,z)] - \boldsymbol{\lambda}_s^{\star}[\mathrm{g}(\cdot,z)] \right\|^2 \right] \right) \mathrm{d}\pi_{\mathsf{Z}}(z) \right.$$
$$\left. + \mathbb{E}\left[ \left\| \mathbf{W}_s^{k,N} - \mathbf{W}_s^{k,\star} \right\|^2 \right] \right\} \mathrm{d}s \ .$$

Then using the Cauchy-Schwarz's inequality, the fact that $\{(\mathbf{W}_t^{k,N})_{t \geq 0}\}_{k=1}^N$ are exchangeable, *i.e.* for any permutation $\tau : \{1, \ldots, N\} \to \{1, \ldots, N\}$, $\{(\mathbf{W}_t^{k,N})_{t \geq 0}\}_{k=1}^N$ has the same distribution as $\{(\mathbf{W}_t^{\tau(k),N})_{t \geq 0}\}_{k=1}^N$ and $\{(\mathbf{W}_t^{k,\star})_{t \geq 0}\}_{k=1}^N$ are independent we have

$$\mathbb{E}\left[ \sup_{t \in [0,T]} \left\| \mathbf{W}_t^{k,N} - \mathbf{W}_t^{k,\star} \right\|^2 \right] \leq 8\mathtt{M}_2^2(1+T)^2 N^{-2\kappa} + 16\mathtt{M}_1^2(1+T)$$

$$\times \int_0^T \left\{ \frac{1}{N} \int_{\mathsf{Z}} \zeta^2(z) \mathrm{d}\pi_{\mathsf{Z}}(z) \sum_{j=1}^N \mathbb{E}\left[ \left\| \mathbf{W}_s^{j,N} - \mathbf{W}_s^{j,\star} \right\|^2 \right] + \mathbb{E}\left[ \left\| \mathbf{W}_s^{k,N} - \mathbf{W}_s^{k,\star} \right\|^2 \right] \right.$$

$$\left. + \int_{\mathsf{Z}} \mathbb{E}\left[ \left\| \frac{1}{N} \sum_{j=1}^N \mathrm{g}(\mathbf{W}_s^{j,\star}, z) - \int_{\mathbb{R}^p} \mathrm{g}(\bar{w}, z) \mathrm{d}\boldsymbol{\lambda}_s^{\star}(\bar{w}) \right\|^2 \right] \mathrm{d}\pi_{\mathsf{Z}}(z) \right\} \mathrm{d}s$$

$$\leq 8\mathtt{M}_2^2(1+T)^2 N^{-2\kappa} + 16\mathtt{M}_1^2(1+T)\left( 1 + \int_{\mathsf{Z}} \zeta^2(z) \mathrm{d}\pi_{\mathsf{Z}}(z) \right) \int_0^T \mathbb{E}\left[ \left\| \mathbf{W}_s^{k,N} - \mathbf{W}_s^{k,\star} \right\|^2 \right] \mathrm{d}s$$

$$+ 16\mathtt{M}_1^2(1+T) N^{-1} \int_0^T \int_{\mathsf{Z}} \mathbb{E}\left[ \left\| \mathrm{g}(\mathbf{W}_s^{k,\star}, z) - \int_{\mathbb{R}^p} \mathrm{g}(\bar{w}, z) \mathrm{d}\boldsymbol{\lambda}_s^{\star}(\bar{w}) \right\|^2 \right] \mathrm{d}\pi_{\mathsf{Z}}(z) \mathrm{d}s$$

$$\leq 8\mathtt{M}_2^2(1+T)^2 N^{-2\kappa} + 16\mathtt{M}_1^2(1+T)\left( 1 + \int_{\mathsf{Z}} \zeta^2(z) \mathrm{d}\pi_{\mathsf{Z}}(z) \right) \int_0^T \mathbb{E}\left[ \left\| \mathbf{W}_s^{k,N} - \mathbf{W}_s^{k,\star} \right\|^2 \right] \mathrm{d}s$$

$$+ 32\mathtt{M}_1^2(1+T)^2 N^{-1}\left( 1 + \int_{\mathsf{Z}} \zeta^2(z) \mathrm{d}\pi_{\mathsf{Z}}(z) \right) \ .$$

We conclude the proof upon combining this result and Grönwall's inequality. $\qquad \square$

### 9.4  Proofs of the main results

In this section we prove Theorem 1, Theorem 2, Theorem 7, Theorem 8. Note that we only need to show Theorem 7 and Theorem 8, since in the case $\eta = 0$, Theorem 7 boils down to Theorem 1 and Theorem 8 to Theorem 2.

