

---

# Quantitative Propagation of Chaos for SGD in Wide Neural Networks

## SUPPLEMENTARY DOCUMENT

---

**Valentin De Bortoli**  
University of Oxford  
debortoli@stats.ox.ac.uk

**Xavier Fontaine**  
Université Paris-Saclay  
fontaine@cmla.ens-cachan.fr

**Alain Durmus**  
Université Paris-Saclay  
alain.durmus@cmla.ens-cachan.fr

**Umut Şimşekli**  
Telecom Paris  
umut.simsekli@telecom-paris.fr

### Contents

<b>S1 Preliminaries</b>	<b>2</b>
S1.1 Notation . . . . .	2
S1.2 Wasserstein distances . . . . .	2
<b>S2 A mean-field modification of Stochastic Gradient Langevin Dynamics</b>	<b>3</b>
S2.1 Presentation of the modified SGLD and its continuous counterpart . . . . .	3
S2.2 Mean field approximation and propagation of chaos for mSGLD . . . . .	4
<b>S3 Technical results</b>	<b>4</b>
<b>S4 Quantitative propagation of chaos</b>	<b>8</b>
S4.1 Existence of strong solutions to the particle SDE . . . . .	8
S4.2 Existence of solutions to the mean-field SDE . . . . .	11
S4.3 Main result . . . . .	12
S4.4 Proofs of the main results . . . . .	14
<b>S5 Existence of invariant measure in the one-dimensional case</b>	<b>14</b>
<b>S6 Links with gradient flow approach</b>	<b>15</b>
<b>S7 Additional Experiments</b>	<b>16</b>

## S1 Preliminaries

### S1.1 Notation

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

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

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

Consider  $(F, d)$  a metric space. Let  $\mathcal{P}(F)$  be the space of probability measures over  $F$  equipped with its Borel  $\sigma$ -field  $\mathcal{B}(F)$ . For any  $\mu \in \mathcal{P}(F)$  and  $f : F \rightarrow \mathbb{R}$ , we say that  $f$  is  $\mu$ -integrable if  $\int_F |f(x)| d\mu(x) < +\infty$ . In this case, we set  $\mu[f] = \int_F f(x) d\mu(x)$ . Let  $\mu_0 \in \mathcal{P}(F)$ . For any  $r \geq 1$ , define  $\mathcal{P}_r(F) = \{\mu \in \mathcal{P}(F) : \int_{\mathbb{R}^p} d(\mu_0, \mu)^r d\mu(x) < +\infty\}$ . 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  $f : (E, \mathcal{E}) \rightarrow (G, \mathcal{G})$  be a measurable function. Then for any measure  $\mu$  on  $\mathcal{E}$  we define its pushforward measure by  $f$ ,  $f_\# \mu$ , for any  $A \in \mathcal{G}$  by  $f_\# \mu(A) = \mu(f^{-1}(A))$ .

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  $\mathbb{S}_p(\mathbb{R})$ .

### S1.2 Wasserstein distances

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

$$\mathcal{W}_r^r(\mu_1, \mu_2) = \inf_{\zeta \in \Lambda(\mu_1, \mu_2)} \left\{ \int_{F \times F} d(x, y)^r d\zeta(x, y) \right\}. \quad (\text{S1})$$

Note that  $\mathcal{W}_r$  is a distance on  $\mathcal{P}_r(F)$  by [1, Theorem 6.18]. In addition  $(\mathcal{P}_r(\mathbb{R}^p), \mathcal{W}_r)$  is a complete separable metric space. For any  $\mu_1, \mu_2 \in \mathcal{P}_p(F)$  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  $\mathcal{C}_{2,T}^p = C([0, T], \mathcal{P}_2(\mathbb{R}^p))$  is a complete separable metric space [2, 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 S1.** *Let  $r \geq 1$ ,  $N \in \mathbb{N}^*$ ,  $\{\alpha_k\}_{k=1}^N \in [0, 1]^N$  with  $\sum_{k=1}^N \alpha_k = 1$ ,  $\{\mu_{k,a}\}_{k=1}^N \in \mathcal{P}(F)^N$  and  $\{\mu_{k,b}\}_{k=1}^N \in \mathcal{P}(F)^N$ . Then, setting  $\nu_i = \sum_{k=1}^N \alpha_k \mu_{k,i}$  with  $i \in \{a, b\}$ , we have*

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

*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

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

□

As a special case of Proposition **S1**, we obtain that for any  $r \geq 1$ ,  $\{w_{k,a}\}_{k=1}^N \in \mathbb{F}^N$  and  $\{w_{k,b}\}_{k=1}^N \in \mathbb{F}^N$ ,

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

As another special case of Proposition **S1**, we obtain that for any  $\mu \in \mathcal{P}_r(\mathbb{F})$  and  $\{w_k\}_{k=1}^N \in \mathbb{F}^N$

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

## S2 A mean-field modification of Stochastic Gradient Langevin Dynamics

### S2.1 Presentation of the modified SGLD and its continuous counterpart

We start by introducing a modified Stochastic Gradient Langevin Dynamics (mSGLD) [3]. In the mean-field regime, this setting was studied in the case  $\beta = 0$  in [4]. We recall that the mean-field  $h : \mathbb{R}^p \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}^p$  and  $\xi : \mathbb{R}^p \times \mathcal{P}(\mathbb{R}^d) \times \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{R}^p$  are given for any  $\mu \in \mathcal{P}(\mathbb{R}^p)$ ,  $w \in \mathbb{R}^p$ ,  $(x, y) \in \mathbb{X} \times \mathbb{Y}$  by

$$\begin{aligned} h(w, \mu) &= - \int_{\mathbb{X} \times \mathbb{Y}} \partial_1 \ell(\mu[F(\cdot, x)], y) \nabla_w F(w, x) d\pi(x, y) - \nabla V(w), \\ \xi(w, \mu, x, y) &= -h(w, \mu) - \partial_1 \ell(\mu[F(\cdot, x)], y) \nabla_w F(w, x) - \nabla V(w). \end{aligned}$$

Let  $(W_0^k)_{k \in \mathbb{N}^*}$  be i.i.d.  $p$  dimensional random variables with distribution  $\mu_0$  and  $\{Z_k^n : k, n \in \mathbb{N}^*\}$  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  $n \in \mathbb{N}$ ,  $k \in \{1, \dots, N\}$ ,

$$\begin{aligned} W_{n+1}^{k,N} &= W_n^{k,N} + \gamma N^{\beta-1} (n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha} \{h(W_n^{k,N}, \nu_n^N) + \xi(W_n^{k,N}, \nu_n^N, X_n, Y_n)\} \\ &\quad + [2\eta\gamma N^{\beta-1} (n + \gamma_{\alpha,\beta}(N)^{-1})^{-\alpha}]^{1/2} Z_{k,n}, \quad (\text{S2}) \end{aligned}$$

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 case  $\eta = 0$ , we obtain (3). In addition, (S2) does not exactly correspond to the usual implementation of SGLD as introduced in [3]. 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 (S2). The scheme presented in (S2) 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 [4].

We now present the continuous model associated with this discrete process in the limit  $\gamma \rightarrow 0$  or  $N \rightarrow +\infty$ . For  $N \in \mathbb{N}^*$ , 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  $k \in \{1, \dots, N\}$  by

$$d\mathbf{W}_t^{k,N} = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^{k,N}, \nu_t^N) dt + \gamma_{\alpha,\beta}(N)^{1/2} \Sigma^{1/2}(\mathbf{W}_t^{k,N}, \nu_t^N) d\mathbf{B}_t^k + \sqrt{2\eta} d\tilde{\mathbf{B}}_t^k \right\}, \quad (\text{S3})$$

where  $\{\mathbf{B}_t^k\}_{t \geq 0} : k \in \mathbb{N}^*\}$  and  $\{\tilde{\mathbf{B}}_t^k\}_{t \geq 0} : k \in \mathbb{N}^*\}$  are two independent families of independent  $p$  dimensional Brownian motions and  $\nu_t^N$  is the empirical probability distribution of the particles defined for any  $t \geq 0$  by  $\nu_t^N = N^{-1} \sum_{k=1}^N \delta_{\mathbf{W}_t^{k,N}}$ . Similarly to Section 2, (S3) is the continuous

counterpart of (S2). Let  $M \in \mathbb{N}^*$ . 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  $k \in \{1, \dots, N\}$  by

$$d\mathbf{W}_t^{k,N} = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N) dt + (\gamma_{\alpha,\beta}(N)/M)^{1/2} \Sigma^{1/2}(\mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N) d\mathbf{B}_t^k + \sqrt{2\eta} d\tilde{\mathbf{B}}_t^k \right\}. \quad (\text{S4})$$

## S2.2 Mean field approximation and propagation of chaos for mSGLD

The following theorems are the extensions of Theorem 1 and Theorem 2 to (S3) for any  $\eta \geq 0$ . Note that in the case  $\eta = 0$ , Theorem S2 boils down to Theorem 1 and Theorem S3 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  $\mathbf{W}_0^*$  given by

$$d\mathbf{W}_t^* = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^*, \boldsymbol{\lambda}_t^*) dt + \sqrt{2\eta} \tilde{\mathbf{B}}_t \right\}, \quad \text{with } \boldsymbol{\lambda}_t^* \text{ the distribution of } \mathbf{W}_t^*. \quad (\text{S5})$$

**Theorem S2.** *Assume A1. Let  $(\mathbf{W}_0^k)_{k \in \mathbb{N}}$  be a sequence of i.i.d.  $\mathbb{R}^p$ -valued random variables with distribution  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^p)$  and set for any  $N \in \mathbb{N}^*$ ,  $\mathbf{W}_0^{1:N} = \{\mathbf{W}_0^k\}_{k=1}^N$ . Then, for any  $m \in \mathbb{N}^*$  and  $T \geq 0$ , there exists  $C_{m,T} \geq 0$  such that for any  $\alpha \in [0, 1)$ ,  $\beta \in [0, 1)$ ,  $M \in \mathbb{N}^*$  and  $N \in \mathbb{N}^*$*

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

with  $(\mathbf{W}_t^{1:m,N}, \mathbf{W}_t^{1:m,*}) = \{(\mathbf{W}_t^{k,N}, \mathbf{W}_t^{k,*})\}_{k=1}^m$ ,  $(\mathbf{W}_t^{1:N})$  is the solution of (S4) starting from  $\mathbf{W}_0^{1:N}$ , and for any  $k \in \{1, \dots, N\}$ ,  $\mathbf{W}_t^{k,*}$  is the solution of (S5) starting from  $\mathbf{W}_0^k$  and Brownian motion  $(\tilde{\mathbf{B}}_t^k)_{t \geq 0}$ .

