



**HAL**  
open science

# Law of Large Numbers for Bayesian two-layer Neural Network trained with Variational Inference

Arnaud Descours, Tom Huix, Arnaud Guillin, Manon Michel, Éric Moulines,  
Boris Nectoux

► **To cite this version:**

Arnaud Descours, Tom Huix, Arnaud Guillin, Manon Michel, Éric Moulines, et al.. Law of Large Numbers for Bayesian two-layer Neural Network trained with Variational Inference. The Thirty Sixth Annual Conference on Learning Theory, Jul 2023, Bangalore, India. pp.4657-4695. hal-04153801

**HAL Id: hal-04153801**

**<https://hal.science/hal-04153801v1>**

Submitted on 6 Jul 2023

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

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

# Law of Large Numbers for Bayesian two-layer Neural Network trained with Variational Inference

**Arnaud Descours**

ARNAUD.DESCOURS@UCA.FR

*Laboratoire de Mathématiques Blaise Pascal, UMR 6620, CNRS, Université Clermont Auvergne, Aubière, France.*

**Tom Huix**

TOM.HUIX@POLYTECHNIQUE.EDU

*Centre de Mathématiques Appliquées, UMR 7641, Ecole polytechnique, France.*

**Arnaud Guillin**

ARNAUD.GUILLIN@UCA.FR

*Laboratoire de Mathématiques Blaise Pascal, UMR 6620, CNRS, Université Clermont Auvergne, Aubière, France.*

**Manon Michel**

MANON.MICHEL@UCA.FR

*Laboratoire de Mathématiques Blaise Pascal, UMR 6620, CNRS, Université Clermont Auvergne, Aubière, France.*

**Éric Moulines**

ERIC.MOULINES@POLYTECHNIQUE.EDU

*Centre de Mathématiques Appliquées, UMR 7641, Ecole polytechnique, France.*

**Boris Nectoux**

BORIS.NECTOUX@UCA.FR

*Laboratoire de Mathématiques Blaise Pascal, UMR 6620, CNRS, Université Clermont Auvergne, Aubière, France.*

**Editors:** Gergely Neu and Lorenzo Rosasco

## Abstract

We provide a rigorous analysis of training by variational inference (VI) of Bayesian neural networks in the two-layer and infinite-width case. We consider a regression problem with a regularized evidence lower bound (ELBO) which is decomposed into the expected log-likelihood of the data and the Kullback-Leibler (KL) divergence between the a priori distribution and the variational posterior. With an appropriate weighting of the KL, we prove a law of large numbers for three different training schemes: (i) the idealized case with exact estimation of a multiple Gaussian integral from the reparametrization trick, (ii) a minibatch scheme using Monte Carlo sampling, commonly known as *Bayes by Backprop*, and (iii) a new and computationally cheaper algorithm which we introduce as *Minimal VI*. An important result is that all methods converge to the same mean-field limit. Finally, we illustrate our results numerically and discuss the need for the derivation of a central limit theorem.

**Keywords:** Bayesian neural networks, variational inference, mean-field, law of large numbers, infinite-width neural networks.

## 1. Introduction

Deep Learning has led to a revolution in machine learning with impressive successes. However, some limitations of DL have been identified and, despite, many attempts, our understanding of DL is still limited. A long-standing problem is the assessment of predictive uncertainty: DL tends to be overconfident in its predictions [Abdar et al. \(2021\)](#), which is a problem in applications such as autonomous driving ([McAllister et al., 2017](#); [Michelmore et al., 2020](#)), medical diagnosis ([Kendall and Gal, 2017](#); [Filos et al., 2019](#)), or finance; cf [Krzywinski and Altman \(2013\)](#); [Ghahramani \(2015\)](#). Therefore, on the one hand, analytical efforts are being made to thoroughly investigate the performance of DL; and on the other hand, many approaches have been proposed to alleviate its shortcomings. The Bayesian paradigm is an attractive way to tackle predictive uncertainty, as it provides a framework for training uncertainty-aware neural networks (NNs) (e.g. [Ghahramani \(2015\)](#); [Blundell et al. \(2015\)](#); [Gal and Ghahramani \(2016\)](#)).

Thanks to a fully probabilistic approach, Bayesian Neural Networks (BNN) combine the impressive neural-network expressivity with the decision-theoretic approach of Bayesian inference, making them capable of providing predictive uncertainty; see [Blundell et al. \(2015\)](#); [Michelmore et al. \(2020\)](#); [McAllister et al. \(2017\)](#); [Filos et al. \(2019\)](#). However, Bayesian inference requires deriving the posterior distribution of the NN weights. This posterior distribution is typically not tractable. A classical approach is to sample the posterior distribution using Markov chain Monte Carlo methods (such as Hamilton-Monte-Carlo methods). There are however long-standing difficulties, such as the proper choice of the prior and fine-tuning of the sampler. Such difficulties often become prohibitive in large-dimensional cases, ([Cobb and Jalaian, 2021](#)). An alternative is to use variational inference, which has a long history ([Hinton and Camp, 1993](#); [MacKay, 1995](#); [MacKay et al., 1995](#)). Simpler methods that do not require exact computation of integrals over the variational posterior were then developed, e.g. first by [Graves \(2011\)](#) thanks to some approximation and then by [Blundell et al. \(2015\)](#) with the *Bayes by Backprop* approach. In the latter, the posterior distribution is approximated by a parametric distribution and a generalisation of the reparametrization trick used by [Kingma and Welling \(2014\)](#) leads to an unbiased estimator of the gradient of the ELBO; see also [Gal and Ghahramani \(2016\)](#); [Louizos and Welling \(2017\)](#); [Khan et al. \(2018\)](#). Despite the successful application of this approach, little is known about the overparameterized limit and appropriate weighting that must be assumed to obtain a nontrivial Bayesian posterior, see [Izmailov et al. \(2021\)](#). Recently, [Huix et al. \(2022\)](#) outlined the importance of balancing in ELBO the integrated log-likelihood term and the KL regularizer, to avoid both overfitting and dominance of the prior. However, a suitable limiting theory has yet to be established, as well as guarantees for the practical implementation of the stochastic gradient descent (SGD) used to estimate the parameters of the variational distribution.

Motivated by the need to provide a solid theoretical framework, asymptotic analysis of NN has gained much interest recently. The main focus has been on the gradient descent algorithm and its variants ([Rotskoff and Vanden-Eijnden, 2018](#); [Chizat and Bach, 2018](#); [Mei et al., 2018](#); [Sirignano and Spiliopoulos, 2020](#); [Descours et al., 2022](#)). In much of these works, a mean-field analysis is performed to characterize the limiting nonlinear evolution of the weights of a two-layer NN, allowing the derivation of a law of large numbers and a central limit theorem for the empirical distribution of neuron weights. A long-term goal of these works is to demonstrate convergence toward a global minimum of these limits for the mean field. Despite some progress in this direction, this is still an open and highly challenging problem; cf [Chizat and Bach \(2018\)](#); [Chizat \(2022\)](#);

Chizat et al. (2022). Nevertheless, this asymptotic analysis is also of interest in its own right, as we show here in the case of variational inference for Bayesian neural networks. Indeed, based on this asymptotic analysis, we develop an efficient and new variant of the stochastic gradient descent (SGD) algorithm for variational inference in BNN that computes only the information necessary to recover the limit behavior.

Our goal, then, is to work at the intersection of analytical efforts to gain theoretical guarantees and insights and of practical methods for a workable variational inference procedure. By adapting the framework developed by Descours et al. (2022), we produce a rigorous asymptotic analysis of BNN trained in a variational setting for a regression task. From the limit equation analysis, we first find that a proper regularisation of the Kullback-Leibler divergence term in relation with the integrated loss leads to their right asymptotic balance. Second, we prove the asymptotic equivalence of the idealized and Bayes-by-Backprop SGD schemes, as both preserve the same core contributions to the limit. Finally, we introduce a computationally more favourable scheme, directly stemming from the effective asymptotic contributions. This scheme is the true mean-field algorithmic approach, as only deriving from non-interacting terms.

More specifically, our contributions are the following:

- We first focus on the idealized SGD algorithm, where the variational expectations of the derivative of the loss from the reparametrization trick of Blundell et al. (2015) are computed exactly. More precisely, we prove that with the number of neurons  $N \rightarrow +\infty$ , the sequence of trajectories of the scaled empirical distributions of the parameters satisfies a law of large numbers. This is the purpose of Theorem 2. The proof is completely new: it establishes directly the limit in the topology inherited by the Wasserstein distance bypassing the highly technical Sobolev space arguments used in Descours et al. (2022).

The idealized SGD requires the computation of some integrals, which in practice prevents a direct application of this algorithm. However, we can prove its convergence to an explicit nonlinear process. These integrals are usually obtained by a Monte Carlo approximation, leading to the *Bayes-by-Backprop* SGD, see Blundell et al. (2015).

- We show for the *Bayes-by-Backprop* SGD (see Theorem 3) that the sequence of trajectories of the scaled empirical distributions of the parameters satisfies the same law of large numbers as that in Theorem 2, which justifies such an approximation procedure. Note that each step of the algorithm involves the simulation of  $O(N)$  Gaussian random variables, which can make the associated gradient evaluation prohibitively expensive.
- A careful analysis of the structure of the limit equation (11) allows us to develop a new algorithm, called *Minimal-VI* SGD, which at each step generates only two Gaussian random variables and for which we prove the same limiting behavior. The key idea here is to keep only those contributions which affect the asymptotic behavior and which can be understood as the mean-field approximation from the uncorrelated degrees of freedom. This is all the more interesting since we observe numerically that the number weights  $N$  required to reach this asymptotic limit is quite small which makes this variant of immediate practical interest.
- We numerically investigate the convergence of the three methods to the common limit behavior on a toy example. We observe that the mean-field method is effective for a small number of neurons ( $N = 300$ ). The differences between the methods are reflected in the variances.

The paper is organized as follows: Section 2 introduces the variational inference in BNN, as well as the SGD schemes commonly considered, namely the idealized and *Bayes-by-backprop* variants. Then, in Section 3 we establish our initial result, the LLN for the idealized SGD. In Section 4 we prove the LLN for the *Bayes-by-backprop* SGD and its variants. We show that both SGD schemes have the same limit behavior. Based on an analysis of the obtained limit equation, we present in Section 5 the new *minimal-VI*. Finally, in Section 6 we illustrate our findings using numerical experiments. The proofs of the mean-field limits, which are original and quite technically demanding, are gathered in the supplementary paper.

**Related works.** Law of Large Numbers (LLN) for mean-field interacting particle systems, have attracted a lot of attentions; see for example Hitsuda and Mitoma (1986); Sznitman (1991); Fernandez and Méléard (1997); Jourdain and Méléard (1998); Delarue et al. (2019); Del Moral and Guionnet (1999); Kurtz and Xiong (2004) and references therein. The use of mean-field particle systems to analyse two-layer neural networks with random initialization have been considered in Mei et al. (2018, 2019), which establish a LLN on the empirical measure of the weights at fixed times - we consider in this paper the trajectory convergence, i.e. the whole empirical measure process (time indexed) converges uniformly w.r.t. Skorohod topology. It enables not only to use the limiting PDE, for example to study the convergence of the weights towards the infimum of the loss function (see Chizat and Bach (2018) for preliminary results), but is also crucial to establish the central limit theorem, see for example Descours et al. (2022). Rotskoff and Vanden-Eijnden (2018) give conditions for global convergence of GD for exact mean-square loss and online stochastic gradient descent (SGD) with mini-batches increasing in size with the number of weights  $N$ . A LLN for the entire trajectory of the empirical measure is also given in Sirignano and Spiliopoulos (2020) for a standard SGD. De Bortoli et al. (2020) establish the propagation of chaos for SGD with different step size schemes. Compared to the existing literature dealing with the SGD empirical risk minimization in two-layer neural networks, Descours et al. (2022) provide the first rigorous proof of the existence of the limit PDE, and in particular its uniqueness, in the LLN.