*Proof of Theorem 7.* Define for any $N \in \mathbb{N}^\star$, $w \in \mathbb{R}^p$, $\mu \in \mathscr{P}_2(\mathbb{R}^p)$ and $t \geq 0$

$$b_N(t, w, \mu) = (t+1)^{-\alpha} h(w, \mu) \, , \ \sigma_N(t, w, \mu) = (t+1)^{-\alpha}((\gamma_{\alpha,\beta}(N)/M)^{1/2} \Sigma^{1/2}(w, \mu), \sqrt{2}\,\mathrm{Id}) \, ,$$
$$b(t, w, \mu) = (t+1)^{-\alpha} h(w, \mu) \, , \quad \sigma(t, w, \mu) = (t+1)^{-\alpha}(0, \sqrt{2}\,\mathrm{Id}) \, ,$$

with $h$ and $\Sigma$ given in (16). Using Proposition 9, we get that **B**1 holds with $\mathtt{M}_1 \leftarrow \mathtt{L}$ and $\gamma_{\alpha,\beta}(N) = \gamma^{1/(1-\alpha)} N^{(\beta-1)/(1-\alpha)}$. In addition, using Proposition 9, **B**2 holds with $\mathtt{M}_2 \leftarrow (\gamma^{1-\alpha}/M)^{1/2} p\mathtt{L}$ and $2\kappa = (1-\beta)/(1-\alpha)$. We conclude using Theorem 12. $\qquad \square$

*Proof of Theorem 8.* Define for any $N \in \mathbb{N}^\star$, $w \in \mathbb{R}^p$, $\mu \in \mathscr{P}_2(\mathbb{R}^p)$ and $t \geq 0$

$$b_N(t, w, \mu) = (t+1)^{-\alpha} h(w, \mu) \, , \ \sigma_N(t, w, \mu) = (t+1)^{-\alpha}((\gamma^{1/(1-\alpha)}/M)^{1/2} \Sigma^{1/2}(w, \mu), \sqrt{2}\,\mathrm{Id}) \, ,$$

with $h$ and $\Sigma$ given in (16). Using Proposition 9, we get that **B**1 holds with $\mathtt{M}_1 \leftarrow \mathtt{L}$. In addition, **B**2 holds with $b = b_N$, $\sigma = \sigma_N$, $\mathtt{M}_2 \leftarrow 0$ and $\kappa = 0$. We conclude using Theorem 12. $\qquad \square$

*Proof of Proposition 4.* We consider only the case $\beta = 1$, the proof for $\beta \in [0, 1)$ following the same lines. Let $M \in \mathbb{N}^\star$. We have for any $N \in \mathbb{N}^\star$ using Proposition 6,

$$\mathcal{W}_2(\boldsymbol{\Upsilon}^N, \delta_{\boldsymbol{\lambda}^\star})^2 \leq \mathbb{E}\left[ \mathcal{W}_2(\boldsymbol{\nu}^N, \boldsymbol{\lambda}^\star)^2 \right]$$

$$\leq N^{-1} \sum_{k=1}^N \mathbb{E}\left[ \mathcal{W}_2(\delta_{(\mathbf{W}_t^{k,N})_{t \geq 0}}, \boldsymbol{\lambda}^\star)^2 \right] \leq N^{-1} \sum_{k=1}^N \mathbb{E}\left[ \mathrm{m}^2((\mathbf{W}_t^{k,N})_{t \geq 0}, (\mathbf{W}_t^{k,\star})_{t \geq 0}) \right] \ . \quad (40)$$

Let $\varepsilon > 0$ and $n_0$ such that $\sum_{n=n_0+1}^{+\infty} 2^{-n} \leq \varepsilon$. Combining (40), Theorem 1 and the Cauchy-Schwarz inequality we get that for any $N \in \mathbb{N}^\star$

$$\mathcal{W}_2(\mathbf{\Upsilon}^N, \delta_{\boldsymbol{\lambda}^\star})^2 \leq 2\varepsilon^2 + \frac{2n_0}{N} \sum_{k=1}^{N} \sum_{n=1}^{n_0} \mathbb{E}\left[\sup_{t \in [0,n]} \|\mathbf{W}_t^{k,N} - \mathbf{W}_t^{k,\star}\|^2\right] \leq 2\varepsilon^2 + 2n_0 N^{-1} \sum_{n=0}^{n_0} C_{1,n} .$$

Therefore, for any $\varepsilon > 0$ there exists $N_0 \in \mathbb{N}^\star$ such that for any $N \in \mathbb{N}^\star$ with $N \geq N_0$, $\mathcal{W}_2(\mathbf{\Upsilon}^N, \delta_{\boldsymbol{\lambda}^\star}) \leq \varepsilon$, which concludes the proof. □

## 10 Existence of invariant measure in the one-dimensional case

In this section we prove Proposition 5.

*Proof of Proposition 5.* Since $V$ is $\eta$-strongly convex it admits a unique minimum at $w_0 \in \mathbb{R}$. Using **A**1-(c), the fact that $V$ is $\eta$-strongly convex and [48, Theorem 2.1.5, Theorem 2.1.7] there exists $\mathtt{M} \geq 0$ such that for any $w \in \mathbb{R}$ we have

$$\eta(w - w_0)^2/2 \leq V(w) - V(w_0) \leq \mathtt{M}(w - w_0)^2/2 . \tag{41}$$

In addition, using Proposition 9, we have for any $\mu \in \mathscr{P}_2(\mathbb{R})$ and $w \in \mathbb{R}$,

$$\bar{\sigma}^2 \leq \Sigma(w, \mu) \leq \mathtt{L}^2 . \tag{42}$$

Recall that for any $\mu \in \mathscr{P}_2(\mathbb{R})$ and $w \in \mathbb{R}$, $h(w, \mu) = \bar{h}(w, \mu) + V'(w)$, with $\bar{h}$ given in (16). Note that for any $w \in [w_0, +\infty)$, $V'(w) \geq 0$ and for any $w \in (-\infty, w_0]$, $V'(w) \leq 0$. Combining this result, Proposition 9, (41) and (42), there exists $\mathtt{m}_1 > 0$ and $c_1 \in \mathbb{R}$ such that for any $\mu \in \mathscr{P}_2(\mathbb{R})$ and $w \in \mathbb{R}$, we have distinguishing the case $w \leq w_0$ and $w > w_0$,