*Proof.* The proof is postponed to Section S4.4 □

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

$$d\mathbf{W}_t^* = (t+1)^{-\alpha} \left\{ h(\mathbf{W}_t^*, \boldsymbol{\lambda}_t^*) dt + (\gamma^{1/(1-\alpha)} \Sigma(\mathbf{W}_t^*, \boldsymbol{\lambda}_t^*)/M)^{1/2} d\mathbf{B}_t + \sqrt{2\eta} d\tilde{\mathbf{B}}_t \right\}, \quad (\text{S6})$$

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

**Theorem S3.** *Let  $\beta = 1$ . Assume A1. Let  $(\mathbf{W}_0^k)_{k \in \mathbb{N}}$  be a sequence of  $\mathbb{R}^p$ -valued random variables with distribution  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^p)$  and assume that for any  $N \in \mathbb{N}^*$ ,  $\mathbf{W}_0^{1:N} = \{\mathbf{W}_0^k\}_{k=1}^N$ . Then, for any  $m \in \mathbb{N}^*$  and  $T \geq 0$ , there exists  $C_{m,T} \geq 0$  such that for any  $\alpha \in [0, 1)$ ,  $M \in \mathbb{N}^*$  and  $N \in \mathbb{N}^*$  we have*

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

with  $(\mathbf{W}_t^{1:m,N}, \mathbf{W}_t^{1:m,*}) = \{(\mathbf{W}_t^{k,N}, \mathbf{W}_t^{k,*})\}_{k=1}^m$ ,  $(\mathbf{W}_t^{1:N})$  is the solution of (S4) starting from  $\mathbf{W}_0^{1:N}$ , and for any  $k \in \{1, \dots, N\}$ ,  $\mathbf{W}_t^{k,*}$  is the solution of (S6) 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 S4.4 □

## S3 Technical results

In this section, we derive technical results needed to establish Theorem 1, Theorem 2, Theorem S2 and Theorem S3. In particular, we are interested in the regularity properties of the mean field  $h$  and the diffusion matrix  $\Sigma$  under A1. We recall that in this setting, for any  $w \in \mathbb{R}^p$ ,  $\mu \in \mathcal{P}(\mathbb{R}^p)$ ,  $(x, y) \in X \times Y$ , we have

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

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

$$\begin{aligned}\xi(w, \mu, x, y) &= -\tilde{h}(w, \mu) - \partial_1 \ell \left( \int_{\mathbb{R}^p} F(\zeta, x) d\mu(\zeta), y \right) \nabla_w F(w, x), \\ \Sigma(w, \mu) &= \int_{\mathbb{X} \times \mathbb{Y}} \{\xi \xi^\top\}(w, \mu, x, y) d\pi(x, y), \quad S(w, \mu) = \Sigma^{1/2}(w, \mu).\end{aligned}\quad (\text{S7})$$

Note that by **A1-(a)**, we obtain the following estimate used in the proof of the results of this Section: for any  $y, y \in \mathbb{R}$

$$|\partial_1 \ell(y, y)| \leq |\partial_1 \ell(0, y)| + \Psi(y) |y| \leq 2\Psi(y) \max(1, |y|). \quad (\text{S8})$$

In addition, note that under **A1-(c)**, there exists  $K \geq 0$  such that for any  $w \in \mathbb{R}^p$

$$\|\nabla^2 V(w)\| + \|\mathbb{D}^3 V(w)\| \leq K, \quad \|\nabla V(w)\| \leq K(1 + \|w\|). \quad (\text{S9})$$

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

$$G(w, x, y) = \{\Phi^4(x) + \Psi^2(y)\}F(w, x). \quad (\text{S10})$$

We now state our main regularity/boundedness proposition.

**Proposition S4.** *Assume **A1**. Then, there exists  $L \geq 0$  such that the following hold.*

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

$$\begin{aligned}\|h(w_1, \mu_1) - h(w_2, \mu_2)\| \\ \leq L \left\{ \|w_1 - w_2\| + \left( \int_{\mathbb{X} \times \mathbb{Y}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 d\pi(x, y) \right)^{1/2} \right\}.\end{aligned}\quad (\text{S11})$$

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

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

$$\begin{aligned}|S_{i,j}(w_1, \mu_1) - S_{i,j}(w_2, \mu_2)| \\ \leq L \left\{ \|w_1 - w_2\| + \left( \int_{\mathbb{X} \times \mathbb{Y}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 d\pi(x, y) \right)^{1/2} \right\}.\end{aligned}\quad (\text{S12})$$

*In addition, we have for any  $\mu \in \mathcal{P}(\mathbb{R}^p)$ ,  $w \in \mathbb{R}^p$  and  $i, j \in \{1, \dots, p\}$ ,  $|S_{i,j}(w, \mu)| \leq L$ .*

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

*Proof.* (a) First, we show that (S11) holds. Note that by the triangle inequality and (S7), we only need to consider  $h \leftarrow \tilde{h}$  and  $h \leftarrow V$ . The case  $h \leftarrow V$  is straightforward using (S9). We now deal with the first case. For any  $w_1, w_2 \in \mathbb{R}^p$  and  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$ , 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 **A1-(a)**, **A1-(b)**, (S7) and (S8) we have for any  $w_1, w_2 \in \mathbb{R}^p$  and  $\mu_1 \in \mathcal{P}(\mathbb{R}^p)$

$$\begin{aligned}\|\tilde{h}(w_1, \mu_1) - \tilde{h}(w_2, \mu_1)\| &\leq \left\| \int_{\mathbb{X} \times \mathbb{Y}} \partial_1 \ell(\mu_1[F(\cdot, x)], y) \nabla_w F(w_1, x) d\pi(x, y) \right. \\ &\quad \left. - \int_{\mathbb{X} \times \mathbb{Y}} \partial_1 \ell(\mu_1[F(\cdot, x)], y) \nabla_w F(w_2, x) d\pi(x, y) \right\| \\ &\leq \int_{\mathbb{X} \times \mathbb{Y}} |\partial_1 \ell(\mu_1[F(\cdot, x)], y)| \Phi(x) d\pi(x, y) \|w_1 - w_2\| \\ &\leq \int_{\mathbb{X} \times \mathbb{Y}} \Psi(y) \Phi(x) (1 + |\mu_1[F(\cdot, x)]|) d\pi(x, y) \|w_1 - w_2\| \\ &\leq 2 \int_{\mathbb{X} \times \mathbb{Y}} \Psi(y) \Phi^2(x) d\pi(x, y) \|w_1 - w_2\|.\end{aligned}\quad (\text{S13})$$

Using **A1-(a)**, **A1-(b)**, **(S7)** and the Cauchy-Schwarz inequality, we also have for any  $w_1 \in \mathbb{R}^p$  and  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$

$$\begin{aligned}
& \|\tilde{h}(\mu_1, w_1) - \tilde{h}(\mu_2, w_1)\| \\
& \leq \left\| \int_{\mathbb{X} \times \mathbb{Y}} \{\partial_1 \ell(\mu_1[F(\cdot, x)], y) \nabla_w F(w_1, x) - \partial_1 \ell(\mu_2[F(\cdot, x)], y) \nabla_w F(w_1, x)\} d\pi(x, y) \right\| \\
& \leq \int_{\mathbb{X} \times \mathbb{Y}} |\partial_1 \ell(\mu_1[F(\cdot, x)], y) - \partial_1 \ell(\mu_2[F(\cdot, x)], y)| \|\nabla_w F(w_1, x)\| d\pi(x, y) \\
& \leq \int_{\mathbb{X} \times \mathbb{Y}} \Psi(y) \|\mu_1[F(\cdot, x)] - \mu_2[F(\cdot, x)]\| \Phi(x) d\pi(x, y) \\
& \leq \left( \int_{\mathbb{X} \times \mathbb{Y}} \Psi^2(y) \Phi^2(x) d\pi(x, y) \right)^{1/2} \left( \int_{\mathbb{X}} \|\mu_1[F(\cdot, x)] - \mu_2[F(\cdot, x)]\|^2 d\pi(x) \right)^{1/2}. \quad (\text{S14})
\end{aligned}$$

Combining **(S10)**, **(S13)**, **(S14)**, the fact that for any  $a, b \geq 0$ ,  $2ab \leq a^2 + b^2$  and **A1-(d)**, we obtain that there exists  $L_1 \geq 0$  such that for any  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$  and  $w_1, w_2 \in \mathbb{R}^p$  we have

$$\begin{aligned}
& \|\tilde{h}(w_1, \mu_1) - \tilde{h}(w_2, \mu_2)\| \\
& \leq L_1 \left\{ \|w_1 - w_2\| + \left( \int_{\mathbb{X} \times \mathbb{Y}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 d\pi(x, y) \right)^{1/2} \right\}.
\end{aligned}$$

In addition, using **A1-(b)** and **(S8)**, we have for any  $w \in \mathbb{R}^p$ ,  $\mu \in \mathcal{P}(\mathbb{R}^p)$ ,  $x \in \mathbb{X}$  and  $y \in \mathbb{Y}$

$$|\partial_1 \ell(\mu[F(\cdot, x)], y) \|\nabla_w F(w, x)\| \leq \Psi(y) \Phi(x) (1 + \Phi(x)) \leq 2\Psi(y) \Phi^2(x). \quad (\text{S15})$$

Therefore, combining this result and **(S7)**, we get that for any  $w \in \mathbb{R}^p$  and  $\mu \in \mathcal{P}(\mathbb{R}^p)$

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

Using the fact that for any  $a, b \geq 0$ ,  $2ab \leq a^2 + b^2$  and **A1-(d)**, there exists  $L_2 \geq 0$  such that for any  $w \in \mathbb{R}^p$  and  $\mu \in \mathcal{P}(\mathbb{R}^p)$ ,

$$\|\tilde{h}(w, \mu)\| \leq L_2 \quad (\text{S16})$$

(b) Second, we first show that there exists  $L_3 \geq 0$  such that for any  $\mu \in \mathcal{P}(\mathbb{R}^p)$ ,  $w \in \mathbb{R}^p$  and  $i, j \in \{1, \dots, p\}$ ,  $|S_{i,j}(w, \mu)| \leq L$ . Let  $i, j \in \{1, \dots, p\}$ . We have for any  $w \in \mathbb{R}^p$  and  $\mu \in \mathcal{P}(\mathbb{R}^p)$

$$|S_{i,j}(w, \mu)| \leq \|S(w, \mu)\| \leq \text{Tr}^{1/2}(\Sigma(w, \mu)). \quad (\text{S17})$$

Similarly to **(S15)**, using **(S7)**, **(S16)**, 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  $\mu \in \mathcal{P}(\mathbb{R}^p)$

$$\text{Tr}(\Sigma(w, \mu)) \leq \int_{\mathbb{X} \times \mathbb{Y}} \|\xi(w, \mu, x, y)\|^2 d\pi(x, y) \leq 2 \int_{\mathbb{X} \times \mathbb{Y}} \{L_2^2 + 2\Psi^2(y) \Phi^4(x)\} d\pi(x, y). \quad (\text{S18})$$

Combining **(S17)**, **(S18)** and **A1-(d)**, there exists  $L_3 \geq 0$  such that for any  $w \in \mathbb{R}^p$  and  $\mu \in \mathcal{P}(\mathbb{R}^p)$ ,  $\max_{1 \leq i, j \leq p} |S_{i,j}(w, \mu)| \leq L_3$ .