We are interested here in deriving a LLN but for Variational Inference (VI) of two-layer Bayesian Neural Networks (BNN), where we consider a regularized version of the Evidence Lower Bound (ELBO).

## 2. Variational inference in BNN: Notations and common SGD schemes

### 2.1. Variational inference and Evidence Lower Bound

**Setting.** Let  $X$  and  $Y$  be subsets of  $\mathbf{R}^n$  ( $n \geq 1$ ) and  $\mathbf{R}$  respectively. For  $N \geq 1$  and  $\mathbf{w} = (w^1, \dots, w^N) \in (\mathbf{R}^d)^N$ , let  $f_{\mathbf{w}}^N : X \rightarrow \mathbf{R}$  be the following two-layer neural network: for  $x \in X$ ,

$$f_{\mathbf{w}}^N(x) := \frac{1}{N} \sum_{i=1}^N s(w^i, x) \in \mathbf{R},$$

where  $s : \mathbf{R}^d \times X \rightarrow \mathbf{R}$  is the activation function. We work in a Bayesian setting, in which we seek a distribution of the latent variable  $\mathbf{w}$  which represents the weights of the neural network. The standard problem in Bayesian inference over complex models is that the posterior distribution is hard to sample. To tackle this problem, we consider Variational Inference, in which we consider a family of distribution  $\mathcal{Q}^N = \{q_{\boldsymbol{\theta}}^N, \boldsymbol{\theta} \in \Xi^N\}$  (where  $\Xi$  is some parameter space) easy to sample. The

objective is to find the best  $q_{\theta}^N \in \mathcal{Q}^N$ , the one closest in KL divergence (denoted  $\mathcal{D}_{\text{KL}}$ ) to the exact posterior. Because we cannot compute the KL, we optimize the evidence lower bound (ELBO), which is equivalent to the KL up to an additive constant.

Denoting by  $\mathfrak{L} : \mathbf{R} \times \mathbf{R} \rightarrow \mathbf{R}_+$  the negative log-likelihood (by an abuse of language, we call this quantity the *loss*), the ELBO (see Blei et al. (2017)) is defined, for  $\theta \in \Xi^N$ ,  $(x, y) \in \mathbf{X} \times \mathbf{Y}$ , by

$$E_{\text{lbo}}(\theta, x, y) := - \int_{(\mathbf{R}^d)^N} \mathfrak{L}(y, f_{\mathbf{w}}^N(x)) q_{\theta}^N(\mathbf{w}) d\mathbf{w} - \mathcal{D}_{\text{KL}}(q_{\theta}^N | P_0^N),$$

where  $P_0^N$  is some prior on the weights of the NN. The ELBO is decomposed into two terms: one corresponding to the Kullback-Leibler (KL) divergence between the variational density and the prior and the other to a marginal likelihood term. It was empirically found that the maximization of the ELBO function is prone to yield very poor inferences (Coker et al., 2022). It is argued in Coker et al. (2022) and Huix et al. (2022) that optimizing the ELBO leads as  $N \rightarrow \infty$  to the collapse of the variational posterior to the prior. Huix et al. (2022) proposed to consider a regularized version of the ELBO, which consists in multiplying the KL term by a parameter which is scaled by the inverse of the number of neurons:

$$E_{\text{lbo}}^N(\theta, x, y) := - \int_{(\mathbf{R}^d)^N} \mathfrak{L}(y, f_{\mathbf{w}}^N(x)) q_{\theta}^N(\mathbf{w}) d\mathbf{w} - \frac{1}{N} \mathcal{D}_{\text{KL}}(q_{\theta}^N | P_0^N), \quad (1)$$

A first objective of this paper is to show that the proposed regularization leads to a stable asymptotic behavior and the effect of both the integrated loss and Kullback-Leibler terms on the limiting behavior are balanced in the limit  $N \rightarrow \infty$ . The maximization of  $E_{\text{lbo}}^N$  is carried out using SGD.

The variational family  $\mathcal{Q}^N$  we consider is a Gaussian family of distributions. More precisely, we assume that for any  $\theta = (\theta^1, \dots, \theta^N) \in \Xi^N$ , the variational distribution  $q_{\theta}^N$  factorizes over the neurons: for all  $\mathbf{w} = (w^1, \dots, w^N) \in (\mathbf{R}^d)^N$ ,  $q_{\theta}^N(\mathbf{w}) = \prod_{i=1}^N q_{\theta^i}^1(w^i)$ , where  $\theta = (m, \rho) \in \Xi := \mathbf{R}^d \times \mathbf{R}$  and  $q_{\theta}^1$  is the probability density function (pdf) of  $\mathcal{N}(m, g(\rho)^2 I_d)$ , with  $g(\rho) = \log(1 + e^{\rho})$ ,  $\rho \in \mathbf{R}$ .

In the following, we simply write  $\mathbf{R}^{d+1}$  for  $\mathbf{R}^d \times \mathbf{R}$ . In addition, following the reparameterisation trick of Blundell et al. (2015),  $q_{\theta}^1(w) dw$  is the pushforward of a reference probability measure with density  $\gamma$  by  $\Psi_{\theta}$  (see more precisely Assumption A1). In practice,  $\gamma$  is the pdf of  $\mathcal{N}(0, I_d)$  and  $\Psi_{\theta}(z) = m + g(\rho)z$ . With these notations, (1) writes

$$E_{\text{lbo}}^N(\theta, x, y) = - \int_{(\mathbf{R}^d)^N} \mathfrak{L}\left(y, \frac{1}{N} \sum_{i=1}^N s(\Psi_{\theta^i}(z^i), x)\right) \gamma(z^1) \dots \gamma(z^N) dz_1 \dots dz_N - \frac{1}{N} \mathcal{D}_{\text{KL}}(q_{\theta}^N | P_0^N).$$

**Loss function and prior distribution.** In this work, we focus on the regression problem, i.e.  $\mathfrak{L}$  is the Mean Square Loss: for  $y_1, y_2 \in \mathbf{R}$ ,  $\mathfrak{L}(y_1, y_2) = \frac{1}{2} |y_1 - y_2|^2$ . We also introduce the function  $\phi : (\theta, z, x) \in \mathbf{R}^{d+1} \times \mathbf{R}^d \times \mathbf{X} \mapsto s(\Psi_{\theta}(z), x)$ . On the other hand, we assume that the prior distribution  $P_0^N$  write, for all  $\mathbf{w} \in (\mathbf{R}^d)^N$ ,  $P_0^N(\mathbf{w}) = \prod_{i=1}^N P_0^1(w^i)$ , where  $P_0^1 : \mathbf{R}^d \rightarrow \mathbf{R}_+$  is the pdf of  $\mathcal{N}(m_0, \sigma_0^2 I_d)$ , and  $\sigma_0 > 0$ . Therefore  $\mathcal{D}_{\text{KL}}(q_{\theta}^N | P_0^N) = \sum_{i=1}^N \mathcal{D}_{\text{KL}}(q_{\theta^i}^1 | P_0^1)$  and, for  $\theta = (m, \rho) \in \mathbf{R}^{d+1}$ ,

$$\mathcal{D}_{\text{KL}}(q_{\theta}^1 | P_0^1) = \int_{\mathbf{R}^d} q_{\theta}^1(x) \log(q_{\theta}^1(x) / P_0^1(x)) dx = \frac{\|m - m_0\|_2^2}{2\sigma_0^2} + \frac{d}{2} \left( \frac{g(\rho)^2}{\sigma_0^2} - 1 \right) + \frac{d}{2} \log \left( \frac{\sigma_0^2}{g(\rho)^2} \right).$$

Note that  $\mathcal{D}_{\text{KL}}$  has at most a quadratic growth in  $m$  and  $\rho$ .

Note that we assume here a Gaussian prior to get an explicit expression of the Kullback-Leibler divergence. Most arguments extend to sufficiently regular densities and are essentially the same for exponential families, using conjugate families for the variational approximation.

## 2.2. Common SGD schemes in backpropagation in a variational setting

**Idealized SGD.** Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. Consider a data set  $\{(x_k, y_k)\}_{k \geq 0}$  i.i.d. w.r.t.  $\pi \in \mathcal{P}(X \times Y)$ , the space of probability measures over  $X \times Y$ . For  $N \geq 1$  and given a learning rate  $\eta > 0$ , the maximization of  $\theta \in \mathbf{R}^{d+1} \mapsto \mathbb{E}_{\text{lbo}}^N(\theta, x, y)$  with a SGD algorithm writes as follows: for  $k \geq 0$  and  $i \in \{1, \dots, N\}$ ,

$$\begin{cases} \theta_{k+1} = \theta_k + \eta \nabla_{\theta} \mathbb{E}_{\text{lbo}}^N(\theta_k, x_k, y_k) \\ \theta_0 \sim \mu_0^{\otimes N}, \end{cases} \quad (2)$$

where  $\mu_0 \in \mathcal{P}(\mathbf{R}^{d+1})$  and  $\theta_k = (\theta_k^1, \dots, \theta_k^N)$ . We now compute  $\nabla_{\theta} \mathbb{E}_{\text{lbo}}^N(\theta, x, y)$ .