$$\int_0^w \{h/\Sigma\}(\tilde{w}, \mu)\mathrm{d}w \geq -\bar{\sigma}^{-2}\mathtt{L}^2 |w| + \int_0^w V'(\tilde{w})/\Sigma(\tilde{w}, \mu)\mathrm{d}\tilde{w}$$

$$\geq -\bar{\sigma}^{-2}\mathtt{L}^2 |w| - \bar{\sigma}^{-2} \sup_{\tilde{w} \in [0,w_0]} |V'(\tilde{w})| \, |w_0| + \int_{w_0}^w V'(\tilde{w})/\Sigma(\tilde{w}, \mu)\mathrm{d}\tilde{w}$$

$$\geq -\bar{\sigma}^{-2}\mathtt{L}^2 |w| - \bar{\sigma}^{-2} \sup_{\tilde{w} \in [0,w_0]} |V'(\tilde{w})| \, |w_0| + (V(w) - V(w_0))\mathtt{L}^{-2} \geq \mathtt{m}_1 w^2 + c_1 . \tag{43}$$

Therefore, we obtain that for any $\mu \in \mathscr{P}_2(\mathbb{R})$, $\int_{\mathbb{R}} \exp[\int_0^w h(\tilde{w}, \mu)/\Sigma(\tilde{w}, \mu)\mathrm{d}\tilde{w}]\mathrm{d}w < +\infty$. Define $H : \mathscr{P}_2(\mathbb{R}) \to \mathscr{P}_2(\mathbb{R})$ such that for any $\mu \in \mathscr{P}_2(\mathbb{R})$, $H(\mu)$ is the probability measure with density $\rho_\mu$ given for any $w \in \mathbb{R}$ by

$$\rho_\mu(w) \propto \bar{\Sigma}^{-1}(w, \mu) \exp\left[-2 \int_0^w h(\tilde{w}, \mu)/\bar{\Sigma}(\tilde{w}, \mu)\mathrm{d}\tilde{w}\right] ,$$

where $\bar{\Sigma}(w, \mu) = \gamma^{1/(1-\alpha)}\Sigma(w, \mu)/M$. Similarly to (43), there exist $\mathtt{m}_2 > 0$ and $c_2 \in \mathbb{R}$ such that for any $\mu \in \mathscr{P}_2(\mathbb{R})$ and $w \in \mathbb{R}$

$$\int_0^w h(\tilde{w}, \mu)/\Sigma(\tilde{w}, \mu)\mathrm{d}\tilde{w} \leq \mathtt{m}_2 w^2 + c_2 . \tag{44}$$

Combining (42), (43) and (44), there exists $\mathtt{m} > 0$ and $c \in \mathbb{R}$ such that for any $\mu \in \mathscr{P}_2(\mathbb{R})$ and $w \in \mathbb{R}$, $\rho_\mu(w) \leq c e^{-\mathtt{m}w^2}$. Using this result, we get that $\sup_{\mu \in \mathscr{P}_2(\mathbb{R})} \int_{\mathbb{R}} w^4 \rho_\mu(w)\mathrm{d}w < +\infty$. Therefore, using [49, Theorem 2.7] we obtain that $H(\mathscr{P}_2(\mathbb{R}))$ is relatively compact in $(\mathscr{P}_2(\mathbb{R}), \mathcal{W}_2)$.

We now show that $H \in \mathrm{C}(\mathscr{P}_2(\mathbb{R}), \mathscr{P}_2(\mathbb{R}))$. Let $\mu \in \mathscr{P}_2(\mathbb{R})$ and $(\mu_n)_{n \in \mathbb{N}} \in \mathscr{P}_2(\mathbb{R})^{\mathbb{N}}$ such that $\lim_{n \to +\infty} \mu_n = \mu$. Using Proposition 9 and the Lebesgue dominated convergence theorem we obtain that for any $w \in \mathbb{R}$, $\lim_{n \to +\infty} \rho_{\mu_n}(w) = \rho_\mu(w)$. Using Scheffé's lemma we get that $\lim_{n \to +\infty} \int_{\mathbb{R}} |\rho_{\mu_n}(w) - \rho_\mu(w)| \, \mathrm{d}w = 0$. Hence, $(H(\mu_n))_{n \in \mathbb{N}}$ weakly converges towards $H(\mu)$.

Let $(H(\mu_{n_k}))_{k \in \mathbb{N}}$ be a converging sequence in $(\mathscr{P}_2(\mathbb{R}), \mathcal{W}_2)$. Therefore, $(H(\mu_{n_k}))_{k \in \mathbb{N}}$ also weakly converges and we obtain that $\lim_{k \to +\infty} \mathcal{W}_2(H(\mu_{n_k}), H(\mu)) = 0$. Since $\{H(\mu_n) : n \in \mathbb{N}\}$ is relatively compact and admits a unique limit point we obtain that $\lim_{n \to +\infty} \mathcal{W}_2(H(\mu_n), H(\mu)) = 0$.