We now show that **(S12)** holds. For any  $w_1, w_2 \in \mathbb{R}^p$ ,  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$  define  $\varphi_\Sigma : [0, 1] \rightarrow \mathbb{S}_p(\mathbb{R})$  for any  $t \in [0, 1]$  by

$$\varphi_\Sigma(t) = \Sigma(tw_1 + (1-t)w_2, t\mu_1 + (1-t)\mu_2). \quad (\text{S19})$$

For ease of notation, the dependency of  $\varphi_\Sigma$  with respect to  $w_1, w_2 \in \mathbb{R}^p$  and  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$  is omitted. In what follows, we show that for any  $w_1, w_2 \in \mathbb{R}^p$ ,  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$ ,  $\varphi_\Sigma \in C^2([0, 1], \mathbb{S}_p(\mathbb{R}))$  and that there exists  $L_4 \geq 0$  such that for any  $t \in [0, 1]$

$$\|\varphi_\Sigma''(t)\| \leq L_4 \left\{ \|w_1 - w_2\| + \left( \int_{\mathbb{X} \times \mathbb{Y}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 d\pi(x, y) \right)^{1/2} \right\}^2,$$

which will conclude the proof of **(S12)** upon using a straightforward adaptation of **[5, Lemma 3.2.3, Theorem 5.2.3]**. We conclude the proof of Proposition **S4** upon letting  $L = \max(L_1, L_2, L_3, L_4)$ .

For any  $t \in [0, 1]$ , let  $\mu_t = \mu_1 + t(\mu_2 - \mu_1) \in \mathcal{P}(\mathbb{R}^p)$  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

$$\begin{aligned} f(t, x, y) &= \partial_1 \ell(\mu_t[F(\cdot, x)], y) \nabla_w F(w_t, x), \\ \tilde{f}(t, x, y) &= \xi(w_t, \mu_t, x, y) = \int_{\mathsf{X} \times \mathsf{Y}} f(t, x, y) d\pi(x, y) - f(t, x, y). \end{aligned} \quad (\text{S20})$$

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

By (S7), (S15) and (S16), we get that for any  $w_1, w_2 \in \mathbb{R}^p$ ,  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$ ,  $(x, y) \in \mathsf{X} \times \mathsf{Y}$  and  $t \in [0, 1]$

$$\|f(t, x, y)\| \leq 2\Psi(y)\Phi^2(x), \quad \|\tilde{f}(t, x, y)\| \leq L_2 + 2\Psi(y)\Phi^2(x). \quad (\text{S21})$$

Using (S20), A1-(a) and A1-(b), we have that for any  $(x, y) \in \mathsf{X} \times \mathsf{Y}$ ,  $f(\cdot, x, y) \in C^1([0, 1], \mathbb{R}^p)$  and for any  $w_1, w_2 \in \mathbb{R}^p$ ,  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$ ,  $(x, y) \in \mathsf{X} \times \mathsf{Y}$  and  $t \in [0, 1]$

$$\begin{aligned} \partial_1 f(t, x, y) &= \partial_1^2 \ell(\mu_t[F(\cdot, x)], y) \nabla_w F(w_t, x) (\mu_2[F(\cdot, x)] - \mu_1[F(\cdot, x)]) \\ &\quad + \partial_1 \ell(\mu_t[F(\cdot, x)], y) \nabla_w^2 F(w_t, x) (w_2 - w_1). \end{aligned} \quad (\text{S22})$$

Using A1-(a), A1-(b), (S10) and (S8), we get that for any  $(x, y) \in \mathsf{X} \times \mathsf{Y}$  and  $t \in [0, 1]$

$$\|\partial_1 f(t, x, y)\| \leq 3\Psi(y)\Phi^2(x) (\|w_2 - w_1\| + \|\mu_1[F(\cdot, x)] - \mu_2[F(\cdot, x)]\|), \quad (\text{S23})$$

Similarly, using (S22), A1-(a) and A1-(b), we have that for any  $(x, y) \in \mathsf{X} \times \mathsf{Y}$ ,  $f(\cdot, x, y) \in C^2([0, 1], \mathbb{R}^p)$  and for any  $w_1, w_2 \in \mathbb{R}^p$ ,  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$ ,  $(x, y) \in \mathsf{X} \times \mathsf{Y}$  and  $t \in [0, 1]$

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

Using A1-(a), A1-(b) and (S8) 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  $t \in [0, 1]$

$$\|\partial_1^2 f(t, x, y)\| \leq 5\Psi(y)\Phi^2(x) \left( \|w_2 - w_1\|^2 + \|\mu_1[F(\cdot, x)] - \mu_2[F(\cdot, x)]\|^2 \right). \quad (\text{S24})$$

Combining (S20), (S23), (S24), A1-(d) and the dominated convergence theorem, we get that for any  $(x, y) \in \mathsf{X} \times \mathsf{Y}$ ,  $\tilde{f}(\cdot, x, y) \in C^2([0, 1], \mathbb{R}^p)$ . In addition, using (S20), (S21), (S23), (S24), 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  $w_1, w_2 \in \mathbb{R}^p$ ,  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$ ,  $(x, y) \in \mathsf{X} \times \mathsf{Y}$  and  $t \in [0, 1]$

$$\begin{aligned} \|\tilde{f}(t, x, y)\| &\leq C (\Phi^4(x) + \Psi^2(y)), \\ \|\partial_1 \tilde{f}(t, x, y)\| &\leq C (\Phi^4(x) + \Psi^2(y)) \chi(w_1, w_2, \mu_1, \mu_2, x), \\ \|\partial_1^2 \tilde{f}(t, x, y)\| &\leq C (\Phi^4(x) + \Psi^2(y)) \chi^2(w_1, w_2, \mu_1, \mu_2, x), \end{aligned} \quad (\text{S25})$$

where

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

Using (S19) and (S7), we have that for any  $w_1, w_2 \in \mathbb{R}^p$ ,  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$ ,  $t \in [0, 1]$

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

Combining this result, (S25) and A1-(d) we get that for any  $w_1, w_2 \in \mathbb{R}^p$  and  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$ ,  $\varphi_\Sigma \in C^2([0, 1], \mathbb{S}_p(\mathbb{R}))$  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  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$ ,  $t \in [0, 1]$  and  $u \in \mathbb{R}^p$  with  $\|u\| = 1$ , we have

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

$$\begin{aligned}
&\leq 2 \int_{\mathbb{X} \times \mathbb{Y}} \|\partial_1 \tilde{f}(t, x, y)\|^2 d\pi(x, y) + 2 \int_{\mathbb{X} \times \mathbb{Y}} \|\partial_1^2 \tilde{f}(t, x, y)\| \|\tilde{f}(t, x, y)\| d\pi(x, y) \\
&\leq C_1 \int_{\mathbb{X} \times \mathbb{Y}} (\Phi^8(x) + \Psi^4(y)) \chi^2(w_1, w_2, x, y) d\pi(x, y) \\
&\leq C_2 \left\{ \|w_1 - w_2\| + \left( \int_{\mathbb{X}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 d\pi(x, y) \right)^{1/2} \right\}^2,
\end{aligned}$$

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

$$\begin{aligned}
\|\varphi''_{\Sigma}(t)\| &= \sup_{u \in \mathbb{R}^p, \|u\|=1} \langle u, \varphi''_{\Sigma}(t)u \rangle \\
&\leq C \left\{ \|w_1 - w_2\| + \left( \int_{\mathbb{X}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 d\pi(x, y) \right)^{1/2} \right\}^2.
\end{aligned}$$

Combining this result and a straightforward adaptation of [5, Lemma 3.2.3, Theorem 5.2.3] we obtain that for any  $w_1, w_2 \in \mathbb{R}^p$ ,  $\mu_1, \mu_2 \in \mathcal{P}(\mathbb{R}^p)$

$$|S_{i,j}(w_1, \mu_1) - S_{i,j}(w_2, \mu_2)| \leq L_4 \left\{ \|w_1 - w_2\| + \left( \int_{\mathbb{X}} \|\mu_1[G(\cdot, x, y)] - \mu_2[G(\cdot, x, y)]\|^2 d\pi(x, y) \right)^{1/2} \right\},$$

with  $L_4 = \sqrt{2C}p$ .

(c) Using (S7), we have for any  $w \in \mathbb{R}^p$  and  $\mu \in \mathcal{P}(\mathbb{R}^p)$

$$\int_{\mathbb{X} \times \mathbb{Y}} \|\xi(w, \mu, x, y)\|^2 d\pi(x, y) = \int_{\mathbb{X} \times \mathbb{Y}} \text{Tr}(\xi \xi^\top(w, \mu, x, y)) d\pi(x, y) = \sum_{i,j=1}^p |S_{i,j}(w, \mu)|^2 \leq p^2 L^2.$$

□

## S4 Quantitative propagation of chaos

### S4.1 Existence of strong solutions to the particle SDE

In this section, for two functions  $A, B : \bigcup_{N \in \mathbb{N}^*} \{1, \dots, N\} \times \mathbb{R}_+ \times (\mathbb{R}^p)^2 \times (\mathcal{P}_2(\mathbb{R}^p))^2 \rightarrow \mathbb{R}$ , 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  $N \in \mathbb{N}^*$ ,  $k \in \{1, \dots, N\}$ ,  $t \in \mathbb{R}_+$ ,  $w_1, w_2 \in \mathbb{R}^p$ ,  $\mu_1, \mu_2 \in \mathcal{P}_2(\mathbb{R}^p)$ ,  $A_N(k, t, w_1, w_2, \mu_1, \mu_2) \leq C B_N(k, t, w_1, w_2, \mu_1, \mu_2)$ , where  $A_N$  and  $B_N$  are the restrictions of  $A$  and  $B$  to  $\{1, \dots, N\} \times \mathbb{R}_+ \times (\mathbb{R}^p)^2 \times (\mathcal{P}_2(\mathbb{R}^p))^2$ .

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

$$d\mathbf{W}_t^{k,N} = b_N(t, \mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N) dt + \sigma_N(t, \mathbf{W}_t^{k,N}, \boldsymbol{\nu}_t^N) d\mathbf{B}_t^k, \quad \boldsymbol{\nu}_t^N = (1/N) \sum_{k=1}^N \delta_{\mathbf{W}_t^{k,N}}, \quad (\text{S26})$$

where  $(\mathbf{B}_t^k)_{k \in \mathbb{N}^*}$  are independent  $r$ -dimensional Brownian motions and where  $(b_N)_{N \in \mathbb{N}^*}$  and  $(\sigma_N)_{N \in \mathbb{N}^*}$  are family of measurable functions such that for any  $N \in \mathbb{N}^*$ ,  $b_N : \mathbb{R}_+ \times \mathbb{R}^p \times \mathcal{P}_2(\mathbb{R}^p) \rightarrow \mathbb{R}^p$  and  $\sigma_N : \mathbb{R}_+ \times \mathbb{R}^p \times \mathcal{P}_2(\mathbb{R}^p) \rightarrow \mathbb{R}^{p \times r}$ . We make the following assumption ensuring the existence and uniqueness of solutions of (S26) for any  $N \in \mathbb{N}^*$ . Consider in the sequel a measurable space  $(Z, \mathcal{Z})$  and a probability measure  $\pi_Z$  on this space.