First, under regularity assumptions on the function  $\phi$  (which will be formulated later, see **A1** and **A3** below) and by assumption on  $\mathfrak{L}$ , we have for all  $i \in \{1, \dots, N\}$  and all  $(x, y) \in X \times Y$ ,

$$\begin{aligned} & \int_{(\mathbf{R}^d)^N} \nabla_{\theta^i} \mathfrak{L}\left(y, \frac{1}{N} \sum_{j=1}^N \phi(\theta^j, z^j, x)\right) \gamma(z^1) \dots \gamma(z^N) dz^1 \dots dz^N \\ &= -\frac{1}{N^2} \sum_{j=1}^N \int_{(\mathbf{R}^d)^N} (y - \phi(\theta^j, z^j, x)) \nabla_{\theta} \phi(\theta^i, z^i, x) \gamma(z^1) \dots \gamma(z^N) dz^1 \dots dz^N \\ &= -\frac{1}{N^2} \left[ \sum_{j=1, j \neq i}^N (y - \langle \phi(\theta^j, \cdot, x), \gamma \rangle) \langle \nabla_{\theta} \phi(\theta^i, \cdot, x), \gamma \rangle + \langle (y - \phi(\theta^i, \cdot, x)) \nabla_{\theta} \phi(\theta^i, \cdot, x), \gamma \rangle \right], \end{aligned} \quad (3)$$

where we have used the notation  $\langle U, \nu \rangle = \int_{\mathbf{R}^q} U(z) \nu(dz)$  for any integrable function  $U : \mathbf{R}^q \rightarrow \mathbf{R}$  w.r.t. a measure  $\nu$  (with a slight abuse of notation, we denote by  $\gamma$  the measure  $\gamma(z) dz$ ). Second, for  $\theta \in \mathbf{R}^{d+1}$ , we have

$$\nabla_{\theta} \mathcal{D}_{\text{KL}}(q_{\theta}^1 | P_0^1) = \begin{pmatrix} \nabla_m \mathcal{D}_{\text{KL}}(q_{\theta}^1 | P_0^1) \\ \partial_{\rho} \mathcal{D}_{\text{KL}}(q_{\theta}^1 | P_0^1) \end{pmatrix} = \begin{pmatrix} \frac{1}{\sigma_0^2} (m - m_0) \\ \frac{d}{d\rho} g'(\rho) g(\rho) - d \frac{g'(\rho)}{g(\rho)} \end{pmatrix}. \quad (4)$$

In conclusion, the SGD (2) writes: for  $k \geq 0$  and  $i \in \{1, \dots, N\}$ ,

$$\begin{cases} \theta_{k+1}^i = \theta_k^i - \frac{\eta}{N^2} \sum_{j=1, j \neq i}^N \left( \langle \phi(\theta_k^j, \cdot, x_k), \gamma \rangle - y_k \right) \langle \nabla_{\theta} \phi(\theta_k^i, \cdot, x_k), \gamma \rangle \\ \quad - \frac{\eta}{N^2} \langle (\phi(\theta_k^i, \cdot, x_k) - y_k) \nabla_{\theta} \phi(\theta_k^i, \cdot, x_k), \gamma \rangle - \frac{\eta}{N} \nabla_{\theta} \mathcal{D}_{\text{KL}}(q_{\theta_k^i}^1 | P_0^1) \\ \theta_0^i \sim \mu_0. \end{cases} \quad (5)$$

We shall call this algorithm *idealised* SGD because it contains an intractable term given by the integral w.r.t.  $\gamma$ . This has motivated the development of methods where this integral is replaced by an

unbiased Monte Carlo estimator (see [Blundell et al. \(2015\)](#)) as detailed below.

**Bayes-by-Backprop SGD.** The second SGD algorithm we study is based on an approximation, for  $i \in \{1, \dots, N\}$ , of  $\int_{(\mathbf{R}^d)^N} (y - \phi(\theta^j, z^j, x)) \nabla_{\theta} \phi(\theta^i, z^i, x) \gamma(z^1) \dots \gamma(z^N) dz^1 \dots dz^N$  (see (3)) by

$$\frac{1}{B} \sum_{\ell=1}^B (y - \phi(\theta^j, \mathbf{Z}^{j,\ell}, x)) \nabla_{\theta} \phi(\theta^i, \mathbf{Z}^{i,\ell}, x) \quad (6)$$

where  $B \in \mathbf{N}^*$  is a fixed integer and  $(\mathbf{Z}^{q,\ell}, q \in \{i, j\}, 1 \leq \ell \leq B)$  is a i.i.d finite sequence of random variables distributed according to  $\gamma(z)dz$ . In this case, for  $N \geq 1$ , given a dataset  $(x_k, y_k)_{k \geq 0}$ , the maximization of  $\theta \in \mathbf{R}^{d+1} \mapsto \mathbb{E}_{\text{Ibo}}^N(\theta, x, y)$  with a SGD algorithm is the following: for  $k \geq 0$  and  $i \in \{1, \dots, N\}$ ,

$$\left\{ \begin{array}{l} \theta_{k+1}^i = \theta_k^i - \frac{\eta}{N^2 B} \sum_{j=1}^N \sum_{\ell=1}^B (\phi(\theta_k^j, \mathbf{Z}_k^{j,\ell}, x_k) - y_k) \nabla_{\theta} \phi(\theta_k^i, \mathbf{Z}_k^{i,\ell}, x_k) - \frac{\eta}{N} \nabla_{\theta} \mathcal{D}_{\text{KL}}(q_{\theta_k^i}^1 | P_0^1) \\ \theta_0^i = (m_0^i, \rho_0^i) \sim \mu_0, \end{array} \right. \quad (7)$$

where  $\eta > 0$  and  $(\mathbf{Z}_k^{j,\ell}, 1 \leq j \leq N, 1 \leq \ell \leq B, k \geq 0)$  is a i.i.d sequence of random variables distributed according to  $\gamma$ .

### 3. Law of large numbers for the idealized SGD

**Assumptions and notations.** When  $E$  is a metric space and  $\mathcal{I} = \mathbf{R}_+$  or  $\mathcal{I} = [0, T]$  ( $T \geq 0$ ), we denote by  $\mathcal{D}(\mathcal{I}, E)$  the Skorohod space of càdlàg functions on  $\mathcal{I}$  taking values in  $E$  and  $\mathcal{C}(\mathcal{I}, E)$  the space of continuous functions on  $\mathcal{I}$  taking values in  $E$ . The evolution of the parameters  $(\{\theta_k^i, i = 1, \dots, N\})_{k \geq 1}$  defined by (5) is tracked through their empirical distribution  $\nu_k^N$  (for  $k \geq 0$ ) and its scaled version  $\mu_t^N$  (for  $t \in \mathbf{R}_+$ ), which are defined as follows:

$$\nu_k^N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_k^i} \quad \text{and} \quad \mu_t^N := \nu_{[Nt]}^N, \quad \text{where the } \theta_k^i \text{'s are defined (5)}. \quad (8)$$

Fix  $T > 0$ . For all  $N \geq 1$ ,  $\mu^N := \{\mu_t^N, t \in [0, T]\}$  is a random element of  $\mathcal{D}([0, T], \mathcal{P}(\mathbf{R}^{d+1}))$ , where  $\mathcal{P}(\mathbf{R}^{d+1})$  is endowed with the weak convergence topology. For  $N \geq 1$  and  $k \geq 1$ , we introduce the following  $\sigma$ -algebras:

$$\mathcal{F}_0^N = \sigma(\theta_0^i, 1 \leq i \leq N) \quad \text{and} \quad \mathcal{F}_k^N = \sigma(\theta_0^i, (x_q, y_q), 1 \leq i \leq N, 0 \leq q \leq k-1). \quad (9)$$

Recall  $q_{\theta}^1 : \mathbf{R}^d \rightarrow \mathbf{R}_+$  be the pdf of  $\mathcal{N}(m, g(\rho)^2 I_d)$  ( $\theta = (m, \rho) \in \mathbf{R}^{d+1}$ ). In this work, we assume the following.

- A1.** There exists a pdf  $\gamma : \mathbf{R}^d \rightarrow \mathbf{R}_+$  such that for all  $\theta \in \mathbf{R}^{d+1}$ ,  $q_{\theta}^1 dx = \Psi_{\theta} \# \gamma dx$ , where  $\{\Psi_{\theta}, \theta \in \mathbf{R}^{d+1}\}$  is a family of  $\mathcal{C}^1$ -diffeomorphisms over  $\mathbf{R}^d$  such that for all  $z \in \mathbf{R}^d$ ,  $\theta \in \mathbf{R}^{d+1} \mapsto \Psi_{\theta}(z)$  is of class  $\mathcal{C}^{\infty}$ . Finally, there exists  $\mathfrak{b} : \mathbf{R}^d \rightarrow \mathbf{R}_+$  such that for all multi-index  $\alpha \in \mathbf{N}^{d+1}$  with  $|\alpha| \geq 1$ , there exists  $C_{\alpha} > 0$ , for all  $z \in \mathbf{R}^d$  and  $\theta = (\theta_1, \dots, \theta_{d+1}) \in \mathbf{R}^{d+1}$ ,

$$|\partial_{\alpha} \Psi_{\theta}(z)| \leq C_{\alpha} \mathfrak{b}(z) \quad \text{with for all } q \geq 1, \langle \mathfrak{b}^q, \gamma \rangle < +\infty, \quad (10)$$

where  $\partial_{\alpha} = \partial_{\theta_1}^{\alpha_1} \dots \partial_{\theta_{d+1}}^{\alpha_{d+1}}$  and  $\partial_{\theta_j}^{\alpha_j}$  is the partial derivatives of order  $\alpha_j$  w.r.t. to  $\theta_j$ .



- A2.** The sequence  $\{(x_k, y_k)\}_{k \geq 0}$  is i.i.d. w.r.t.  $\pi \in \mathcal{P}(X \times Y)$ . The set  $X \times Y \subset \mathbf{R}^d \times \mathbf{R}$  is compact. For all  $k \geq 0$ ,  $(x_k, y_k) \perp\!\!\!\perp \mathcal{F}_k^N$ , where  $\mathcal{F}_k^N$  is defined in (9).
- A3.** The activation function  $s : \mathbf{R}^d \times X \rightarrow \mathbf{R}$  belongs to  $\mathcal{C}_b^\infty(\mathbf{R}^d \times X)$  (the space of smooth functions over  $\mathbf{R}^d \times X$  whose derivatives of all order are bounded).
- A4.** The initial parameters  $(\theta_0^i)_{i=1}^N$  are i.i.d. w.r.t.  $\mu_0 \in \mathcal{P}(\mathbf{R}^{d+1})$  which has compact support.

Note that **A1** is satisfied when  $\gamma$  is the pdf of  $\mathcal{N}(0, I_d)$  and  $\Psi_\theta(z) = m + g(\rho)z$ , with  $\mathfrak{b}(z) = 1 + |z|$ . With these assumptions, for every fixed  $T > 0$ , the sequence  $(\{\theta_k^i, i = 1, \dots, N\})_{k=0, \dots, \lfloor NT \rfloor}$  defined by (5) is a.s. bounded:

**Lemma 1 (Uniform bound on the parameters)** *Assume **A1**→**A4**. Then, there exists  $C > 0$  such that a.s. for all  $T > 0$ ,  $N \geq 1$ ,  $i \in \{1, \dots, N\}$ , and  $0 \leq k \leq \lfloor NT \rfloor$ ,  $|\theta_k^i| \leq Ce^{[C(2+T)]T}$ .*

Lemma 1 implies that a.s. for all  $T > 0$  and  $N \geq 1$ ,  $\mu^N \in \mathcal{D}([0, T], \mathcal{P}(\Theta_T))$ , where

$$\Theta_T = \{\theta \in \mathbf{R}^{d+1}, |\theta| \leq Ce^{[C(2+T)]T}\}.$$

**Law of large numbers for  $(\mu^N)_{N \geq 1}$  defined in (8).** The first main result of this work is the following.

**Theorem 2** *Assume **A1**→**A4**. Let  $T > 0$ . Then, the sequence  $(\mu^N)_{N \geq 1} \subset \mathcal{D}([0, T], \mathcal{P}(\Theta_T))$  defined in (8) converges in probability to the unique deterministic solution  $\bar{\mu} \in \mathcal{C}([0, T], \mathcal{P}(\Theta_T))$  to the following measure-valued evolution equation:  $\forall f \in \mathcal{C}^\infty(\Theta_T)$  and  $\forall t \in [0, T]$ ,*

$$\begin{aligned} \langle f, \bar{\mu}_t \rangle - \langle f, \mu_0 \rangle &= -\eta \int_0^t \int_{X \times Y} \langle \phi(\cdot, \cdot, x) - y, \bar{\mu}_s \otimes \gamma \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \bar{\mu}_s \otimes \gamma \rangle \pi(dx, dy) ds \\ &\quad - \eta \int_0^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \bar{\mu}_s \rangle ds. \end{aligned} \quad (11)$$