Hence $H \in \mathrm{C}(\mathscr{P}_2(\mathbb{R}), \mathscr{P}_2(\mathbb{R}))$. Therefore, since $H \in \mathrm{C}(\mathscr{P}_2(\mathbb{R}), \mathscr{P}_2(\mathbb{R}))$ and $H(\mathscr{P}_2(\mathbb{R}))$ is relatively compact in $\mathscr{P}_2(\mathbb{R})$ Schauder's theorem [50, Appendix] implies that $H$ admits a fixed point.

Let $\mu \in \mathscr{P}_2(\mathbb{R})$ be a fixed point of $H$. We now show that $\mu$ is an invariant probability distribution for (8). Let $(\mathbf{W}_t^\mu)_{t \geq 0}$ such that $\mathbf{W}_0^\mu$ has distribution $\mu$ and strong solution to the following SDE

$$\mathrm{d}\mathbf{W}_t^\mu = h(t, \mu)\mathrm{d}t + \gamma^{1/(1-\alpha)}\Sigma(\mathbf{W}_t^\mu, \mu)\mathrm{d}\mathbf{B}_t \ . \tag{45}$$

An invariant distribution for (45) is given by $H(\mu)$, see [51]. Hence, since $\mu = H(\mu)$, for any $t \geq 0$, $\mathbf{W}_t^\mu$ has distribution $\mu$ and $(\mathbf{W}_t^\mu)_{t \geq 0}$ is a strong solution to (8). Therefore, $\mu$ is an invariant probability measure for (8) which concludes the proof. $\qquad\square$

## 11  Links with gradient flow approach

**Case $\beta \in [0, 1)$**  We now focus on the mean-field distribution $\boldsymbol{\lambda}^\star$. Note that the trajectories of $(\mathbf{W}_t^{k,\star})_{t \geq 0}$ for any $k \in \mathbb{N}^\star$ are deterministic conditionally to $\mathbf{W}_0^{k,\star}$. Using Itô's formula, we obtain that for any function $f \in \mathrm{C}^2(\mathbb{R}^p)$ with compact support and $t \geq 0$

$$\int_{\mathbb{R}^p} f(\tilde{w})\mathrm{d}\boldsymbol{\lambda}_t^\star(\tilde{w}) = \int_{\mathbb{R}^p} f(\tilde{w})\mathrm{d}\mu_0(\tilde{w}) + \int_0^t \int_{\mathbb{R}^p} (s+1)^{-\alpha}\langle h(\tilde{w}, \boldsymbol{\lambda}_s^\star), \nabla f(\tilde{w})\rangle \mathrm{d}\boldsymbol{\lambda}_s^\star(\tilde{w}) \ . \tag{46}$$

Therefore, if for any $t \geq 0$, $\boldsymbol{\lambda}_t^\star$ admits a density $\boldsymbol{\rho}_t^\star$ such that $(\boldsymbol{\rho}_t^\star)_{t \geq 0} \in \mathrm{C}^1(\mathbb{R}_+ \times \mathbb{R}^p, \mathbb{R})$ we obtain that $(\boldsymbol{\rho}_t)_{t \geq 0}$ satisfies the following evolution equation for any $t > 0$ and $w \in \mathbb{R}^p$

$$\partial_t \boldsymbol{\rho}_t^\star(w) = -(t+1)^{-\alpha}\mathrm{div}(\bar{h}(\cdot, \boldsymbol{\rho}_t^\star)\boldsymbol{\rho}_t^\star)(w) \ ,$$

with for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$ with density $\rho$, $h(w, \mu) = \bar{h}(w, \rho)$. In the case $\alpha = 0$, it is well-known, see [22, 23, 28], that $(\boldsymbol{\rho}_t^\star)_{t \geq 0}$ is a Wasserstein gradient flow for the functional $\mathscr{R}^\star : \mathscr{P}_2^c(\mathbb{R}^p) \to \mathbb{R}$ given for any $\rho \in \mathscr{P}_2^c(\mathbb{R}^p)$

$$\mathscr{R}^\star(\rho) = \int_{\mathsf{X} \times \mathsf{Y}} \ell\left(\int_{\mathbb{R}^p} F(\tilde{w}, x)\rho(\tilde{w})\mathrm{d}\tilde{w}, y\right)\mathrm{d}\pi(x, y) \ , \tag{47}$$

where $\mathscr{P}_2^c(\mathbb{R}^p)$ is the set of probability density satisfying $\int_{\mathbb{R}^p} \|\tilde{w}\|^2 \rho(\tilde{w})\mathrm{d}\tilde{w} < +\infty$.