**B1.** *There exist a measurable function  $g : \mathbb{R}^p \times Z \rightarrow \mathbb{R}$ ,  $M_1 \geq 0$  and  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^p)$  such that for any  $N \in \mathbb{N}^*$ , the following hold.*

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

$$\|g(w_1, z) - g(w_2, z)\| \leq \zeta(z) \|w_1 - w_2\|, \quad \|g(w_1, z)\| \leq \zeta(z), \quad \text{with } \int_Z \zeta^2(z) d\pi_Z(z) < +\infty.$$



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

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

$$\begin{aligned} & \sup_{t \geq 0} \{ \|b_N(t, w_1, \mu_1) - b_N(t, w_2, \mu_2)\| + \|\sigma_N(t, w_1, \mu_1) - \sigma_N(t, w_2, \mu_2)\| \} \\ & \leq M_1 \left\{ \|w_1 - w_2\| + \left( \int_{\mathbb{Z}} |\mu_1[g(\cdot, z)] - \mu_2[g(\cdot, z)]|^2 d\pi_{\mathbb{Z}}(z) \right)^{1/2} \right\}, \\ & \sup_{t \geq 0} \{ \|b_N(t, 0, \mu_0)\| + \|\sigma_N(t, 0, \mu_0)\| \} \leq M_1. \end{aligned}$$

**B2.** There exist  $M_2 \geq 0$ ,  $\kappa > 0$ ,  $b \in C(\mathbb{R}_+ \times \mathbb{R}^p \times \mathcal{P}_2(\mathbb{R}^p), \mathbb{R}^p)$  and  $\sigma \in C(\mathbb{R}_+ \times \mathbb{R}^p \times \mathcal{P}_2(\mathbb{R}^p), \mathbb{R}^{p \times r})$  such that

$$\sup_{t \geq 0, w \in \mathbb{R}^p, \mu \in \mathcal{P}_2(\mathbb{R}^p)} \{ \|b_N(t, w, \mu) - b(t, w, \mu)\| + \|\sigma_N(t, w, \mu) - \sigma(t, w, \mu)\| \} \leq M_2 N^{-\kappa}.$$

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

$$\|b_N(t, w, \mu)\| + \|\sigma_N(t, w, \mu)\| \lesssim \left[ 1 + \|w\| + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) d\mu(\tilde{w}) \right)^{1/2} \right], \quad (\text{S27})$$

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

**Theorem S5.** Assume **B1**. Then for any  $N \in \mathbb{N}^*$ , (S26) admits a unique strong solution. If in addition, there exists  $m \geq 1$  such that  $\sup_{N \in \mathbb{N}^*} \sup_{k \in \{1, \dots, N\}} \mathbb{E}[\|\mathbf{W}_0^{k, N}\|^{2m}] < +\infty$ , then for any  $T \geq 0$ , there exists  $C \geq 0$  such that

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

*Proof.* First, we show that for any  $N \in \mathbb{N}^*$ , (S26) admits a unique strong solution. Let  $\tilde{b}_N : \mathbb{R}_+ \times (\mathbb{R}^p)^N \rightarrow (\mathbb{R}^p)^N$  and  $\tilde{\sigma}_N : \mathbb{R}_+ \times (\mathbb{R}^p)^N \rightarrow (\mathbb{R}^{p \times r})^N$  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

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

Let  $w_1^{1:N}, w_2^{1:N} \in (\mathbb{R}^p)^N$ . Using **B1**, Proposition **S1** and that for any  $a, b \geq 0$ ,  $(a + b)^{1/2} \leq a^{1/2} + b^{1/2}$ , we have

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

Similarly, we have  $\|\sigma_N(t, w_1^{k, N}, \nu^{N, w_1}) - \sigma_N(t, w_2^{k, N}, \nu^{N, w_2})\| \lesssim \|w_1^{1:N} - w_2^{1:N}\|$ . Therefore, we obtain that for any  $N \in \mathbb{N}^*$ ,  $\tilde{b}_N$  and  $\tilde{\sigma}_N$  are Lipschitz-continuous and using [6, Theorem 2.9], we get that there exists a unique strong solution to (S26). Let  $m \geq 1$  and assume that  $\sup_{N \in \mathbb{N}^*} \sup_{k \in \{1, \dots, N\}} \mathbb{E}[\|\mathbf{W}_0^{k, N}\|^{2m}] < +\infty$ , we now show that for any  $T \geq 0$ , there exists  $C \geq 0$  such that

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

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

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

Combining this result with (S27), 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

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

$$\begin{aligned}
&\lesssim \left[ 1 + \|w\| + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) d\mu(\tilde{w}) \right)^{1/2} \right] \|\nabla V_m(w)\| \\
&\quad + \left[ 1 + \|w\| + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) d\mu(\tilde{w}) \right)^{1/2} \right]^2 \|\nabla^2 V_m(w)\| \\
&\lesssim \left[ 1 + \|w\| + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) d\mu(\tilde{w}) \right)^{1/2} \right] \|w\|^{2m-1} \\
&\quad + \left[ 1 + \|w\|^2 + \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) d\mu(\tilde{w}) \right] \|w\|^{2m-2} \\
&\lesssim 1 + \|w\|^{2m} + \left( \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^2) d\mu(\tilde{w}) \right)^m \lesssim 1 + \|w\|^{2m} + \int_{\mathbb{R}^p} (1 + \|\tilde{w}\|^{2m}) d\mu(\tilde{w}) .
\end{aligned} \tag{S28}$$

Now let  $\tau_n^N = \inf\{t \geq 0 : \|\mathbf{W}_t^{k,N}\| \geq n \text{ for some } k \in \{1, \dots, N\}\}$ . Using Itô's lemma, (S28) and (S26), we have

$$\begin{aligned}
\mathbb{E} \left[ V_m(\mathbf{W}_{t \wedge \tau_n^N}^{k,N}) \right] &= \mathbb{E} \left[ V_m(\mathbf{W}_{0 \wedge \tau_n^N}^{k,N}) \right] + \mathbb{E} \left[ \int_0^{t \wedge \tau_n^N} \langle \nabla V_m(\mathbf{W}_s^{k,N}), b_N(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N) \rangle ds \right] \\
&\quad + (1/2) \mathbb{E} \left[ \int_0^{t \wedge \tau_n^N} \langle \nabla^2 V_m(\mathbf{W}_s^{k,N}), \sigma_N \sigma_N^\top(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N) \rangle ds \right] \\
&\lesssim \mathbb{E} \left[ V_m(\mathbf{W}_{0 \wedge \tau_n^N}^{k,N}) \right] + \mathbb{E} \left[ \int_0^{t \wedge \tau_n^N} \left\{ V_m(\mathbf{W}_s^{k,N}) + (1/N) \sum_{j=1}^N V_m(\mathbf{W}_s^{j,N}) \right\} ds \right]
\end{aligned}$$

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

$$\begin{aligned}
&\mathbb{E} \left[ V_m(\mathbf{W}_t^{k,N}) + (1/N) \sum_{j=1}^N V_m(\mathbf{W}_t^{j,N}) \right] \\
&\lesssim \mathbb{E} \left[ V_m(\mathbf{W}_0^{k,N}) + (1/N) \sum_{j=1}^N V_m(\mathbf{W}_0^{j,N}) \right] + \int_0^t \mathbb{E} \left[ V_m(\mathbf{W}_s^{k,N}) + (1/N) \sum_{j=1}^N V_m(\mathbf{W}_s^{j,N}) \right] ds .
\end{aligned}$$

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

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

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

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

Using Jensen's inequality, Burkholder-Davis-Gundy's inequality [7, IV.42], (S27) and the fact that for any  $(a_j)_{j \in \{1, \dots, M\}}$  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  $m \in \mathbb{N}^*$

$$\begin{aligned}
&\mathbb{E} \left[ \sup_{t \in [0, T]} \|\mathbf{W}_t^{k,N}\|^{2m} \right] \\
&\lesssim \mathbb{E} \left[ \sup_{t \in [0, T]} \left\| \int_0^t b_N(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N) ds \right\|^{2m} \right] + \mathbb{E} \left[ \sup_{t \in [0, T]} \left\| \int_0^t \sigma_N^{1/2}(s, \mathbf{W}_s^{k,N}, \boldsymbol{\nu}_s^N) d\mathbf{B}_s \right\|^{2m} \right]
\end{aligned}$$

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

which concludes the proof.  $\square$

## S4.2 Existence of solutions to the mean-field SDE

The following result is based on [8, Theorem 1.1] showing, under **B1** and **B2**, the existence of strong solutions and pathwise uniqueness for non-homogeneous McKean-Vlasov SDE with non-constant covariance matrix:

$$d\mathbf{W}_t^* = b(t, \mathbf{W}_t^*, \boldsymbol{\lambda}_t^*) dt + \sigma(t, \mathbf{W}_t^*, \boldsymbol{\lambda}_t^*) d\mathbf{B}_t, \quad (\text{S29})$$

where  $b$  and  $\sigma$  are given in **B2** and where for any  $t \geq 0$ ,  $\mathbf{W}_t^*$  has distribution  $\boldsymbol{\lambda}_t^* \in \mathcal{P}_2(\mathbb{R}^p)$ ,  $(\mathbf{B}_t)_{t \geq 0}$  is a  $r$  dimensional Brownian motion and  $\mathbf{W}_0^*$  has distribution  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^p)$ .

**Proposition S6.** *Assume **B1** and **B2**. Let  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^p)$ . Then, there exists a  $(\mathcal{F}_t)_{t \geq 0}$ -adapted continuous process  $(\mathbf{W}_t^*)_{t \geq 0}$  which is the unique strong solution of (S29) satisfying for any  $T \geq 0$ ,  $\sup_{t \in [0, T]} \mathbb{E}[\|\mathbf{W}_t^*\|^2] < +\infty$ .*

*Proof.* Let  $\delta \geq 0$  and  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^p)$ . Note that we only need to show that (S29) admits a strong solution up to  $\delta > 0$ . First, using [6, Theorem 2.9], note that for any  $(\boldsymbol{\mu}_t)_{t \in [0, \delta]} \in \mathcal{C}_{2, \delta}^p$  the SDE,

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

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

$$\|b(t, w_1, \boldsymbol{\mu}_t) - b(t, w_2, \boldsymbol{\mu}_t)\| + \|\sigma(t, w_1, \boldsymbol{\mu}_t) - \sigma(t, w_2, \boldsymbol{\mu}_t)\| \leq M_1 \|w_1 - w_2\|. \quad (\text{S30})$$

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

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  $\Phi_\delta : \mathcal{C}_{2, \delta}^p \rightarrow \mathcal{C}_{2, \delta}^p$  such that the unique fixed point  $(\boldsymbol{\lambda}_t^*)_{t \in [0, \delta]}$  is a weak solution of (S29).

Considering  $(\mathbf{W}_t^{\boldsymbol{\lambda}^*})_{t \in [0, \delta]}$ , we obtain the unique strong solution of (S29) on  $[0, \delta]$ .