The proof of Theorem 2 is given in Appendix A. We stress here the most important steps and used techniques. In a first step, we derive an identity satisfied by  $(\mu^N)_{N \geq 1}$ , namely the pre-limit equation (28); see Sec. A.1. Then we show in Sec. A.2.2 that  $(\mu^N)_{N \geq 1}$  is relatively compact in  $\mathcal{D}([0, T], \mathcal{P}(\Theta_T))$ . To do so, we check that the sequence  $(\mu^N)_{N \geq 1}$  satisfies all the required assumptions of (Jakubowski, 1986, Theorem 3.1) when  $E = \mathcal{P}(\Theta_T)$  there. In Sec. A.2.3 we prove that every limit point of  $(\mu^N)_{N \geq 1}$  satisfies the limit equation (11). Then, in Section A.2.4, we prove that there is a unique solution of the measure-valued equation (11). To prove the uniqueness of the solution of (11), we use techniques developed in Piccoli et al. (2015) which are based on a representation formula for solution to measure-valued equations (Villani, 2003, Theorem 5.34) together with estimates in Wasserstein distances between two solutions of (11) derived in Piccoli and Rossi (2016). In Section A.2.4, we also conclude the proof of Theorem 2. Compared to (Descours et al., 2022, Theorem 1), the fact that  $(\{\theta_k^i, i = 1, \dots, N\})_{k=0, \dots, \lfloor NT \rfloor}$  defined by (5) are a.s. bounded allows to use different and more straightforward arguments to prove (i) the relative compactness in  $\mathcal{D}([0, T], \mathcal{P}(\Theta_T))$  of  $(\mu^N)_{N \geq 1}$  (defined in (8)) (ii) the continuity property of the operator  $m \mapsto \Lambda_t[f](m)$  defined in (35) w.r.t. the topology of  $\mathcal{D}([0, T], \mathcal{P}(\Theta_T))$  and (iii)  $(\mu^N)_{N \geq 1}$  has limit points in  $\mathcal{C}([0, T], \mathcal{P}(\Theta_T))$ . Step (ii) is necessary in order to pass to the limit  $N \rightarrow +\infty$

in the pre-limit equation and Step (iii) is crucial since we prove that there is at most one solution of (11) in  $\mathcal{C}([0, T], \mathcal{P}(\Theta_T))$ . It is worthwhile to emphasize that, as  $N \rightarrow \infty$ , the effects of the integrated loss and of the KL terms are balanced, as conjectured in [Huix et al. \(2022\)](#).

To avoid further technicalities, we have chosen what may seem restrictive assumptions on the data or the activation function. Note however that it readily extends to unbounded set  $X$ , and also unbounded  $Y$  assuming that  $\pi$  as polynomial moments of sufficiently high order. Also, RELU (or more easily leaky RELU) may be considered by using weak derivatives (to consider the singularity at 0), and a priori moment bounds on the weights.

#### 4. LLN for the *Bayes-by-Backprop* SGD

The sequence  $\{\theta_k^i, i \in \{1, \dots, N\}\}_{k=0, \dots, [NT]}$  defined recursively by the algorithm (7) is in general not bounded, since  $\nabla_\theta \phi(\theta, Z, x)$  is not necessarily bounded if  $Z \sim \gamma(s)dz$ . Therefore, we cannot expect Lemma 1 to hold for  $\{\theta_k^i, i \in \{1, \dots, N\}\}_{k=0, \dots, [NT]}$  set by (7). Thus, the sequence  $\{\theta_k^i, i \in \{1, \dots, N\}\}_{k=0, \dots, [NT]}$  is considered on the whole space  $\mathbf{R}^{d+1}$ .

**Wasserstein spaces and results.** For  $N \geq 1$ , and  $k \geq 1$ , we set

$$\mathcal{F}_k^N = \sigma\left(\theta_0^i, Z_q^{j,\ell}, (x_q, y_q), 1 \leq i, j \leq N, 1 \leq \ell \leq B, 0 \leq q \leq k-1\right). \quad (12)$$

In addition to **A1**→**A4** (where in **A2**, when  $k \geq 1$ ,  $\mathcal{F}_k^N$  is now the one defined in (12)), we assume:

**A5.** The sequences  $(Z_k^{j,\ell}, 1 \leq j \leq N, 1 \leq \ell \leq B, k \geq 0)$  and  $((x_k, y_k), k \geq 0)$  are independent.

In addition, for  $k \geq 0$ ,  $((x_k, y_k), Z_k^{j,\ell}, 1 \leq j \leq N, 1 \leq \ell \leq B) \perp\!\!\!\perp \mathcal{F}_k^N$ .

Note that the last statement of **A5** implies the last statement of **A2**. We introduce the scaled empirical distribution of the parameters of the algorithm (7), i.e. for  $k \geq 0$  and  $t \geq 0$ :

$$\nu_k^N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_k^i} \quad \text{and} \quad \mu_t^N := \nu_{[Nt]}^N, \quad \text{where the } \theta_k^i \text{'s are defined (7)}. \quad (13)$$

One can no longer rely on the existence of a compact subset  $\Theta_T \subset \mathbf{R}^{d+1}$  such that a.s.  $(\mu^N)_{N \geq 1} \subset \mathcal{D}([0, T], \mathcal{P}(\Theta_T))$ , where  $\mu^N = \{t \geq 0 \mapsto \mu_t^N\}$  is defined in (13). For this reason, we will work in Wasserstein spaces  $\mathcal{P}_q(\mathbf{R}^{d+1})$ ,  $q \geq 0$ , which, we recall, are defined by

$$\mathcal{P}_q(\mathbf{R}^{d+1}) = \left\{ \nu \in \mathcal{P}(\mathbf{R}^{d+1}), \int_{\mathbf{R}^{d+1}} |\theta|^q \nu(d\theta) < +\infty \right\}. \quad (14)$$

These spaces are endowed with the Wasserstein metric  $W_q$ , see e.g. [\(Santambrogio, 2015, Chapter 5\)](#) for more materials on Wasserstein spaces. For all  $q \geq 0$ ,  $(\mu^N)_{N \geq 1} \subset \mathcal{D}(\mathbf{R}_+, \mathcal{P}_q(\mathbf{R}^{d+1}))$ . The second main results of this work is a LLN for  $(\mu^N)_{N \geq 1}$  defined in (13).

**Theorem 3** *Assume **A1**→**A5**. Let  $\gamma_0 > 1 + \frac{d+1}{2}$ . Then, the sequence  $(\mu^N)_{N \geq 1}$  defined in (13) converges in probability in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$  to a deterministic element  $\bar{\mu} \in \mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ , where  $\bar{\mu} \in \mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$  is the unique solution in  $\mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$  to the following measure-valued evolution equation:  $\forall f \in \mathcal{C}_b^\infty(\mathbf{R}^{d+1})$  and  $\forall t \in \mathbf{R}_+$ ,*

$$\begin{aligned} \langle f, \bar{\mu}_t \rangle - \langle f, \mu_0 \rangle &= -\eta \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \bar{\mu}_s \otimes \gamma \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \bar{\mu}_s \otimes \gamma \rangle \pi(dx, dy) ds \\ &\quad - \eta \int_0^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \bar{\mu}_s \rangle ds. \end{aligned} \quad (15)$$

Theorem 3 is proved in the appendix B. Since  $\{\theta_k^i, i \in \{1, \dots, N\}\}_{k=0, \dots, \lfloor NT \rfloor}$  defined by (7) is not bounded in general, we work in the space  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ . The proof of Theorem 3 is more involved than that of Theorem 2, and generalizes the latter to the case where the parameters of the SGD algorithm are unbounded. We prove that  $(\mu^N)_{N \geq 1}$  (defined in (13)) is relatively compact in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ . To this end we now use (Jakubowski, 1986, Theorem 4.6). The compact containment, which is the purpose of Lemma 20, is not straightforward since  $\mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1})$  is not compact contrary to Theorem 2 where we used the compactness of  $\mathcal{P}(\Theta_T)$ . More precisely, the compact containment of  $(\mu^N)_{N \geq 1}$  relies on a characterization of the compact subsets of  $\mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1})$  (see Proposition 18) and moment estimates on  $\{\theta_k^i, i \in \{1, \dots, N\}\}_{k=0, \dots, \lfloor NT \rfloor}$  (see Lemma 17). We also mention that contrary to what is done in the proof of Theorem 2, we do not show that every limit point of  $(\mu^N)_{N \geq 1}$  in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$  is continuous in time but we still manage to prove that they all satisfy (15). Then, using the duality formula for the  $W_1$ -distance together with rough estimates on the jumps of  $t \mapsto \langle f, \mu_t^N \rangle$  (for  $f$  uniformly Lipschitz over  $\mathbf{R}^{d+1}$ ), we then show that every limit point of  $(\mu^N)_{N \geq 1}$  in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$  belongs a.s. to  $\mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$ . Again this is important since we have uniqueness of (15) in  $\mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$ .

We conclude this section with the following important uniqueness result.

**Proposition 4** *Under the assumptions of Theorems 2 and 3, the solution to (11) is independent of  $T$  and is equal to the solution to (15).*

This uniqueness result states that both idealized and *Bayes-by-backprop* SGD have the same limiting behavior. It is also noteworthy that the mini-batch  $B$  is held fixed  $B$ . The effect of batch size can be seen at the level of the central limit theorem, which we leave for future work.