**Case $\beta = 1$**  Focusing on $(\boldsymbol{\lambda}_t^\star)_{t \geq 0}$, we no longer obtain that $(\boldsymbol{\lambda}_t^\star)_{t \geq 0}$ is a gradient flow for (47). Indeed, using Itô's formula, we have the following evolution equation for any $f \in \mathrm{C}_c^2(\mathbb{R}^p)$ and $t \geq 0$

$$\int_{\mathbb{R}^p} f(\tilde{w})\mathrm{d}\boldsymbol{\lambda}_t^\star(\tilde{w}) = \int_{\mathbb{R}^p} f(\tilde{w})\mathrm{d}\mu_0(\tilde{w}) + \int_0^t \int_{\mathbb{R}^p} (s+1)^{-\alpha}\langle h(\tilde{w}, \boldsymbol{\lambda}_s^\star), \nabla f(\tilde{w})\rangle \mathrm{d}\boldsymbol{\lambda}_s^\star(\tilde{w})$$
$$+ \int_0^t \int_{\mathbb{R}^p} (s+1)^{-\alpha}\mathrm{Tr}(\Sigma(\tilde{w}, \boldsymbol{\lambda}_s^\star)\nabla^2 f(\tilde{w}))\mathrm{d}\tilde{w} \ . \tag{48}$$

We higlight that the additional term in (48) from (46) corresponds to some entropic regularization of the risk $\mathscr{R}^\star$. Indeed, if for any $w \in \mathbb{R}^p$ and $\mu \in \mathscr{P}(\mathbb{R}^p)$, $\Sigma = \beta \, \mathrm{Id}$ then, in the case $\alpha = 0$, we obtain that $(\boldsymbol{\rho}_t^\star)_{t \geq 0}$ is a gradient flow for $\rho \mapsto U^\star(\rho) + \beta\mathrm{Ent}(\rho)$, where $\mathrm{Ent} : \mathsf{K}_2 \to \mathbb{R}$ is given for any $\rho \in \mathsf{K}_2$ by

$$\mathrm{Ent}(\rho) = -\int_{\mathbb{R}^p} \rho(x)\log(\rho(x))\mathrm{d}x \ .$$

This second regime emphasizes that large stepsizes act as an implicit regularization procedure for SGD.

## 12  Additional Experiments

In this section we present additional experiments illustrating the convergence results of the empirical measures. Contrary to the main document we illustrate our results with histograms of the weights of the first and second layers of the network, with a large number of different values of the parameters $\alpha$, $\beta$ and $N$.

**Setting.** In order to perform the following experiments we implemented a two-layer fully connected neural network on PyTorch. The input layer has the size of the input data, *i.e.*, $N_{\text{input}} = 28 \times 28$ units in the case of the MNIST dataset [41] and $N_{\text{input}} = 32 \times 32 \times 3$ in the case of the CIFAR-10 dataset [42]. We use a varying number of $N$ units in the hidden layer and the output layer has 10 units corresponding to the 10 possible labels of the classification tasks. We use a ReLU activation function and the cross-entropy loss.

The linear layers' weights are initialized with PyTorch default initialization function which is a uniform initialization between $-1/N_{\text{input}}^{1/2}$ and $1/N_{\text{input}}^{1/2}$. In all our experiments, if not specified, we consider an initialization $\mathbf{W}_0^{1:N}$ with distribution $\mu_0^{\otimes N}$ where $\mu_0$ is the uniform distribution on $[-0.04, 0.04]$.

In order to train the network we use SGD as described in Section 2 with an initial learning rate of $\gamma N^\beta$. In the case where $\alpha > 0$ we decrease this stepsize at each iteration to have a learning rate of $\gamma N^\beta (n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha}$. All experiments on the MNIST dataset are run for a finite time horizon $T = 100$ and the ones on the CIFAR-10 dataset are run for $T = 10000$. The average runtime of the experiments for $N = 50000$ on the MNIST dataset is one day and the experiments on the CIFAR-10 dataset run during two days. The experiments were run on a cluster of 24 CPUs with 126Go of RAM.

All the histograms represented below correspond to the first coordinate of the weights' vector.

**Experiments.** Figure 4 shows that the empirical distributions of the weights converge as the number of hidden units $N$ goes to infinity. Those figures illustrate also the fact that we obtain two different limiting distributions one for $\beta < 1$ (represented on the 3 first figures) and one for $\beta = 1$ (on the last figure). The results presented on Figure 5 illustrate the same fact, one the second layer. This means that the results we stated in Section 3 are also true for the weights of the second layer, thanks to the procedure described for example in [28].

On Figure 6 and Figure 7 we show the results of the exact same experiments but this time using decreasing stepsizes and a parameter $\alpha = 0.25$. Once again our experiments illustrate the convergence of the empirical distributions to some limiting distribution, and we can also identify two regimes. Note that the limiting distribution satisfying (46) or (48) (depending on the value of $\beta$), it depends on the parameter $\alpha$. Therefore the limiting distribution obtained in the case where $\alpha = 0.25$ is different from the one obtained when $\alpha = 0$. This is particularly visible in the case where $\beta = 1$ (as shown in green on Figure 6 and Figure 7).