Let  $\delta > 0$ . Denote  $(\boldsymbol{\lambda}_t^\mu)_{t \in [0, \delta]} \in \mathcal{P}_2(\mathbb{R}^p)^{[0, \delta]}$  such that for any  $t \in [0, \delta]$ ,  $\boldsymbol{\lambda}_t^\mu$  is the distribution of  $\mathbf{W}_t^\mu$  with initial condition  $\mathbf{W}_0^*$  with distribution  $\boldsymbol{\lambda}_0^\mu = \mu_0$ . In addition, using (S1), (S27), (S30), **B1**, **B2**, 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  $t \geq s$ ,

$$\begin{aligned}
\mathcal{W}_2(\boldsymbol{\lambda}_t^\mu, \boldsymbol{\lambda}_s^\mu)^2 &\leq \mathbb{E} \left[ \|\mathbf{W}_t^\mu - \mathbf{W}_s^\mu\|^2 \right] \\
&\leq 2\mathbb{E} \left[ \left\| \int_s^t b(u, \mathbf{W}_u^\mu, \boldsymbol{\mu}_u) du \right\|^2 \right] + 2\mathbb{E} \left[ \left\| \int_s^t \sigma(u, \mathbf{W}_u^\mu, \boldsymbol{\mu}_u) d\mathbf{B}_u \right\|^2 \right] \\
&\leq 2(t-s) \int_s^t \mathbb{E} \left[ \|b(u, \mathbf{W}_u^\mu, \boldsymbol{\mu}_u)\|^2 \right] du + 2 \int_s^t \mathbb{E} \left[ \text{Tr}(\sigma \sigma^\top(u, \mathbf{W}_u^\mu, \boldsymbol{\mu}_u)) \right] du \\
&\leq 4(t-s) \int_s^t \left\{ \|b(u, 0, \boldsymbol{\mu}_u)\|^2 + M_1^2 \mathbb{E} \left[ \|\mathbf{W}_u^\mu\|^2 \right] \right\} du
\end{aligned}$$

$$\begin{aligned}
& + 4 \int_s^t \left\{ \|\sigma(u, 0, \boldsymbol{\mu}_u)\|^2 + M_1^2 \mathbb{E} \left[ \|\mathbf{W}_u^\mu\|^2 \right] \right\} du \\
& \leq 4(1 + \delta)(t - s) \left[ M_1^2 \sup_{t \in [0, \delta]} \mathbb{E}[\|\mathbf{W}_t^\mu\|^2] + \sup_{t \in [0, \delta]} \left\{ \|b(t, 0, \boldsymbol{\mu}_t)\|^2 + \|\sigma(t, 0, \boldsymbol{\mu}_t)\|^2 \right\} \right] \\
& \leq C(t - s) \left\{ 1 + \sup_{t \in [0, \delta]} \mathbb{E}[\|\mathbf{W}_t^\mu\|^2] \right\}.
\end{aligned}$$

Therefore,  $(\boldsymbol{\lambda}_t^\mu)_{t \in [0, \delta]} \in \mathcal{C}_{2, \delta}^p$ . Let  $\Phi_\delta : \mathcal{C}_{2, \delta}^p \rightarrow \mathcal{C}_{2, \delta}^p$  given for any  $(\boldsymbol{\mu}_t)_{t \in [0, \delta]} \in \mathcal{C}_{2, \delta}^p$  by  $\Phi_\delta((\boldsymbol{\mu}_t)_{t \in [0, \delta]}) = (\boldsymbol{\lambda}_t^\mu)_{t \in [0, \delta]}$ . Let  $(\boldsymbol{\mu}_{1, t})_{t \in [0, \delta]}, (\boldsymbol{\mu}_{2, t})_{t \in [0, \delta]} \in \mathcal{C}_{2, \delta}^p$ , using (S1), (S30), B1, B2, 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  $t \in [0, \delta]$

$$\begin{aligned}
\mathbb{E} \left[ \|\mathbf{W}_t^{\mu_1} - \mathbf{W}_t^{\mu_2}\|^2 \right] & \leq 2\mathbb{E} \left[ \left\| \int_0^t \{b(s, \mathbf{W}_s^{\mu_1}, \boldsymbol{\mu}_{1, s}) - b(s, \mathbf{W}_s^{\mu_2}, \boldsymbol{\mu}_{2, s})\} ds \right\|^2 \right] \\
& \quad + 2\mathbb{E} \left[ \left\| \int_0^t \{\sigma(s, \mathbf{W}_s^{\mu_1}, \boldsymbol{\mu}_{1, s}) - \sigma(s, \mathbf{W}_s^{\mu_2}, \boldsymbol{\mu}_{2, s})\} d\mathbf{B}_s \right\|^2 \right] \\
& \leq 2\delta \int_0^t \mathbb{E} \left[ \|b(s, \mathbf{W}_s^{\mu_1}, \boldsymbol{\mu}_{1, s}) - b(s, \mathbf{W}_s^{\mu_2}, \boldsymbol{\mu}_{2, s})\|^2 \right] ds \\
& \quad + 2 \int_0^t \mathbb{E} \left[ \|\sigma(s, \mathbf{W}_s^{\mu_1}, \boldsymbol{\mu}_{1, s}) - \sigma(s, \mathbf{W}_s^{\mu_2}, \boldsymbol{\mu}_{2, s})\|^2 \right] ds \\
& \leq 4M_1^2(1 + \delta) \int_0^t \left\{ \mathbb{E} \left[ \|\mathbf{W}_s^{\mu_1} - \mathbf{W}_s^{\mu_2}\|^2 \right] + \int_Z \zeta^2(z) d\pi_Z(z) \mathcal{W}_{2, \delta}^2(\boldsymbol{\mu}_{1, s}, \boldsymbol{\mu}_{2, s}) \right\} ds \\
& \leq 4M_1^2\delta(1 + \delta) \int_Z \zeta^2(z) d\pi_Z(z) \mathcal{W}_{2, \delta}^2(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2) + 4M_1^2(1 + \delta) \int_0^t \mathbb{E} \left[ \|\mathbf{W}_s^{\mu_1} - \mathbf{W}_s^{\mu_2}\|^2 \right] ds \\
& \leq 4M_1^2\delta(1 + \delta) \exp \left[ 4M_1^2(1 + \delta)\delta \int_Z \zeta^2(z) d\pi_Z(z) \right] \mathcal{W}_{2, \delta}^2(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2).
\end{aligned}$$

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

$$\begin{aligned}
\mathcal{W}_{2, \delta}^2(\Phi_\delta(\boldsymbol{\mu}_1), \Phi_\delta(\boldsymbol{\mu}_2)) & \leq \sup_{t \in [0, \delta]} \mathbb{E} \left[ \|\mathbf{W}_t^{\mu_1} - \mathbf{W}_t^{\mu_2}\|^2 \right] \\
& \leq 4M_1^2\delta(1 + \delta) \exp \left[ 4M_1^2(1 + \delta)\delta \int_Z \zeta^2(z) d\pi_Z(z) \right] \mathcal{W}_{2, \delta}^2(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2).
\end{aligned}$$

Hence, for  $\delta > 0$  small enough,  $\Phi_\delta$  is contractive and since  $C([0, \delta], \mathcal{P}_2(\mathbb{R}^p))$  is a complete metric space, we get, using Picard fixed point theorem, that there exists a unique  $(\boldsymbol{\lambda}_t^*)_{t \in [0, \delta]} \in C([0, \delta], \mathcal{P}_2(\mathbb{R}^p))$  such that,  $\Phi_\delta(\boldsymbol{\lambda}^*) = \boldsymbol{\lambda}^*$ . For this  $\boldsymbol{\lambda}^*$ , we have that  $(\mathbf{W}_t^{\boldsymbol{\lambda}^*})_{t \in [0, \delta]}$  is a strong solution to (S29). We have shown that (S29) admits a strong solution for any initial condition  $\boldsymbol{\mu}_0 \in \mathcal{P}_2(\mathbb{R}^p)$ .

We now show that pathwise uniqueness holds for (S29). Let  $(\mathbf{W}_t^1)_{t \in [0, \delta]}$  and  $(\mathbf{W}_t^2)_{t \in [0, \delta]}$  be two strong solutions of (S29) 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  $\Phi_\delta$  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 (S30) with  $\boldsymbol{\mu} \leftarrow \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$  and since pathwise uniqueness holds for (S30), we get that  $(\mathbf{W}_t^1)_{t \in [0, \delta]} = (\mathbf{W}_t^2)_{t \in [0, \delta]}$ .

□

### S4.3 Main result

**Theorem S7.** Assume B1 and B2. For any  $N \in \mathbb{N}^*$ , let  $(\mathbf{W}_t^{1:N})_{t \geq 0}$  be a strong solution of (S26) and for any  $N \in \mathbb{N}^*$  and  $k \in \{1, \dots, N\}$ , let  $(\mathbf{W}_t^{k,*})_{t \geq 0}$  be a strong solution of (S29) with

*Brownian motion*  $(\mathbf{B}_t^k)_{t \geq 0}$ . Assume that there exists  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^p)$  such that for any  $N \in \mathbb{N}^*$ ,  $\mathbf{W}_0^{1:N} = \mathbf{W}_0^{*,1:N}$  has distribution  $\mu_0^{\otimes N}$ . Then for any  $T \geq 0$ ,  $N \in \mathbb{N}^*$  and  $k \in \{1, \dots, N\}$

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

*Proof.* Let  $T \geq 0$ . For any  $N \in \mathbb{N}^*$ ,  $t \geq 0$ , let  $\nu_t^{*,N} = (1/N) \sum_{j=1}^N \delta_{\mathbf{W}_t^{*,j}}$ . Using **B1**, **B2**, 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  $N \in \mathbb{N}^*$  and  $k \in \{1, \dots, N\}$

$$\begin{aligned} \mathbb{E} \left[ \sup_{t \in [0, T]} \left\| \mathbf{W}_t^{k, N} - \mathbf{W}_t^{k, *} \right\|^2 \right] &\leq 2\mathbb{E} \left[ \sup_{t \in [0, T]} \left\| \int_0^t (b_N(s, \mathbf{W}_s^{k, N}, \nu_s^N) - b(s, \mathbf{W}_s^{k, *}, \lambda_s^*)) ds \right\|^2 \right] \\ &\quad + 2\mathbb{E} \left[ \sup_{t \in [0, T]} \left\| \int_0^t (\sigma_N(s, \mathbf{W}_s^{k, N}, \nu_s^N) - \sigma(s, \mathbf{W}_s^{k, *}, \lambda_s^*)) 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}, \nu_s^N) - b(s, \mathbf{W}_s^{k, *}, \lambda_s^*) \right\|^2 \right] ds \\ &\quad + 2\mathbb{E} \left[ \left\| \int_0^T (\sigma_N(s, \mathbf{W}_s^{k, N}, \nu_s^N) - \sigma(s, \mathbf{W}_s^{k, *}, \lambda_s^*)) 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}, \nu_s^N) - b(s, \mathbf{W}_s^{k, *}, \lambda_s^*) \right\|^2 \right] \right. \\ &\quad \left. + \mathbb{E} \left[ \left\| \sigma_N(s, \mathbf{W}_s^{k, N}, \nu_s^N) - \sigma(s, \mathbf{W}_s^{k, *}, \lambda_s^*) \right\|^2 \right] \right\} ds \\ &\leq 8\mathbf{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}, \nu_s^N) - b(s, \mathbf{W}_s^{k, *}, \lambda_s^*) \right\|^2 \right] \right. \\ &\quad \left. + \mathbb{E} \left[ \left\| \sigma(s, \mathbf{W}_s^{k, N}, \nu_s^N) - \sigma(s, \mathbf{W}_s^{k, *}, \lambda_s^*) \right\|^2 \right] \right\} ds \\ &\leq 8\mathbf{M}_2^2(1+T)^2 N^{-2\kappa} + 8\mathbf{M}_1^2(1+T) \\ &\quad \times \int_0^T \left\{ \int_{\mathbb{Z}} \mathbb{E} \left[ \left\| \nu_s^N[g(\cdot, z)] - \lambda_s^*[g(\cdot, z)] \right\|^2 \right] d\pi_{\mathbb{Z}}(z) + \mathbb{E} \left[ \left\| \mathbf{W}_s^{k, N} - \mathbf{W}_s^{k, *} \right\|^2 \right] \right\} ds \\ &\leq 8\mathbf{M}_2^2(1+T)^2 N^{-2\kappa} + 16\mathbf{M}_1^2(1+T) \\ &\quad \times \int_0^T \left\{ \int_{\mathbb{Z}} \left( \mathbb{E} \left[ \left\| \nu_s^N[g(\cdot, z)] - \nu_s^{*,N}[g(\cdot, z)] \right\|^2 \right] + \mathbb{E} \left[ \left\| \nu_s^{*,N}[g(\cdot, z)] - \lambda_s^*[g(\cdot, z)] \right\|^2 \right] \right) d\pi_{\mathbb{Z}}(z) \right. \\ &\quad \left. + \mathbb{E} \left[ \left\| \mathbf{W}_s^{k, N} - \mathbf{W}_s^{k, *} \right\|^2 \right] \right\} ds. \end{aligned}$$