## 5. The Minimal-VI SGD algorithm

The idea behind the *Bayes-by-Backprop* SGD stems from the fact that there are integrals wrt  $\gamma$  in the loss function that cannot be computed in practice and it is quite natural up to a reparameterization trick, to replace these integrals by a Monte Carlo approximation (with i.i.d. gaussian random variables). To devise a new cheaper algorithm based on the only terms impacting the asymptotic limit, we directly analyse the limit equation (11) and remark that it can be rewritten as,  $\forall f \in \mathcal{C}^\infty(\Theta_T)$  and  $\forall t \in [0, T]$ ,

$$\begin{aligned} & \langle f, \bar{\mu}_t \rangle - \langle f, \mu_0 \rangle \\ &= -\eta \int_0^t \int_{\mathbf{X} \times \mathbf{Y} \times (\mathbf{R}^d)^2} \langle \phi(\cdot, z_1, x) - y, \bar{\mu}_s \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, z_2, x), \bar{\mu}_s \rangle \gamma^{\otimes 2}(dz_1 dz_2) \pi(dx, dy) ds \\ & \quad - \eta \int_0^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \bar{\mu}_s \rangle ds. \end{aligned}$$

Thus, the integration over  $\gamma^{\otimes 2}$  can be considered as that over  $\pi$ , i.e., we can consider them as two more data variables that only need to be sampled at each new step. In this case, the SGD (7) becomes: for  $k \geq 0$  and  $i \in \{1, \dots, N\}$ ,

$$\begin{cases} \theta_{k+1}^i = \theta_k^i - \frac{\eta}{N^2} \sum_{j=1}^N (\phi(\theta_k^j, Z_k^1, x_k) - y_k) \nabla_\theta \phi(\theta_k^i, Z_k^2, x_k) - \frac{\eta}{N} \nabla_\theta \mathcal{D}_{\text{KL}}(q_{\theta_k^i}^1 | P_0^1) \\ \theta_0^i = (m_0^i, \rho_0^i) \sim \mu_0, \end{cases} \quad (16)$$

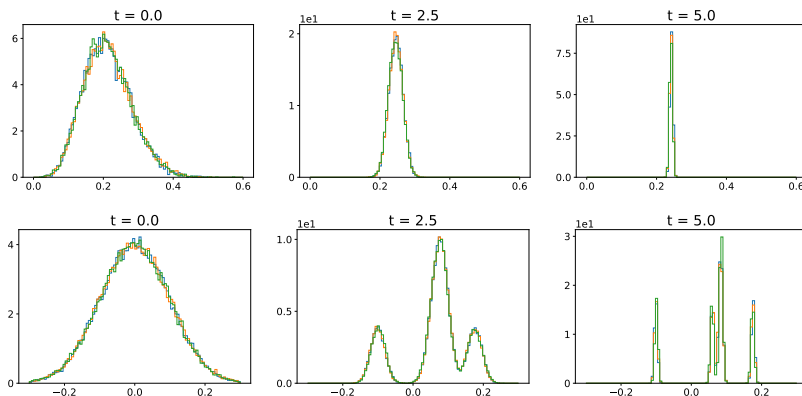


Figure 1: Histograms of  $\{F(\theta_{[NT]}^i), i = 1, \dots, N\}$ , at different times (initialization ( $t = 0$ ), half ( $t = 2.5$ ) and end of training ( $T = 5$ )), when  $N = 10000$ . First line:  $F(\theta) = \|m\|_2$ , where  $\theta = (m, \rho) \in \mathbf{R}^d \times \mathbf{R}$ . Second line:  $F(\theta) = m \in \mathbf{R}^d$ . Idealized (blue), *Bayes-by-Backprop* (orange) and *Minimal-VI* (green).

where  $\eta > 0$  and  $(Z_k^p, p \in \{1, 2\}, k \geq 0)$  is a i.i.d sequence of random variables distributed according to  $\gamma^{\otimes 2}$ . We call this backpropagation scheme *minimal-VI SGD* which is much cheaper in terms of computational complexity, with the same limiting behavior as we now discuss.

We introduce the  $\sigma$ -algebra for  $N, k \geq 1$ :

$$\mathcal{F}_k^N = \sigma\left(\theta_0^i, Z_q^p, (x_q, y_q), 1 \leq i \leq N, p \in \{1, 2\}, 0 \leq q \leq k-1\right). \quad (17)$$

In addition to **A1**→**A4** (where in **A2**,  $\mathcal{F}_k^N$  is now the one defined above in (17) when  $k \geq 1$ ), the following assumption

**A6.** The sequences  $(Z_k^p, p \in \{1, 2\}, k \geq 0)$  and  $((x_k, y_k), k \geq 0)$  are independent. In addition, for  $k \geq 0$ ,  $((x_k, y_k), Z_k^p, p \in \{1, 2\}) \perp\!\!\!\perp \mathcal{F}_k^N$ , where  $\mathcal{F}_k^N$  is defined in (17).

Set for  $k \geq 0$  and  $t \geq 0$ ,  $\nu_k^N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_k^i}$  and  $\mu_t^N := \nu_{[Nt]}^N$ , where the  $\theta_k^i$ 's are defined in (16). The last main result of this work states that the sequence  $(\mu^N)_{N \geq 1}$  satisfies the same law of large numbers when  $N \rightarrow +\infty$  as the one satisfied by (13), whose proof will be omitted as it is the same as the one made for Theorem 3.

**Theorem 5** *Assume **A1**→**A4** and **A6**. Then, the sequence of  $(\mu^N)_{N \geq 1}$  satisfies all the statements of Theorem 3.*

## 6. Numerical experiments

In this section we illustrate the theorems 2, 3, and 5 using the following toy model. We set  $d = 5$ . Given  $\theta^* \in \mathbf{R}^d$  (drawn from a normal distribution and scaled to the unit norm), we draw i.i.d observations as follows: Given  $x \sim \mathcal{U}([-1, 1]^d)$ , we draw  $y = \tanh(x^\top \theta^*) + \epsilon$ , where  $\epsilon$  is

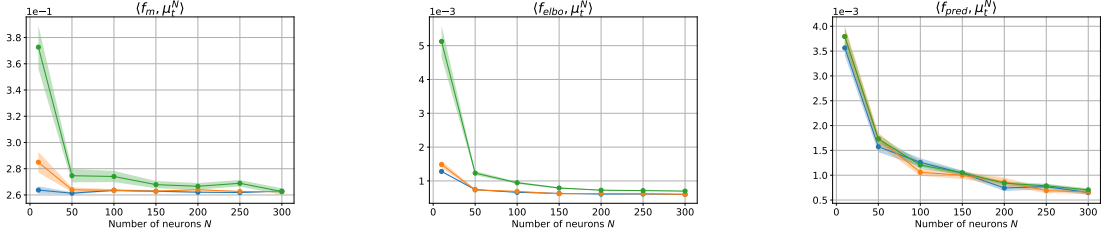


Figure 2: Convergence of  $\langle f, \mu_T^N \rangle$  to  $\langle f, \bar{\mu}_T \rangle$ , for the idealized (blue), *Bayes-by-Backprop* (orange) and *Minimal-VI* (green) SGD algorithms over 50 realizations.

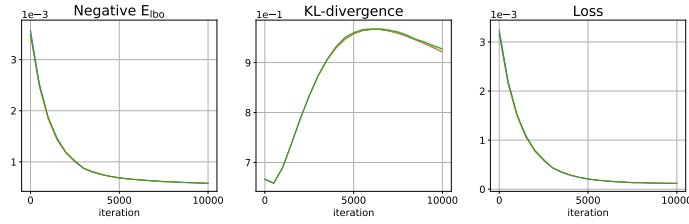


Figure 3: Decay of the negative ELBO (left) and its two components (KL (middle), loss (right)) during the training process done by the idealized (blue), *Bayes-by-Backprop* (orange) and *Minimal-VI* (green) SGD algorithms, for  $N = 10000$ .

zero mean with variance  $10^{-4}$ . The initial distribution of parameters is centered around the prior:  $\theta_0 \sim (\mathcal{N}(m_0, 0.01I_d) \times \mathcal{N}(g^{-1}(\sigma_0), 0.01))^{\otimes N}$ , with  $m_0 = 0$  and  $\sigma_0 = 0.2$ . Since the idealized algorithm cannot be implemented exactly, a mini-batch of size 100 is used as a proxy for the following comparisons of the different algorithms. For the algorithm (7) SGD we set  $B = 1$ .

**Evolution and limit of the distribution** Fig. 1 displays the histograms of  $\{F(\theta_{[Nt]}^i), i = 1, \dots, N\}$  ( $F(\theta) = \|m\|_2, g(\rho)$  or  $m$ , where  $\theta = (m, \rho) \in \mathbf{R}^d \times \mathbf{R}$ ), for  $N = 10000$ , at initialization, halfway through training, and at the end of training. The empirical distributions illustrated by these histograms are very similar over the course of training. It can be seen that for  $N = 10000$  the limit of the mean field is reached.

**Convergence with respect to the numbers of neurons.** We investigate here the speed of convergence of  $\mu_t^N$  to  $\bar{\mu}_t$  (as  $N \rightarrow +\infty$ ), when tested against test functions  $f$ . More precisely, we fix a time  $T$  (end of training) and Figure 2 represents the empirical mean of  $\langle f, \mu_T^N \rangle$  over 50 realizations. The test functions used for this experiment are  $f_m(\theta) = \|m\|_2$ ,  $f_{\text{Elbo}}(\theta) = -\hat{E}_{\text{Elbo}}(\theta)^N$  where  $\hat{E}_{\text{Elbo}}$  is the empirical  $E_{\text{Elbo}}^N$  (see (1)) computed with 100 samples of  $(x, y)$  and  $(z^1, \dots, z^N)$ . Finally,  $f_{\text{pred}}(\theta) = \hat{\mathbb{E}}_x \left[ \hat{\mathbb{V}}_{w \sim q_{\theta}^N} [f_w^N(x)]^{1/2} \right]$  where  $\hat{\mathbb{E}}$  and  $\hat{\mathbb{V}}$  denote respectively the empirical mean and the empirical variance over 100 samples. All algorithms are converging to the same limit and are performing similarly even with a limited number of neurons ( $N = 300$  in this example).

**Convergence with respect to time.** This section illustrates the training process of a BNN with a given number of neurons  $N = 10000$ . In Figure 3, we plot the negative ELBO on a test set and its two components, the loss and the KL-divergence terms. Figure 3 shows that the BNN is able to learn

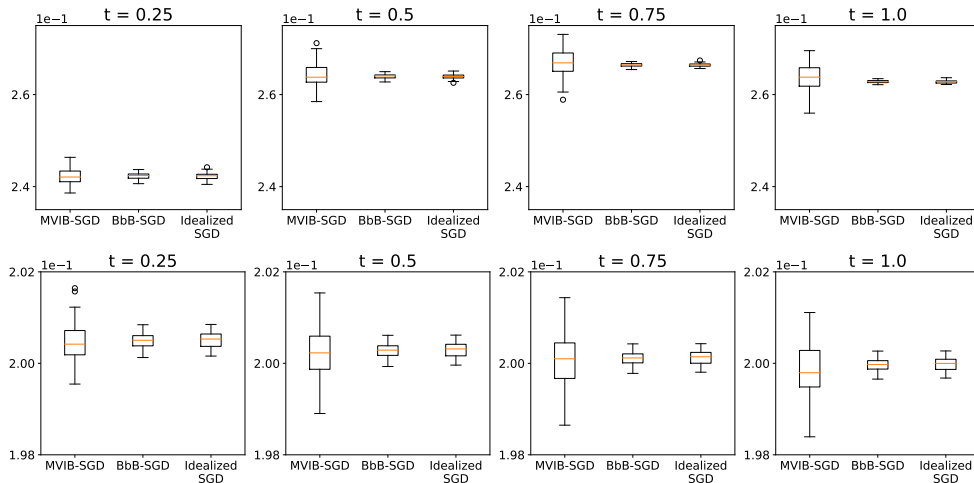


Figure 4: Boxplots for 50 runs of  $\langle f, \mu_t^N \rangle$  for the three SGD schemes for  $f(\theta) = \|m\|_2$  on the first line and  $f(\theta) = g(\rho)$  on the second line. MVIB-SGD: *Minimal-VI* SGD. BbB-SGD: *Bayes-by-Backprop* SGD.

on this specific task and all algorithms exhibit a similar performance. It illustrates the trajectorial convergence of  $\{\mu_t^N, t \in [0, T]\}_{N \geq 1}$  to  $\{\bar{\mu}_t, t \in [0, T]\}$  as  $N \rightarrow +\infty$ .