We now study the role of the batch size $M$ on the convergence toward the mean-field regime. Figure 8 illustrates the convergence of the empirical measures in the case where $\beta < 1$ (here $\beta = 0.75$) of the weights of the hidden layer of the neural network, for a fixed number of neurons $N = 10000$ for different batch sizes $M$. We indeed observe convergence with $M$.
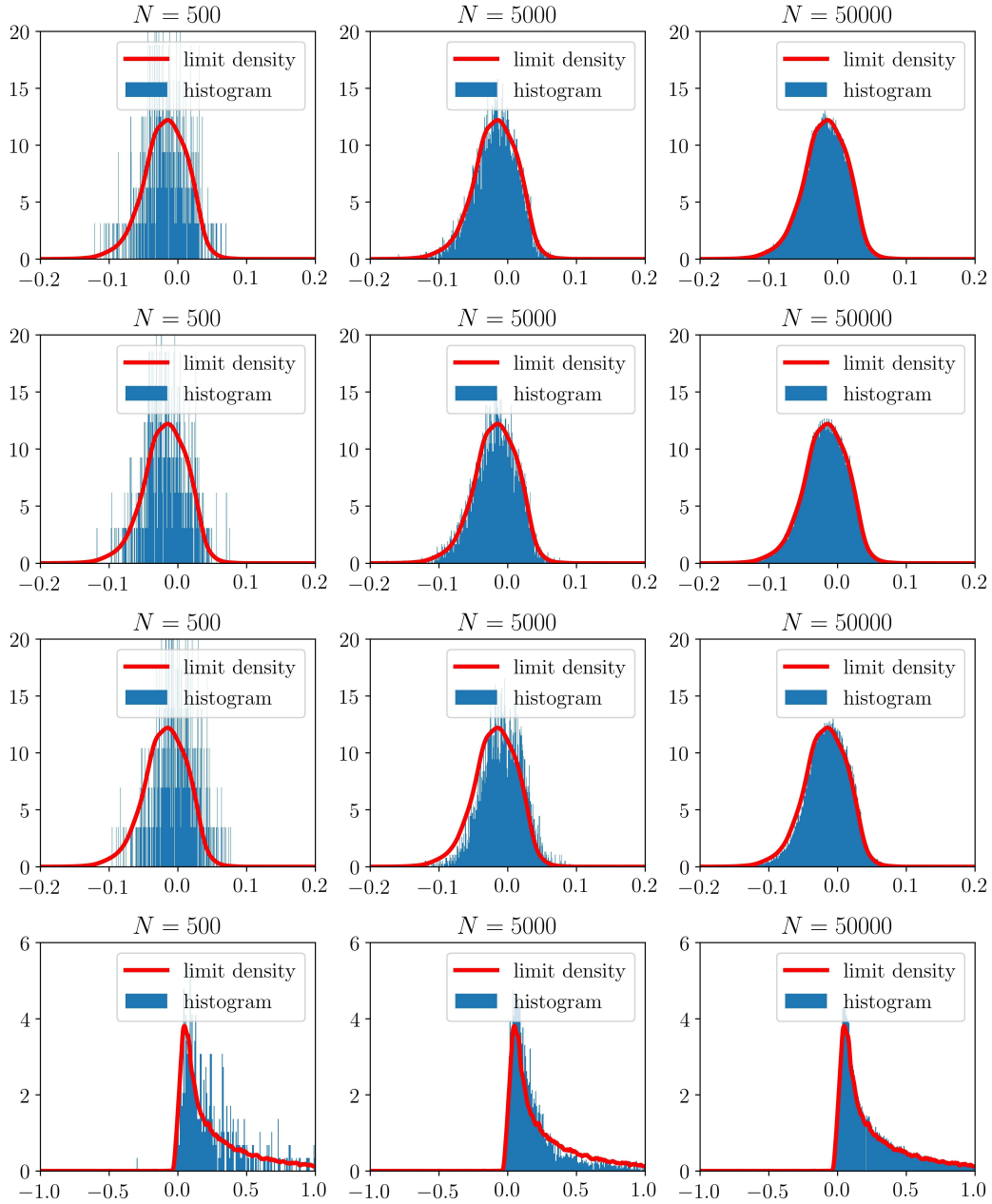
Figure 4: Convergence of the weights of the first layer as $N \to +\infty$ for $\alpha = 0$ and $M = 100$. The first line corresponds to $\beta = 0.25$, the second to $\beta = 0.5$, the third to $\beta = 0.75$ and the last line to $\beta = 1.0$.
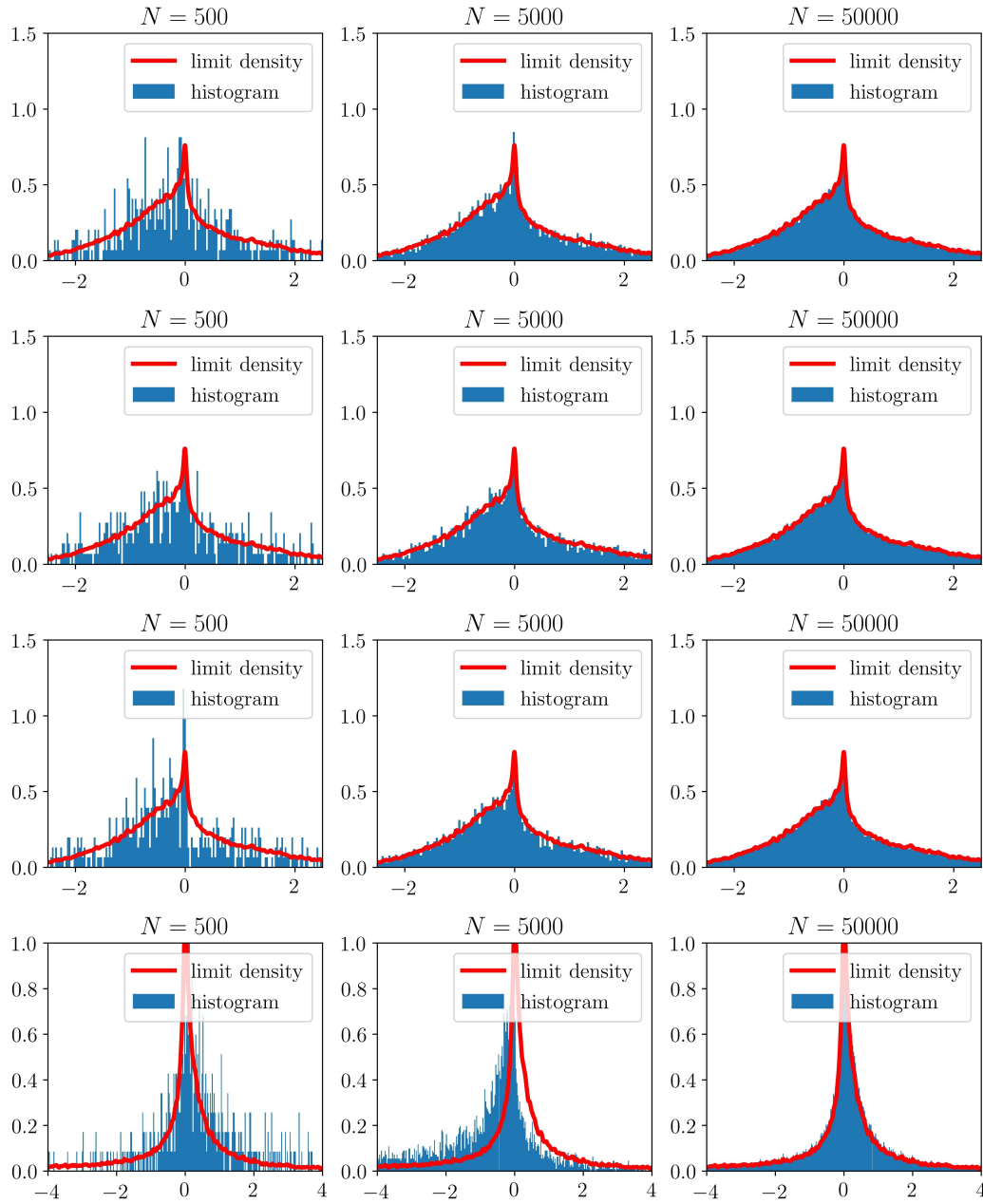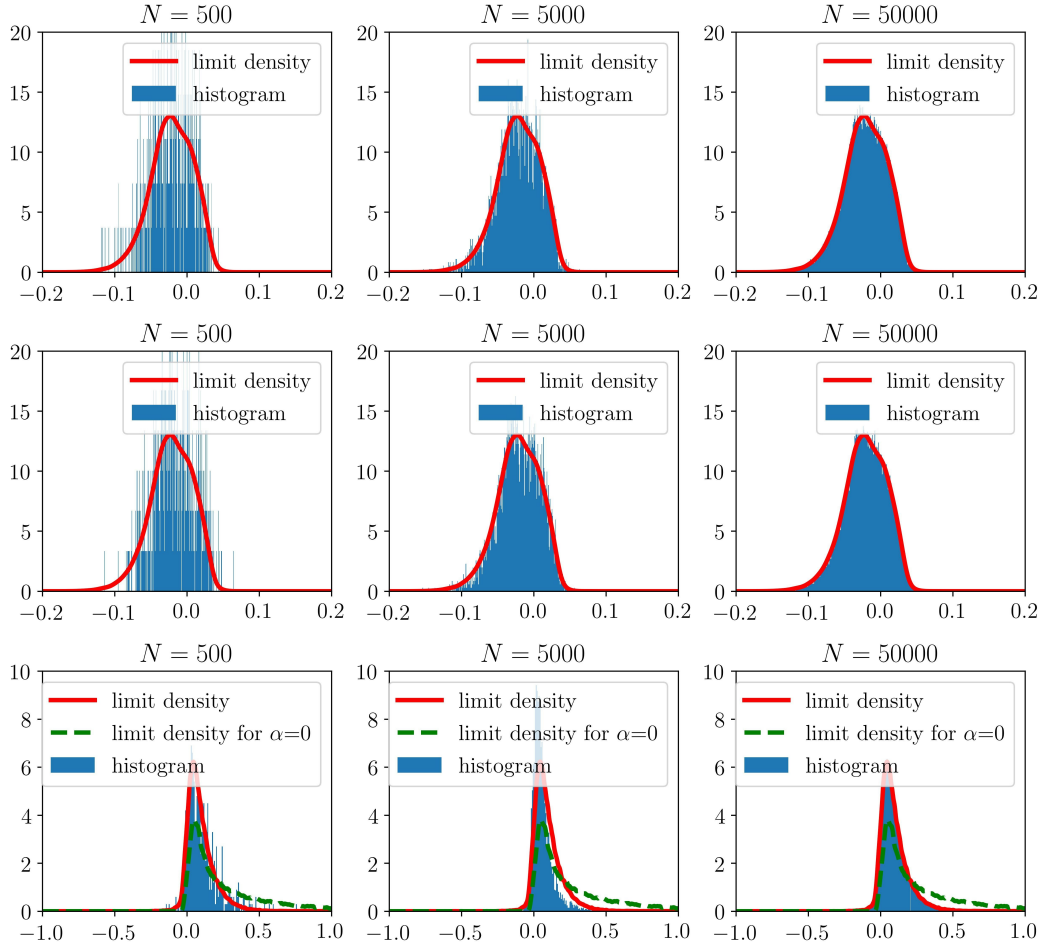
Figure 5: Convergence of the weights of the first layer as $N \to +\infty$ for $\alpha = 0$ and $M = 100$. The first line corresponds to $\beta = 0.25$, the second to $\beta = 0.5$, the third to $\beta = 0.75$ and the last line to $\beta = 1.0$.