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  $\tau : \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ ,  $\{(\mathbf{W}_t^{k, N})_{t \geq 0}\}_{k=1}^N$  has the same distribution as  $\{(\mathbf{W}_t^{\tau(k), N})_{t \geq 0}\}_{k=1}^N$  and  $\{(\mathbf{W}_t^{k, *})_{t \geq 0}\}_{k=1}^N$  are independent we have

$$\begin{aligned} \mathbb{E} \left[ \sup_{t \in [0, T]} \left\| \mathbf{W}_t^{k, N} - \mathbf{W}_t^{k, *} \right\|^2 \right] &\leq 8\mathbf{M}_2^2(1+T)^2 N^{-2\kappa} + 16\mathbf{M}_1^2(1+T) \\ &\quad \times \int_0^T \left\{ \frac{1}{N} \int_{\mathbb{Z}} \zeta^2(z) d\pi_{\mathbb{Z}}(z) \sum_{j=1}^N \mathbb{E} \left[ \left\| \mathbf{W}_s^{j, N} - \mathbf{W}_s^{j, *} \right\|^2 \right] + \mathbb{E} \left[ \left\| \mathbf{W}_s^{k, N} - \mathbf{W}_s^{k, *} \right\|^2 \right] \right. \\ &\quad \left. + \int_{\mathbb{Z}} \mathbb{E} \left[ \left\| \frac{1}{N} \sum_{j=1}^N g(\mathbf{W}_s^{j, *}, z) - \int_{\mathbb{R}^p} g(\bar{w}, z) d\lambda_s^*(\bar{w}) \right\|^2 \right] d\pi_{\mathbb{Z}}(z) \right\} ds \end{aligned}$$

$$\begin{aligned}
&\leq 8M_2^2(1+T)^2N^{-2\kappa} + 16M_1^2(1+T) \left(1 + \int_Z \zeta^2(z)d\pi_Z(z)\right) \int_0^T \mathbb{E} \left[ \|\mathbf{W}_s^{k,N} - \mathbf{W}_s^{k,*}\|^2 \right] ds \\
&\quad + 16M_1^2(1+T)N^{-1} \int_0^T \int_Z \mathbb{E} \left[ \left\| g(\mathbf{W}_s^{k,*}, z) - \int_{\mathbb{R}^p} g(\bar{w}, z) d\lambda_s^*(\bar{w}) \right\|^2 \right] d\pi_Z(z) ds \\
&\leq 8M_2^2(1+T)^2N^{-2\kappa} + 16M_1^2(1+T) \left(1 + \int_Z \zeta^2(z)d\pi_Z(z)\right) \int_0^T \mathbb{E} \left[ \|\mathbf{W}_s^{k,N} - \mathbf{W}_s^{k,*}\|^2 \right] ds \\
&\quad + 32M_1^2(1+T)^2N^{-1} \left(1 + \int_Z \zeta^2(z)d\pi_Z(z)\right) .
\end{aligned}$$

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

#### S4.4 Proofs of the main results

In this section we prove Theorem 1, Theorem 2, Theorem S2, Theorem S3. Note that we only need to show Theorem S2 and Theorem S3, since in the case  $\eta = 0$ , Theorem S2 boils down to Theorem 1 and Theorem S3 to Theorem 2.

*Proof of Theorem S2.* Define for any  $N \in \mathbb{N}^*$ ,  $w \in \mathbb{R}^p$ ,  $\mu \in \mathcal{P}_2(\mathbb{R}^p)$  and  $t \geq 0$

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

with  $h$  and  $\Sigma$  given in (S7). Using Proposition S4, we get that **B1** holds with  $M_1 \leftarrow L$  and  $\gamma_{\alpha, \beta}(N) = \gamma^{1/(1-\alpha)} N^{(\beta-1)/(1-\alpha)}$ . In addition, using Proposition S4, **B2** holds with  $M_2 \leftarrow (\gamma^{1-\alpha}/M)^{1/2} pL$  and  $2\kappa = (1-\beta)/(1-\alpha)$ . We conclude using Theorem S7.  $\square$

*Proof of Theorem S3.* Define for any  $N \in \mathbb{N}^*$ ,  $w \in \mathbb{R}^p$ ,  $\mu \in \mathcal{P}_2(\mathbb{R}^p)$  and  $t \geq 0$

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

with  $h$  and  $\Sigma$  given in (S7). Using Proposition S4, we get that **B1** holds with  $M_1 \leftarrow L$ . In addition, **B2** holds with  $b = b_N$ ,  $\sigma = \sigma_N$ ,  $M_2 \leftarrow 0$  and  $\kappa = 0$ . We conclude using Theorem S7.  $\square$

*Proof of Proposition 4.* We consider only the case  $\beta = 1$ , the proof for  $\beta \in [0, 1)$  following the same lines. Let  $M \in \mathbb{N}^*$ . We have for any  $N \in \mathbb{N}^*$  using Proposition S1,

$$\begin{aligned}
\mathcal{W}_2(\Upsilon^N, \delta_{\lambda^*})^2 &\leq \mathbb{E} [\mathcal{W}_2(\nu^N, \lambda^*)^2] \\
&\leq N^{-1} \sum_{k=1}^N \mathbb{E} [\mathcal{W}_2(\delta_{(\mathbf{W}_t^{k,N})_{t \geq 0}}, \lambda^*)^2] \leq N^{-1} \sum_{k=1}^N \mathbb{E} [\text{m}^2((\mathbf{W}_t^{k,N})_{t \geq 0}, (\mathbf{W}_t^{k,*})_{t \geq 0})] . \quad (\text{S31})
\end{aligned}$$

Let  $\varepsilon > 0$  and  $n_0$  such that  $\sum_{n=n_0+1}^{+\infty} 2^{-n} \leq \varepsilon$ . Combining (S31), Theorem 1 and the Cauchy-Schwarz inequality we get that for any  $N \in \mathbb{N}^*$

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

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

## S5 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 **A1-(c)**, the fact that  $V$  is  $\eta$ -strongly convex and [9, Theorem 2.1.5, Theorem 2.1.7] there exists  $M \geq 0$  such that for any  $w \in \mathbb{R}$  we have

$$\eta(w - w_0)^2/2 \leq V(w) - V(w_0) \leq M(w - w_0)^2/2 . \quad (\text{S32})$$

In addition, using Proposition S4, we have for any  $\mu \in \mathcal{P}_2(\mathbb{R})$  and  $w \in \mathbb{R}$ ,

$$\bar{\sigma}^2 \leq \Sigma(w, \mu) \leq L^2. \quad (\text{S33})$$

Recall that for any  $\mu \in \mathcal{P}_2(\mathbb{R})$  and  $w \in \mathbb{R}$ ,  $h(w, \mu) = \bar{h}(w, \mu) + V'(w)$ , with  $\bar{h}$  given in (S7). 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 S4, (S32) and (S33), there exists  $m_1 > 0$  and  $c_1 \in \mathbb{R}$  such that for any  $\mu \in \mathcal{P}_2(\mathbb{R})$  and  $w \in \mathbb{R}$ , we have distinguishing the case  $w \leq w_0$  and  $w > w_0$ ,

$$\begin{aligned} \int_0^w \{h/\Sigma\}(\tilde{w}, \mu) d\tilde{w} &\geq -\bar{\sigma}^{-2} L^2 |w| + \int_0^w V'(\tilde{w})/\Sigma(\tilde{w}, \mu) d\tilde{w} \\ &\geq -\bar{\sigma}^{-2} L^2 |w| - \bar{\sigma}^{-2} \sup_{\tilde{w} \in [0, w_0]} |V'(\tilde{w})| |w_0| + \int_{w_0}^w V'(\tilde{w})/\Sigma(\tilde{w}, \mu) d\tilde{w} \\ &\geq -\bar{\sigma}^{-2} L^2 |w| - \bar{\sigma}^{-2} \sup_{\tilde{w} \in [0, w_0]} |V'(\tilde{w})| |w_0| + (V(w) - V(w_0))L^{-2} \geq m_1 w^2 + c_1. \end{aligned} \quad (\text{S34})$$

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

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

where  $\bar{\Sigma}(w, \mu) = \gamma^{1/(1-\alpha)} \Sigma(w, \mu)/M$ . Similarly to (S34), there exist  $m_2 > 0$  and  $c_2 \in \mathbb{R}$  such that for any  $\mu \in \mathcal{P}_2(\mathbb{R})$  and  $w \in \mathbb{R}$

$$\int_0^w h(\tilde{w}, \mu)/\Sigma(\tilde{w}, \mu) d\tilde{w} \leq m_2 w^2 + c_2. \quad (\text{S35})$$

Combining (S33), (S34) and (S35), there exists  $m > 0$  and  $c \in \mathbb{R}$  such that for any  $\mu \in \mathcal{P}_2(\mathbb{R})$  and  $w \in \mathbb{R}$ ,  $\rho_\mu(w) \leq ce^{-mw^2}$ . Using this result, we get that  $\sup_{\mu \in \mathcal{P}_2(\mathbb{R})} \int_{\mathbb{R}} w^4 \rho_\mu(w) dw < +\infty$ . Therefore, using [10, Theorem 2.7] we obtain that  $H(\mathcal{P}_2(\mathbb{R}))$  is relatively compact in  $(\mathcal{P}_2(\mathbb{R}), \mathcal{W}_2)$ .

We now show that  $H \in C(\mathcal{P}_2(\mathbb{R}), \mathcal{P}_2(\mathbb{R}))$ . Let  $\mu \in \mathcal{P}_2(\mathbb{R})$  and  $(\mu_n)_{n \in \mathbb{N}} \in \mathcal{P}_2(\mathbb{R})^{\mathbb{N}}$  such that  $\lim_{n \rightarrow +\infty} \mu_n = \mu$ . Using Proposition S4 and the Lebesgue dominated convergence theorem we obtain that for any  $w \in \mathbb{R}$ ,  $\lim_{n \rightarrow +\infty} \rho_{\mu_n}(w) = \rho_\mu(w)$ . Using Scheffé's lemma we get that  $\lim_{n \rightarrow +\infty} \int_{\mathbb{R}} |\rho_{\mu_n}(w) - \rho_\mu(w)| dw = 0$ . 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  $(\mathcal{P}_2(\mathbb{R}), \mathcal{W}_2)$ . Therefore,  $(H(\mu_{n_k}))_{k \in \mathbb{N}}$  also weakly converges and we obtain that  $\lim_{k \rightarrow +\infty} \mathcal{W}_2(H(\mu_{n_k}), H(\mu)) = 0$ . Since  $\{H(\mu_n) : n \in \mathbb{N}\}$  is relatively compact and admits a unique limit point we obtain that  $\lim_{n \rightarrow +\infty} \mathcal{W}_2(H(\mu_n), H(\mu)) = 0$ .