**Behavior around the limit  $\bar{\mu}$ .** On Figure 4, we plot the boxplots of  $\langle f, \mu_t^N \rangle$  for 50 realizations and  $N = 10000$ , at different times of the training. *Minimal-VI* scheme (which is computationally cheaper as explained in 5) exhibit a larger variance than the other algorithms.

## 7. Conclusion

By establishing the limit behavior of the idealized SGD for the variational inference of BNN with the weighting suggested by Huix et al. (2022), we have rigorously shown that the most-commonly used in practice *Bayes-by-Backprop* scheme indeed exhibits the same limit behavior. Furthermore, the analysis of the limit equation led us to validate the correct scaling of the KL divergence term in with respect to the loss. Notably, the mean-field limit dynamics has also helped us to devise a far less costly new SGD algorithm, the *Minimal-VI*. This scheme shares the same limit behavior, but only stems from the non-vanishing asymptotic contributions, hence the reduction of the computational cost. Aside from confirming the analytical results, the first simulations presented here show that the three algorithms, while having the same limit, may differ in terms of variance. Thus, deriving a CLT result and discussing the right trade-off between computational complexity and variance will be done in future work. Also, on a more general level regarding uncertainty quantification, an interesting question is to analyse the impact of the correct scaling of the KL divergence term on the error calibration and how to apply the same analysis in the context of deep ensembles.

## Acknowledgments

A.D. is grateful for the support received from the Agence Nationale de la Recherche (ANR) of the French government through the program "Investissements d’Avenir" (16-IDEX-0001 CAP 20-25) A.G. is supported by the French ANR under the grant ANR-17-CE40-0030 (project *EFI*) and the Institut Universitaire de France. M.M. acknowledges the support of the the French ANR under the grant ANR-20-CE46-0007 (*SuSa* project). B.N. is supported by the grant IA20Nectoux from the Projet I-SITE Clermont CAP 20-25. E.M. and T.H. acknowledge the support of ANR-CHIA-002, "Statistics, computation and Artificial Intelligence"; Part of the work has been developed under the auspice of the Lagrange Center for Mathematics and Calculus

## References

- Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76:243–297, 2021.
- P. Billingsley. *Convergence of Probability Measures*. John Wiley & Sons, 2nd edition, 1999.
- D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017. doi: 10.1080/01621459.2017.1285773. URL <https://doi.org/10.1080/01621459.2017.1285773>.
- C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra. Weight uncertainty in neural network. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1613–1622, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/blundell15.html>.
- L. Chizat. Mean-field langevin dynamics: Exponential convergence and annealing, 2022. URL <https://arxiv.org/abs/2202.01009>.
- L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/a1afc58c6ca9540d057299ec3016d726-Paper.pdf>.
- L. Chizat, M. Colombo, X. Fernandez-Real, and A. Figalli. Infinite-width limit of deep linear neural networks, 2022. URL <https://arxiv.org/abs/2211.16980>.
- A. D. Cobb and B. Jalaian. Scaling hamiltonian monte carlo inference for bayesian neural networks with symmetric splitting. In Cassio de Campos and Marloes H. Maathuis, editors, *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, volume 161 of *Proceedings of Machine Learning Research*, pages 675–685. PMLR, 27–30 Jul 2021. URL <https://proceedings.mlr.press/v161/cobb21a.html>.

- Beau Coker, Wessel P. Bruinsma, David R. Burt, Weiwei Pan, and Finale Doshi-Velez. Wide mean-field bayesian neural networks ignore the data. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 5276–5333. PMLR, 2022.
- V. De Bortoli, A. Durmus, X. Fontaine, and U. Simsekli. Quantitative propagation of chaos for SGD in wide neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 278–288. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/02e74f10e0327ad868d138f2b4fdd6f0-Paper.pdf>.
- P. Del Moral and A. Guionnet. Central limit theorem for nonlinear filtering and interacting particle systems. *The Annals of Applied Probability*, 9(2):275–297, 1999.
- F. Delarue, D. Lacker, and K. Ramanan. From the master equation to mean field game limit theory: a central limit theorem. *Electronic Journal of Probability*, 24:1–54, 2019.
- A. Descours, A. Guillin, M. Michel, and B. Nectoux. Law of large numbers and central limit theorem for wide two-layer neural networks: the mini-batch and noisy case. *arXiv preprint arXiv:2207.12734*, 2022.
- S. Ethier and T. Kurtz. *Markov Processes: Characterization and Convergence*, volume 282. John Wiley & Sons, 2009.
- B. Fernandez and S. Méléard. A Hilbertian approach for fluctuations on the McKean-Vlasov model. *Stochastic Processes and their Applications*, 71(1):33–53, 1997.
- A. Filos, S. Farquhar, A. N. Gomez, T. Rudner, Z. Kenton, L. Smith, M. Alizadeh, A. De Kroon, and Y. Gal. A systematic comparison of bayesian deep learning robustness in diabetic retinopathy tasks. *arXiv preprint arXiv:1912.10481*, 2019.
- Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/gal16.html>.
- Z. Ghahramani. Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553):452–459, 2015.
- A. Graves. Practical variational inference for neural networks. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL <https://proceedings.neurips.cc/paper/2011/file/7eb3c8be3d411e8ebfab08eba5f49632-Paper.pdf>.
- Geoffrey Hinton and Drew Van Camp. Keeping neural networks simple by minimizing the description length of the weights. In *in Proc. of the 6th Ann. ACM Conf. on Computational Learning Theory*, pages 5–13. ACM Press, 1993.



- M. Hitsuda and I. Mitoma. Tightness problem and stochastic evolution equation arising from fluctuation phenomena for interacting diffusions. *Journal of Multivariate Analysis*, 19(2):311–328, 1986.
- T. Huix, S. Majewski, A. Durmus, E. Moulines, and A. Korba. Variational inference of overparameterized bayesian neural networks: a theoretical and empirical study, 2022. URL <https://arxiv.org/abs/2207.03859>.
- P. Izmailov, S. Vikram, M. D. Hoffman, and A. G. G. Wilson. What are bayesian neural network posteriors really like? In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 4629–4640. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/izmailov21a.html>.
- J. Jacod and A. Shiryaev. *Skorokhod Topology and Convergence of Processes*. Springer, 1987.
- A. Jakubowski. On the skorokhod topology. In *Annales de l’IHP Probabilités et statistiques*, volume 22, pages 263–285, 1986.
- B. Jourdain and S. Méléard. Propagation of chaos and fluctuations for a moderate model with smooth initial data. *Annales de l’Institut Henri Poincaré (B) Probability and Statistics*, 34(6):727–766, 1998.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *arXiv preprint arXiv:1703.04977*, 2017.
- M. Khan, D. Nielsen, V. Tangkaratt, W. Lin, Y. Gal, and A. Srivastava. Fast and scalable Bayesian deep learning by weight-perturbation in Adam. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2611–2620. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/khan18a.html>.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *Proceedings of the 2nd International Conference on Learning Representations*, 2014.
- M. Krzywinski and N. Altman. Importance of being uncertain. *Nature methods*, 10(9):809–811, 2013.
- T. Kurtz and J. Xiong. A stochastic evolution equation arising from the fluctuations of a class of interacting particle systems. *Communications in Mathematical Sciences*, 2(3):325–358, 2004.
- C. Louizos and M. Welling. Multiplicative normalizing flows for variational Bayesian neural networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 2218–2227. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/louizos17a.html>.
- David JC MacKay. Probable networks and plausible predictions—a review of practical bayesian methods for supervised neural networks. *Network: computation in neural systems*, 6(3):469, 1995.

- David JC MacKay et al. Ensemble learning and evidence maximization. In *Proc. Nips*, volume 10, page 4083. Citeseer, 1995.
- R. McAllister, Y. Gal, A. Kendall, M. van der Wilk, A. Shah, R. Cipolla, and A. Weller. Concrete problems for autonomous vehicle safety: Advantages of bayesian deep learning. In *IJCAI*, 2017.
- S. Mei, A. Montanari, and P-M. Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- S. Mei, T. Misiakiewicz, and A. Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*, pages 2388–2464. PMLR, 2019.
- R. Micheltore, M. Wicker, L. Laurenti, L. Cardelli, Y. Gal, and M. Kwiatkowska. Uncertainty quantification with statistical guarantees in end-to-end autonomous driving control. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7344–7350, 2020. doi: 10.1109/ICRA40945.2020.9196844.
- V.M. Panaretos and Y. Zemel. *An Invitation to Statistics in Wasserstein Space*. Springer Nature, 2020.
- B. Piccoli and F. Rossi. On properties of the generalized Wasserstein distance. *Archive for Rational Mechanics and Analysis*, 222(3):1339–1365, 2016.
- B. Piccoli, F. Rossi, and E. Trélat. Control to flocking of the kinetic Cucker–Smale model. *SIAM Journal on Mathematical Analysis*, 47(6):4685–4719, 2015.
- G.M. Rotskoff and E. Vanden-Eijnden. Trainability and accuracy of neural networks: An interacting particle system approach. *Preprint arXiv:1805.00915, to appear in Comm. Pure App. Math.*, 2018.
- F. Santambrogio. *Optimal Transport for Applied Mathematicians*, volume 55. Springer, 2015.
- J. Sirignano and K. Spiliopoulos. Mean field analysis of neural networks: A law of large numbers. *SIAM Journal on Applied Mathematics*, 80(2):725–752, 2020.
- A-S. Sznitman. Topics in propagation of chaos. In *Ecole d’Eté de Probabilités de Saint-Flour XIX — 1989*, pages 165–251. Springer, 1991. ISBN 978-3-540-46319-1.
- C. Villani. *Topics in Optimal Transportation*, volume 58 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 2003.
- C. Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.

## Appendix A. Proof of Theorem 2

For simplicity, we prove the theorem 2 when  $T = 1$ , and we denote  $\Theta_1$  simply by  $\Theta$ . In this section we assume **A1–A4**.