Figure 6: Convergence of the weights of the first layer as $N \to +\infty$ for $\alpha = 0.25$ and $M = 100$. The first line corresponds to $\beta = 0.5$, the second to $\beta = 0.75$ and the last line to $\beta = 1.0$.
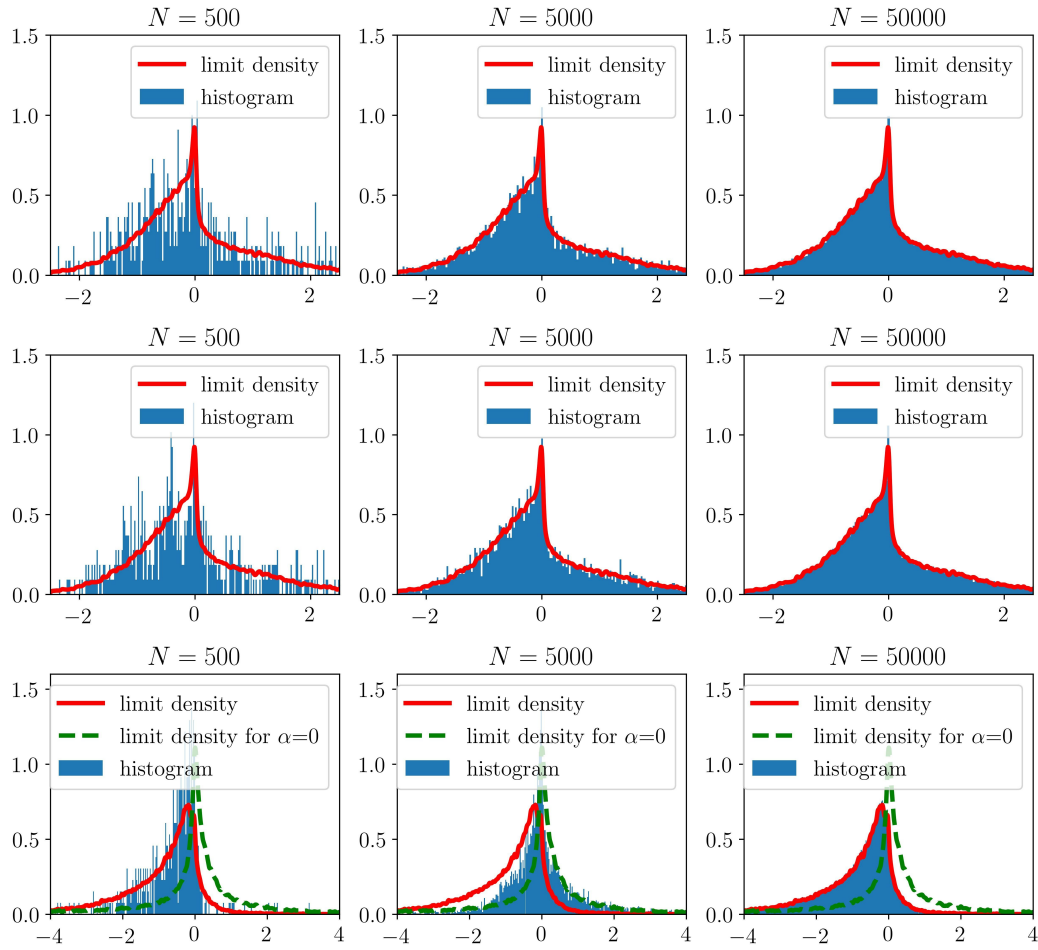
Figure 7: Convergence of the weights of the first layer as $N \to +\infty$ for $\alpha = 0.25$ and $M = 100$. The first line corresponds to $\beta = 0.5$, the second to $\beta = 0.75$ and the last line to $\beta = 1.0$.
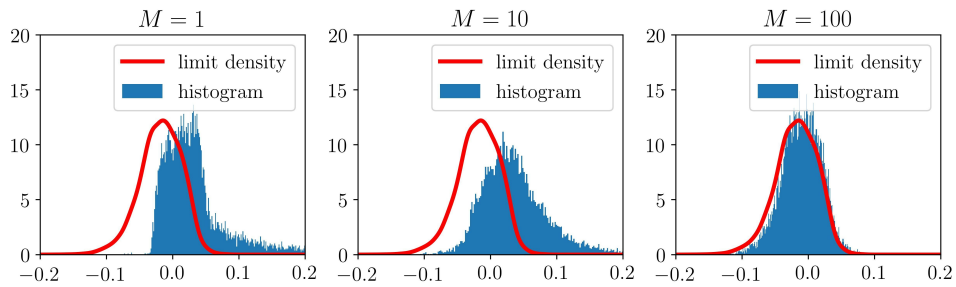


Figure 8: Convergence of the weights as $M \to \infty$