Hence  $H \in C(\mathcal{P}_2(\mathbb{R}), \mathcal{P}_2(\mathbb{R}))$ . Therefore, since  $H \in C(\mathcal{P}_2(\mathbb{R}), \mathcal{P}_2(\mathbb{R}))$  and  $H(\mathcal{P}_2(\mathbb{R}))$  is relatively compact in  $\mathcal{P}_2(\mathbb{R})$  Schauder's theorem [11, Appendix] implies that  $H$  admits a fixed point.

Let  $\mu \in \mathcal{P}_2(\mathbb{R})$  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

$$d\mathbf{W}_t^\mu = h(t, \mu) dt + \gamma^{1/(1-\alpha)} \Sigma(\mathbf{W}_t^\mu, \mu) d\mathbf{B}_t. \quad (\text{S36})$$

An invariant distribution for (S36) is given by  $H(\mu)$ , see [12]. 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.  $\square$

## S6 Links with gradient flow approach

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

$$\int_{\mathbb{R}^p} f(\tilde{w}) d\lambda_t^*(\tilde{w}) = \int_{\mathbb{R}^p} f(\tilde{w}) d\mu_0(\tilde{w}) + \int_0^t \int_{\mathbb{R}^p} (s+1)^{-\alpha} \langle h(\tilde{w}, \lambda_s^*), \nabla f(\tilde{w}) \rangle d\lambda_s^*(\tilde{w}). \quad (\text{S37})$$

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

$$\partial_t \rho_t^*(w) = -(t+1)^{-\alpha} \operatorname{div}(\bar{h}(\cdot, \rho_t^*) \rho_t^*)(w) ,$$

with for any  $w \in \mathbb{R}^p$  and  $\mu \in \mathcal{P}(\mathbb{R}^p)$  with density  $\rho$ ,  $h(w, \mu) = \bar{h}(w, \rho)$ . In the case  $\alpha = 0$ , it is well-known, see [13, 4, 14], that  $(\rho_t^*)_{t \geq 0}$  is a Wasserstein gradient flow for the functional  $\mathcal{R}^* : \mathcal{P}_2^c(\mathbb{R}^p) \rightarrow \mathbb{R}$  given for any  $\rho \in \mathcal{P}_2^c(\mathbb{R}^p)$

$$\mathcal{R}^*(\rho) = \int_{\mathbb{X} \times \mathbb{Y}} \ell \left( \int_{\mathbb{R}^p} F(\tilde{w}, x) \rho(\tilde{w}) d\tilde{w}, y \right) d\pi(x, y) , \quad (\text{S38})$$

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

**Case  $\beta = 1$**  Focusing on  $(\lambda_t^*)_{t \geq 0}$ , we no longer obtain that  $(\lambda_t^*)_{t \geq 0}$  is a gradient flow for (S38). Indeed, using Itô's formula, we have the following evolution equation for any  $f \in C_c^2(\mathbb{R}^p)$  and  $t \geq 0$

$$\begin{aligned} \int_{\mathbb{R}^p} f(\tilde{w}) d\lambda_t^*(\tilde{w}) &= \int_{\mathbb{R}^p} f(\tilde{w}) d\mu_0(\tilde{w}) + \int_0^t \int_{\mathbb{R}^p} (s+1)^{-\alpha} \langle h(\tilde{w}, \lambda_s^*), \nabla f(\tilde{w}) \rangle d\lambda_s^*(\tilde{w}) \\ &\quad + \int_0^t \int_{\mathbb{R}^p} (s+1)^{-\alpha} \operatorname{Tr}(\Sigma(\tilde{w}, \lambda_s^*) \nabla^2 f(\tilde{w})) d\tilde{w} . \end{aligned} \quad (\text{S39})$$

We highlight that the additional term in (S39) from (S37) corresponds to some entropic regularization of the risk  $\mathcal{R}^*$ . Indeed, if for any  $w \in \mathbb{R}^p$  and  $\mu \in \mathcal{P}(\mathbb{R}^p)$ ,  $\Sigma = \beta \operatorname{Id}$  then, in the case  $\alpha = 0$ , we obtain that  $(\rho_t^*)_{t \geq 0}$  is a gradient flow for  $\rho \mapsto U^*(\rho) + \beta \operatorname{Ent}(\rho)$ , where  $\operatorname{Ent} : \mathbb{K}_2 \rightarrow \mathbb{R}$  is given for any  $\rho \in \mathbb{K}_2$  by

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

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

## S7 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 [15] and  $N_{\text{input}} = 32 \times 32 \times 3$  in the case of the CIFAR-10 dataset [16]. 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 S1 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 S2 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 [13].

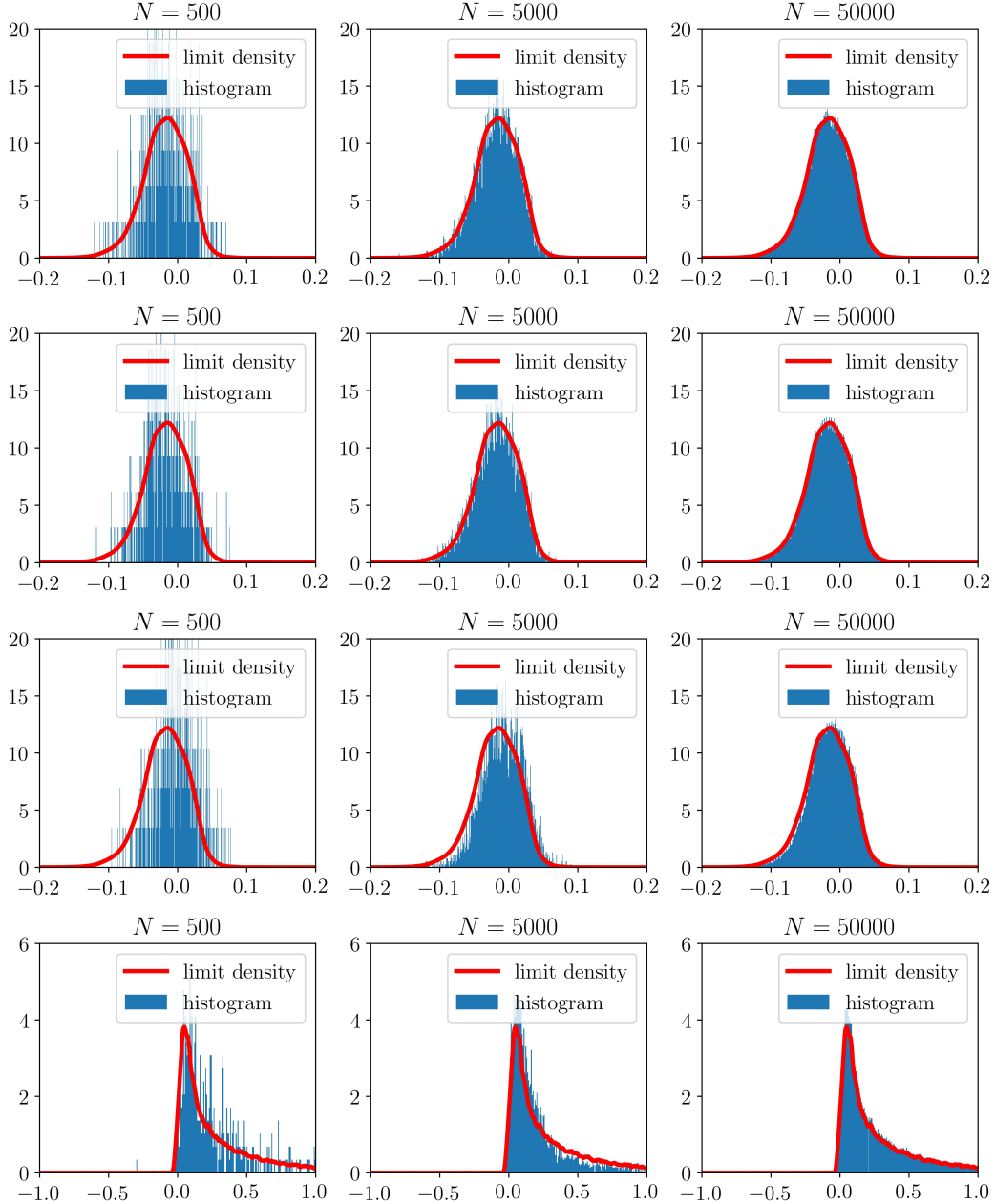


Figure S1: Convergence of the weights of the first layer as  $N \rightarrow +\infty$  for  $\alpha = 0$  and  $M = 100$ . The first line corresponds to  $\beta = 0.25$ , the second to  $\beta = 0.5$ , the third to  $\beta = 0.75$  and the last line to  $\beta = 1.0$ .