### A.1. Pre-limit equation (28) and error terms in (28)

#### A.1.1. DERIVATION OF THE PRE-LIMIT EQUATION

The aim of this section is to establish the so-called pre-limit equation (28), which will be our starting point to derive Equation (11). Let  $N \geq 1$ ,  $k \in \{0, \dots, N\}$ , and  $f \in \mathcal{C}^\infty(\Theta)$ . Recall that by Lemma 1 and since  $0 \leq k \leq N$ , a.s.  $\theta_k^i \in \Theta$ , and thus a.s.  $f(\theta_k^i)$  is well-defined. The Taylor-Lagrange formula yields

$$\begin{aligned} \langle f, \nu_{k+1}^N \rangle - \langle f, \nu_k^N \rangle &= \frac{1}{N} \sum_{i=1}^N f(\theta_{k+1}^i) - f(\theta_k^i) \\ &= \frac{1}{N} \sum_{i=1}^N \nabla_\theta f(\theta_k^i) \cdot (\theta_{k+1}^i - \theta_k^i) + \frac{1}{2N} \sum_{i=1}^N (\theta_{k+1}^i - \theta_k^i)^T \nabla^2 f(\widehat{\theta}_k^i) (\theta_{k+1}^i - \theta_k^i), \end{aligned}$$

where, for all  $i \in \{1, \dots, N\}$ ,  $\widehat{\theta}_k^i \in (\theta_k^i, \theta_{k+1}^i) \subset \Theta$ . Using (5), we then obtain

$$\begin{aligned} \langle f, \nu_{k+1}^N \rangle - \langle f, \nu_k^N \rangle &= -\frac{\eta}{N^3} \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\langle \phi(\theta_k^j, \cdot, x_k), \gamma \rangle - y_k) \langle \nabla_\theta f(\theta_k^i) \cdot \nabla_\theta \phi(\theta_k^i, \cdot, x_k), \gamma \rangle \\ &\quad - \frac{\eta}{N^2} \langle (\phi(\cdot, \cdot, x_k) - y_k) \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x_k), \nu_k^N \otimes \gamma \rangle \\ &\quad - \frac{\eta}{N} \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \nu_k^N \rangle + \mathbf{R}_k^N[f], \end{aligned} \quad (18)$$

where

$$\mathbf{R}_k^N[f] := \frac{1}{2N} \sum_{i=1}^N (\theta_{k+1}^i - \theta_k^i)^T \nabla^2 f(\widehat{\theta}_k^i) (\theta_{k+1}^i - \theta_k^i). \quad (19)$$

Let us define

$$\begin{aligned} \mathbf{D}_k^N[f] &:= \mathbf{E} \left[ -\frac{\eta}{N^3} \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\langle \phi(\theta_k^j, \cdot, x_k), \gamma \rangle - y_k) \langle \nabla_\theta f(\theta_k^i) \cdot \nabla_\theta \phi(\theta_k^i, \cdot, x_k), \gamma \rangle \middle| \mathcal{F}_k^N \right] \\ &\quad - \mathbf{E} \left[ \frac{\eta}{N^2} \langle (\phi(\cdot, \cdot, x_k) - y_k) \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x_k), \nu_k^N \otimes \gamma \rangle \middle| \mathcal{F}_k^N \right]. \end{aligned} \quad (20)$$

Note that using (45) and (47) together with the fact that  $|\nabla_\theta f(\theta_k^i)| \leq \sup_{\theta \in \Theta} |\nabla_\theta f(\theta)|$ , the integrand in (20) is integrable and thus  $\mathbf{D}_k^N[f]$  is well defined. Using the fact that  $(x_k, y_k) \perp\!\!\!\perp \mathcal{F}_k^N$  by **A2** and that  $\{\theta_k^i, i = 1, \dots, N\}$  is  $\mathcal{F}_k^N$ -measurable by (5), we have:

$$\begin{aligned} \mathbf{D}_k^N[f] &= -\frac{\eta}{N^3} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \int_{\mathbb{X} \times \mathbb{Y}} (\langle \phi(\theta_k^j, \cdot, x), \gamma \rangle - y) \langle \nabla_\theta f(\theta_k^i) \cdot \nabla_\theta \phi(\theta_k^i, \cdot, x), \gamma \rangle \pi(dx, dy) \\ &\quad - \frac{\eta}{N^2} \int_{\mathbb{X} \times \mathbb{Y}} \langle (\phi(\cdot, \cdot, x) - y) \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \nu_k^N \otimes \gamma \rangle \pi(dx, dy). \end{aligned} \quad (21)$$

Introduce also

$$\begin{aligned} \mathbf{M}_k^N[f] &:= -\frac{\eta}{N^3} \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\langle \phi(\theta_k^j, \cdot, x_k), \gamma \rangle - y_k) \langle \nabla_{\theta} f(\theta_k^i) \cdot \nabla_{\theta} \phi(\theta_k^i, \cdot, x_k), \gamma \rangle \\ &\quad - \frac{\eta}{N^2} (\langle \phi(\cdot, \cdot, x_k) - y_k \rangle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x_k), \nu_k^N \otimes \gamma) - \mathbf{D}_k^N[f]. \end{aligned}$$

Note that  $\mathbf{E}[\mathbf{M}_k^N[f] | \mathcal{F}_k^N] = 0$ . Equation (18) then writes

$$\langle f, \nu_{k+1}^N \rangle - \langle f, \nu_k^N \rangle = \mathbf{D}_k^N[f] + \mathbf{M}_k^N[f] - \frac{\eta}{N} \langle \nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \nu_k^N \rangle + \mathbf{R}_k^N[f]. \quad (22)$$

Notice also that

$$\begin{aligned} \mathbf{D}_k^N[f] &= -\frac{\eta}{N^3} \sum_{i=1}^N \sum_{j=1}^N \int_{\mathbf{X} \times \mathbf{Y}} (\langle \phi(\theta_k^j, \cdot, x), \gamma \rangle - y) \langle \nabla_{\theta} f(\theta_k^i) \cdot \nabla_{\theta} \phi(\theta_k^i, \cdot, x), \gamma \rangle \pi(\mathrm{d}x, \mathrm{d}y) \\ &\quad + \frac{\eta}{N^3} \sum_{i=1}^N \int_{\mathbf{X} \times \mathbf{Y}} (\langle \phi(\theta_k^i, \cdot, x), \gamma \rangle - y) \langle \nabla_{\theta} f(\theta_k^i) \cdot \nabla_{\theta} \phi(\theta_k^i, \cdot, x), \gamma \rangle \pi(\mathrm{d}x, \mathrm{d}y) \\ &\quad - \frac{\eta}{N^2} \int_{\mathbf{X} \times \mathbf{Y}} (\langle \phi(\cdot, \cdot, x) - y \rangle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \nu_k^N \otimes \gamma) \pi(\mathrm{d}x, \mathrm{d}y) \\ &= -\frac{\eta}{N} \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \nu_k^N \otimes \gamma \rangle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \nu_k^N \otimes \gamma \rangle \pi(\mathrm{d}x, \mathrm{d}y) \\ &\quad + \frac{\eta}{N^2} \int_{\mathbf{X} \times \mathbf{Y}} \left\langle (\langle \phi(\cdot, \cdot, x), \gamma \rangle - y) \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \gamma \rangle, \nu_k^N \right\rangle \pi(\mathrm{d}x, \mathrm{d}y) \\ &\quad - \frac{\eta}{N^2} \int_{\mathbf{X} \times \mathbf{Y}} (\langle \phi(\cdot, \cdot, x) - y \rangle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \nu_k^N \otimes \gamma) \pi(\mathrm{d}x, \mathrm{d}y). \end{aligned} \quad (23)$$

Now, we define for  $t \in [0, 1]$ :

$$\mathbf{D}_t^N[f] := \sum_{k=0}^{\lfloor Nt \rfloor - 1} \mathbf{D}_k^N[f], \quad \mathbf{R}_t^N[f] := \sum_{k=0}^{\lfloor Nt \rfloor - 1} \mathbf{R}_k^N[f], \quad \text{and} \quad \mathbf{M}_t^N[f] := \sum_{k=0}^{\lfloor Nt \rfloor - 1} \mathbf{M}_k^N[f]. \quad (24)$$

We can rewrite  $\mathbf{D}_t^N[f]$  has follows:

$$\mathbf{D}_t^N[f] = \sum_{k=0}^{\lfloor Nt \rfloor - 1} \int_{\frac{k}{N}}^{\frac{k+1}{N}} N \mathbf{D}_{\lfloor Ns \rfloor}^N[f] \mathrm{d}s = N \int_0^t \mathbf{D}_{\lfloor Ns \rfloor}^N[f] \mathrm{d}s - N \int_{\frac{\lfloor Nt \rfloor}{N}}^t \mathbf{D}_{\lfloor Ns \rfloor}^N[f] \mathrm{d}s.$$

Since  $\nu_{\lfloor Ns \rfloor}^N = \mu_s^N$  (by definition, see (8)), we have, using also (23) with  $k = \lfloor Ns \rfloor$ ,

$$\begin{aligned} \mathbf{D}_t^N[f] &= -\eta \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \mu_s^N \otimes \gamma \rangle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \rangle \pi(\mathrm{d}x, \mathrm{d}y) \mathrm{d}s \\ &\quad + \frac{\eta}{N} \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \left\langle (\langle \phi(\cdot, \cdot, x), \gamma \rangle - y) \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \gamma \rangle, \mu_s^N \right\rangle \pi(\mathrm{d}x, \mathrm{d}y) \mathrm{d}s \\ &\quad - \frac{\eta}{N} \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \left\langle (\langle \phi(\cdot, \cdot, x) - y \rangle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma) \right\rangle \pi(\mathrm{d}x, \mathrm{d}y) \mathrm{d}s - \mathbf{V}_t^N[f], \end{aligned} \quad (25)$$

where

$$\begin{aligned} \mathbf{V}_t^N[f] &:= -\eta \int_{\lfloor \frac{Nt}{N} \rfloor}^t \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \mu_s^N \otimes \gamma \rangle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \rangle \pi(dx, dy) ds \\ &\quad + \frac{\eta}{N} \int_{\lfloor \frac{Nt}{N} \rfloor}^t \int_{\mathbf{X} \times \mathbf{Y}} \left\langle \langle \phi(\cdot, \cdot, x) - y, \gamma \rangle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \gamma \rangle, \mu_s^N \right\rangle \pi(dx, dy) ds \\ &\quad - \frac{\eta}{N} \int_{\lfloor \frac{Nt}{N} \rfloor}^t \int_{\mathbf{X} \times \mathbf{Y}} \left\langle (\phi(\cdot, \cdot, x) - y) \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \right\rangle \pi(dx, dy) ds. \end{aligned}$$

On the other hand, we also have for  $t \in [0, 1]$ ,

$$\sum_{k=0}^{\lfloor Nt \rfloor - 1} -\frac{\eta}{N} \langle \nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \nu_k^N \rangle = -\eta \int_0^{\lfloor \frac{Nt}{N} \rfloor} \langle \nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \mu_s^N \rangle ds. \quad (26)$$

We finally set:

$$\mathbf{W}_t^N[f] := -\mathbf{V}_t^N[f] + \eta \int_{\lfloor \frac{Nt}{N} \rfloor}^t \langle \nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \mu_s^N \rangle ds. \quad (27)$$

Since  $\langle f, \mu_t^N \rangle - \langle f, \mu_0^N \rangle = \sum_{k=0}^{\lfloor Nt \rfloor - 1} \langle f, \nu_{k+1}^N \rangle - \langle f, \nu_k^N \rangle$ , we deduce from (22), (24), (25), (26) and (27), the so-called pre-limit equation satisfied by  $\mu^N$ : for  $N \geq 1$ ,  $t \in [0, 1]$ , and  $f \in \mathcal{C}^{\infty}(\Theta)$ ,

$$\begin{aligned} \langle f, \mu_t^N \rangle - \langle f, \mu_0^N \rangle &= -\eta \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \mu_s^N \otimes \gamma \rangle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \rangle \pi(dx, dy) ds \\ &\quad - \eta \int_0^t \langle \nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \mu_s^N \rangle ds \\ &\quad + \frac{\eta}{N} \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \left\langle \langle \phi(\cdot, \cdot, x) - y, \gamma \rangle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \gamma \rangle, \mu_s^N \right\rangle \pi(dx, dy) ds \\ &\quad - \frac{\eta}{N} \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \left\langle (\phi(\cdot, \cdot, x) - y) \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \right\rangle \pi(dx, dy) ds \\ &\quad + \mathbf{M}_t^N[f] + \mathbf{W}_t^N[f] + \mathbf{R}_t^N[f]. \end{aligned} \quad (28)$$

#### A.1.2. THE LAST FIVE TERMS IN (28) ARE ERROR TERMS

The purpose of this section is to show that the last five terms appearing in the r.h.s. of (28) are error terms when  $N \rightarrow +\infty$ . For  $J \in \mathbf{N}^*$  and  $f \in \mathcal{C}^J(\Theta)$ , set  $\|f\|_{\mathcal{C}^J(\Theta)} := \sum_{|k| \leq J} \|\partial_k f\|_{\infty, \Theta}$ , where  $\|g\|_{\infty, \Theta} = \sup_{\theta \in \Theta} |g(\theta)|$  for  $g: \Theta \rightarrow \mathbf{R}^m$ .