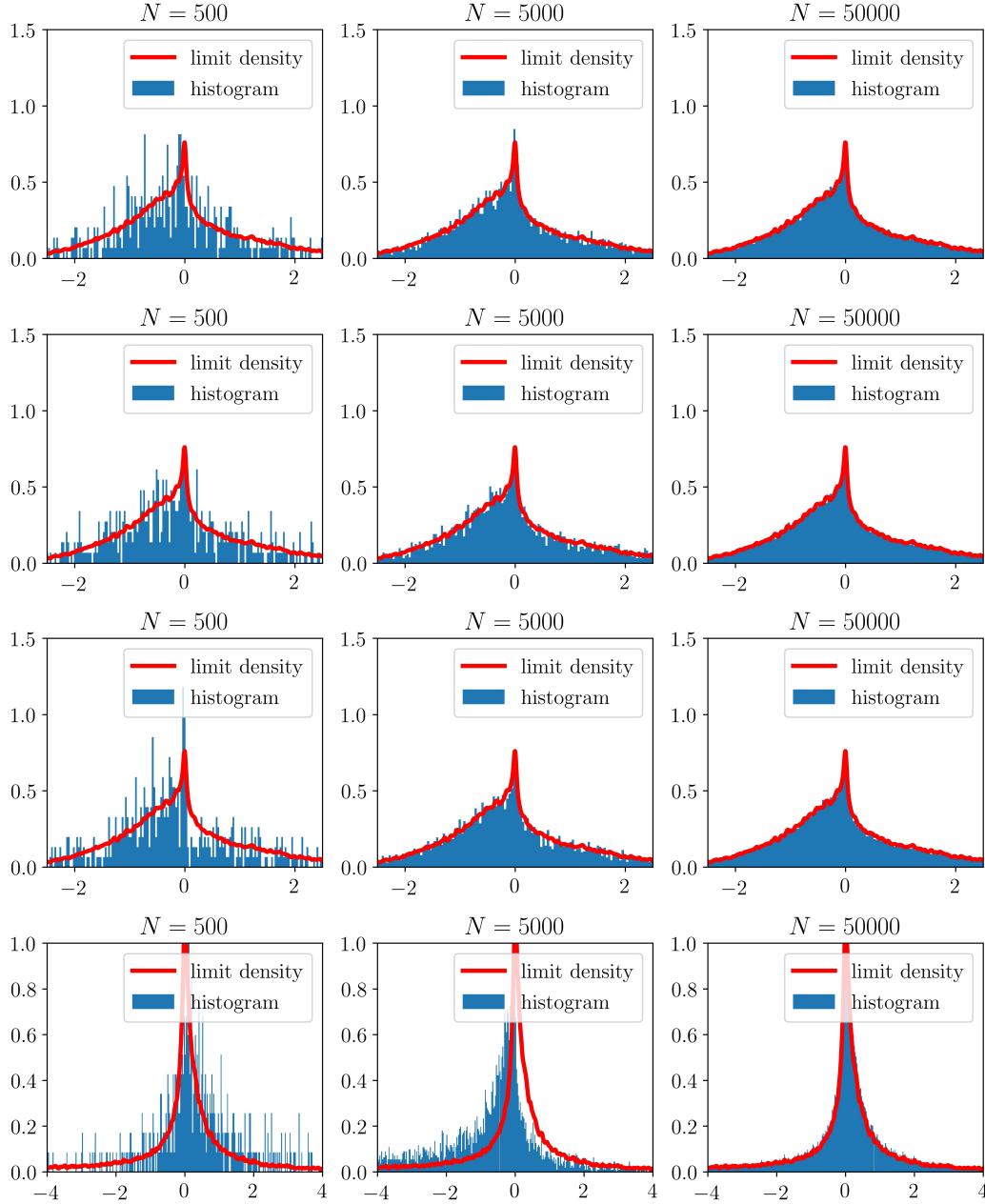


Figure S2: Convergence of the weights of the first layer as  $N \rightarrow +\infty$  for  $\alpha = 0$  and  $M = 100$ . The first line corresponds to  $\beta = 0.25$ , the second to  $\beta = 0.5$ , the third to  $\beta = 0.75$  and the last line to  $\beta = 1.0$ .

On Figure S3 and Figure S4 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 (S37) or (S39) (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 S3 and Figure S4).

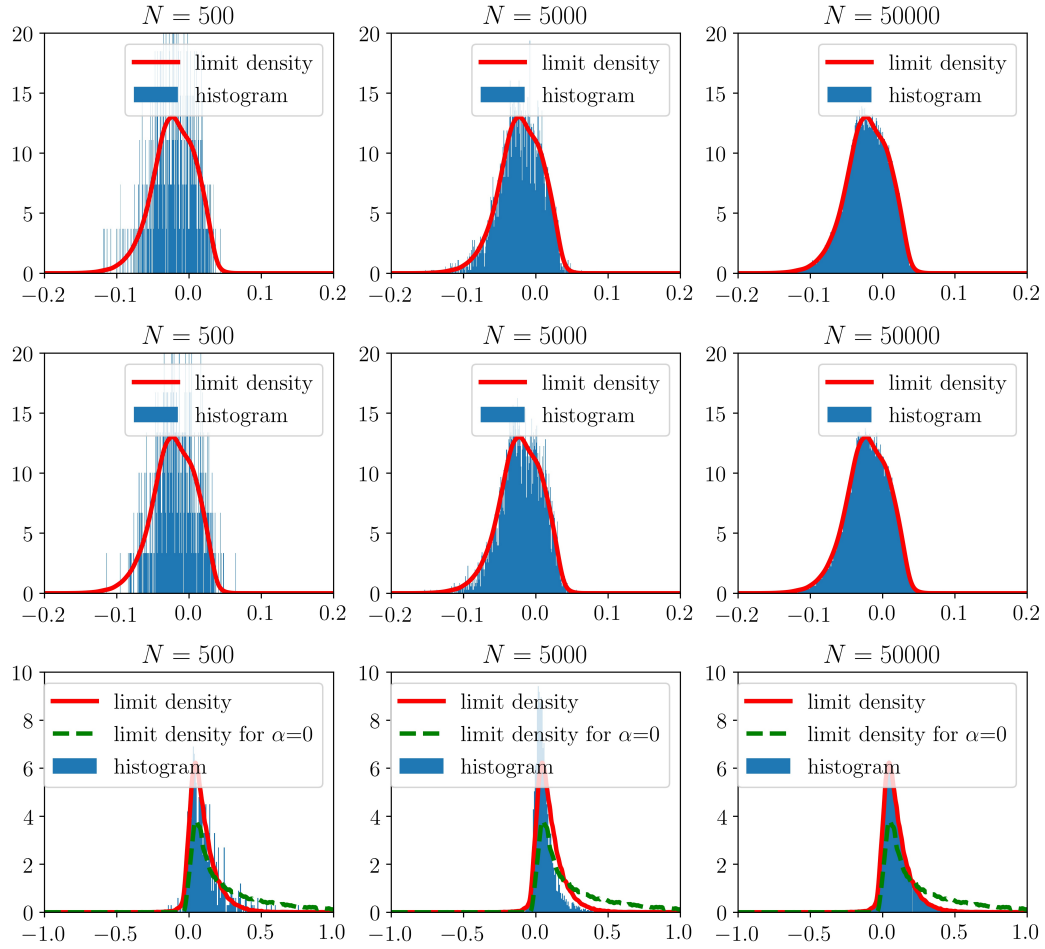


Figure S3: Convergence of the weights of the first layer as  $N \rightarrow +\infty$  for  $\alpha = 0.25$  and  $M = 100$ . The first line corresponds to  $\beta = 0.5$ , the second to  $\beta = 0.75$  and the last line to  $\beta = 1.0$ .

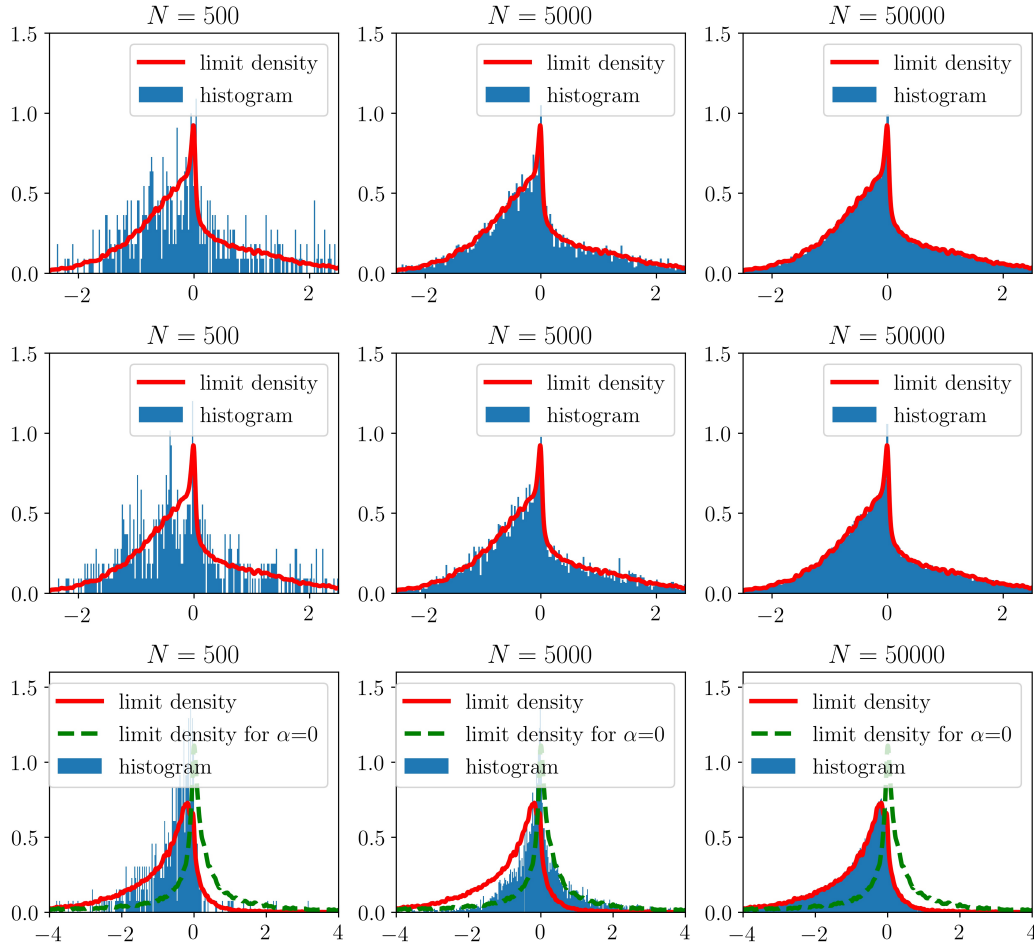


Figure S4: Convergence of the weights of the first layer as  $N \rightarrow +\infty$  for  $\alpha = 0.25$  and  $M = 100$ . The first line corresponds to  $\beta = 0.5$ , the second to  $\beta = 0.75$  and the last line to  $\beta = 1.0$ .

We now study the role of the batch size  $M$  on the convergence toward the mean-field regime. Figure S5 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  $M$ .

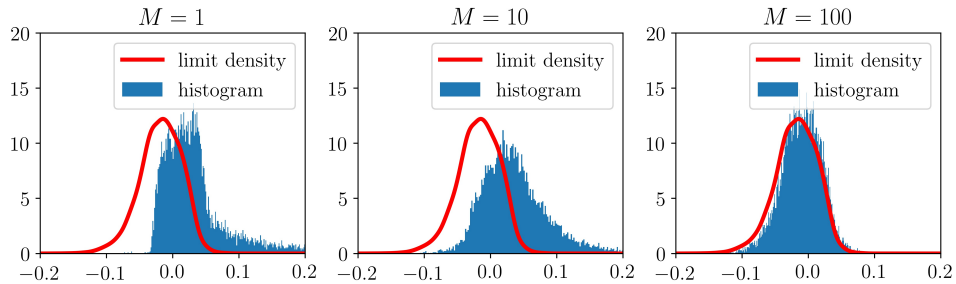


Figure S5: Convergence of the weights as  $M \rightarrow \infty$

## References

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

- [2] A. Kechris, *Classical Descriptive Set Theory*. Graduate Texts in Mathematics, Springer New York, 2012.
- [3] 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.
- [4] 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.
- [5] D. W. Stroock and S. R. S. Varadhan, *Multidimensional diffusion processes*. Classics in Mathematics, Springer-Verlag, Berlin, 2006. Reprint of the 1997 edition.
- [6] 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.
- [7] 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.
- [8] A. Sznitman, “Topics in propagation of chaos,” in *Ecole d’été de probabilités de Saint-Flour XIX—1989*, pp. 165–251, Springer, 1991.
- [9] Y. Nesterov, *Introductory lectures on convex optimization*, vol. 87 of *Applied Optimization*. Kluwer Academic Publishers, Boston, MA, 2004. A basic course.
- [10] 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.
- [11] F. F. Bonsall and K. Vedak, *Lectures on some fixed point theorems of functional analysis*. No. 26, Tata Institute of Fundamental Research Bombay, 1962.
- [12] J. Kent, “Time-reversible diffusions,” *Advances in Applied Probability*, vol. 10, pp. 819–835, 12 1978.
- [13] 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.
- [14] J. Sirignano and K. Spiliopoulos, “Mean field analysis of neural networks,” *arXiv preprint arXiv:1805.01053*, 2018.
- [15] Y. LeCun and C. Cortes, “MNIST handwritten digit database,” 2010.
- [16] A. Krizhevsky, G. Hinton, *et al.*, “Learning multiple layers of features from tiny images,” 2009.