**Lemma 6 (Error terms)** *Assume A1→A4. Then, there exists  $C > 0$  such that a.s. for all  $f \in \mathcal{C}^{\infty}(\Theta)$  and  $N \geq 1$ ,*

1.  $\frac{\eta}{N} \int_0^1 \int_{\mathbf{X} \times \mathbf{Y}} \left| \left\langle \langle \phi(\cdot, \cdot, x) - y, \gamma \rangle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \gamma \rangle, \mu_s^N \right\rangle \right| \pi(dx, dy) ds \leq C \|f\|_{\mathcal{C}^1(\Theta)} / N.$
2.  $\frac{\eta}{N} \int_0^1 \int_{\mathbf{X} \times \mathbf{Y}} \left| \left\langle (\phi(\cdot, \cdot, x) - y) \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \right\rangle \right| \pi(dx, dy) ds \leq C \|f\|_{\mathcal{C}^1(\Theta)} / N.$

$$3. \sup_{t \in [0,1]} |\mathbf{W}_t^N[f]| + \sup_{t \in [0,1]} |\mathbf{R}_t^N[f]| \leq C \|f\|_{\mathcal{C}^2(\Theta)}/N.$$

Finally,  $\sup_{t \in [0,1]} \mathbf{E}[|\mathbf{M}_t^N[f]|] \leq C \|f\|_{\mathcal{C}^1(\Theta)}/\sqrt{N}$ .

**Proof** All along the proof,  $C > 0$  denotes a positive constant independent of  $N \geq 1, k \in \{0, \dots, N-1\}, (s, t) \in [0, 1]^2, (x, y) \in \mathbf{X} \times \mathbf{Y}, \theta \in \Theta, z \in \mathbf{R}^d$ , and  $f \in \mathcal{C}^\infty(\Theta)$  which can change from one occurrence to another. Using (47), the Cauchy-Schwarz inequality, and the fact that  $\nabla_\theta f$  is bounded over  $\Theta$  imply:

$$|\langle \nabla_\theta f(\theta) \cdot \nabla_\theta \phi(\theta, \cdot, x), \gamma \rangle| \leq \langle |\nabla_\theta f(\theta) \cdot \nabla_\theta \phi(\theta, \cdot, x)|, \gamma \rangle \leq C \|f\|_{\mathcal{C}^1(\Theta)}. \quad (29)$$

Combining (45) and (29), we obtain:

$$\int_0^1 \int_{\mathbf{X} \times \mathbf{Y}} \left| \left\langle \phi(\cdot, \cdot, x) - y, \gamma \right\rangle \left\langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \mu_s^N \right\rangle \right| \pi(dx, dy) ds \leq C \|f\|_{\mathcal{C}^1(\Theta)}$$

and

$$\int_0^1 \int_{\mathbf{X} \times \mathbf{Y}} \left| \left\langle (\phi(\cdot, \cdot, x) - y) \nabla_m f \cdot \nabla_m \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \right\rangle \right| \pi(dx, dy) ds \leq C \|f\|_{\mathcal{C}^1(\Theta)},$$

which proves Items 1 and 2.

Let us now prove Item 3. By (45) and (29),  $\sup_{t \in [0,1]} |\mathbf{V}_t^N[f]| \leq C \|f\|_{\mathcal{C}^1(\Theta)}/N$ . On the other hand, because  $f \in \mathcal{C}^\infty(\Theta)$  and  $\theta \mapsto \nabla_\theta \mathcal{D}_{\text{KL}}(q_\theta^1 | P_0^1)$  is continuous (see (4)) over  $\Theta$  which is compact, it holds,  $\|\nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q_\theta^1 | P_0^1)\|_{\infty, \Theta} < +\infty$ . Hence, it holds:

$$\sup_{t \in [0,1]} \left| \int_{\frac{\lfloor Nt \rfloor}{N}}^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q_\theta^1 | P_0^1), \mu_s^N \rangle ds \right| \leq C \|f\|_{\mathcal{C}^1(\Theta)}/N.$$

Using (27), it then holds  $\sup_{t \in [0,1]} |\mathbf{W}_t^N[f]| \leq C \|f\|_{\mathcal{C}^1(\Theta)}/N$ . Since  $f \in \mathcal{C}^\infty(\Theta)$ , we have, by (19), for  $N \geq 1$  and  $0 \leq k \leq N-1$ ,  $|\mathbf{R}_k^N[f]| \leq \|f\|_{\mathcal{C}^2(\Theta)} \frac{C}{N} \sum_{i=1}^N |\theta_{k+1}^i - \theta_k^i|^2$ . By (48) and Lemma 1,  $|\theta_{k+1}^i - \theta_k^i|^2 \leq C/N^2$  and consequently, one has:

$$|\mathbf{R}_k^N[f]| \leq C \|f\|_{\mathcal{C}^2(\Theta)}/N^2. \quad (30)$$

Hence, for all  $t \in [0, 1]$ ,  $|\mathbf{R}_t^N[f]| \leq C \|f\|_{\mathcal{C}^2(\Theta)}/N$ . This proves Item 3.

Let us now prove the last item in Lemma 6. Let  $t \in [0, 1]$ . We have, by (24),

$$|\mathbf{M}_t^N[f]|^2 = \sum_{k=0}^{\lfloor Nt \rfloor - 1} |\mathbf{M}_k^N[f]|^2 + 2 \sum_{k < j} \mathbf{M}_k^N[f] \mathbf{M}_j^N[f].$$

For all  $0 \leq k < j < \lfloor Nt \rfloor$ ,  $\mathbf{M}_k^N[f]$  is  $\mathcal{F}_j^N$ -measurable (see (9)), and since  $\mathbf{E}[\mathbf{M}_j^N[f] | \mathcal{F}_j^N] = 0$ , one deduces that  $\mathbf{E}[\mathbf{M}_k^N[f] \mathbf{M}_j^N[f]] = \mathbf{E}[\mathbf{M}_k^N[f] \mathbf{E}[\mathbf{M}_j^N[f] | \mathcal{F}_j^N]] = 0$ . Hence,  $\mathbf{E}[|\mathbf{M}_t^N[f]|^2] = \sum_{k=0}^{\lfloor Nt \rfloor - 1} \mathbf{E}[|\mathbf{M}_k^N[f]|^2]$ . By (45) and (29), one has a.s. for all  $0 \leq k \leq N-1$ ,

$$|\mathbf{M}_k^N[f]| \leq C \|f\|_{\mathcal{C}^1(\Theta)}/N. \quad (31)$$

Hence,  $\mathbf{E}[|\mathbf{M}_t^N[f]|^2] \leq C \|f\|_{\mathcal{C}^1(\Theta)}/N$ , which proves the last inequality in Lemma 6.  $\blacksquare$

## A.2. Convergence to the limit equation as $N \rightarrow +\infty$

In this section we prove the relative compactness of  $(\mu^N)_{N \geq 1}$  in  $\mathcal{D}([0, 1], \mathcal{P}(\Theta))$ . We then show that any of its limit points satisfies the limit equation (11).

### A.2.1. WASSERSTEIN SPACES AND DUALITY FORMULA

In this section we recall some basic results which will be used throughout this work on the space  $\mathcal{P}(\mathcal{S})$  when  $(\mathcal{S}, d)$  is a Polish space. First when endowed with the weak convergence topology,  $\mathcal{P}(\mathcal{S})$  is a Polish space (Billingsley, 1999, Theorem 6.8). In addition,  $\mathcal{P}_q(\mathcal{S}) = \{\nu \in \mathcal{P}(\mathcal{S}), \int_{\mathcal{S}} d(w_0, w)^q \nu(dw) < +\infty\}$ , where  $w_0 \in \mathcal{S}$  is arbitrary (note that this space was defined previously in (14) when  $\mathcal{S} = \mathbf{R}^{d+1}$ ) when endowed with the  $W_q$  metric is also a Polish space (Villani, 2009, Theorem 6.18). Recall also the duality formula for the  $W_1$ -distance on  $\mathcal{P}_1(\mathcal{S})$  (see e.g (Villani, 2009, Remark 6.5)):

$$W_1(\mu, \nu) = \sup \left\{ \left| \int_{\mathcal{S}} f(w) d\mu(w) - \int_{\mathcal{S}} f(w) \nu(dw) \right|, \|f\|_{\text{Lip}} \leq 1 \right\}. \quad (32)$$

Finally, when  $\mathcal{K} \subset \mathbf{R}^{d+1}$  is compact, the convergence in  $W_q$ -distance is equivalent to the usual weak convergence on  $\mathcal{P}(\mathcal{K})$  (see e.g. (Villani, 2009, Corollary 6.13)).

### A.2.2. RELATIVE COMPACTNESS

The main result of this section is to prove that  $(\mu^N)_{N \geq 1}$  is relatively compact in  $\mathcal{D}([0, 1], \mathcal{P}(\Theta))$ , which is the purpose of Proposition 8 below. To this end, we need to prove that for all  $f \in \mathcal{C}^\infty(\Theta)$ , every sequence  $(\langle f, \mu_t^N \rangle)_{N \geq 1}$  satisfies some regularity conditions, which is the purpose of the next result.

**Lemma 7 (Regularity condition)** *Assume A1→A4. Then there exists  $C > 0$  such that a.s. for all  $f \in \mathcal{C}^\infty(\Theta)$ ,  $0 \leq r < t \leq 1$ , and  $N \geq 1$ :*

$$|\langle f, \mu_t^N \rangle - \langle f, \mu_r^N \rangle| \leq C \|f\|_{\mathcal{C}^2(\Theta)} \left[ |t - r| + \frac{|t - r|}{N} + \frac{1}{N} \right]. \quad (33)$$

**Proof** Let  $f \in \mathcal{C}^\infty(\Theta)$  and let  $N \geq 1$  and  $0 \leq r < t \leq 1$ . In the following  $C > 0$  is a positive constant independent of  $f \in \mathcal{C}^\infty(\Theta)$ ,  $N \geq 1$ , and  $0 \leq r < t \leq 1$ , which can change from one occurrence to another. From (28), we have

$$\begin{aligned} \langle f, \mu_t^N \rangle - \langle f, \mu_r^N \rangle &= \mathbf{A}_{r,t}^N[f] - \eta \int_r^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \mu_s^N \rangle ds \\ &\quad + \mathbf{M}_t^N[f] - \mathbf{M}_r^N[f] + \mathbf{W}_t^N[f] - \mathbf{W}_r^N[f] + \mathbf{R}_t^N[f] - \mathbf{R}_r^N[f], \end{aligned} \quad (34)$$

where

$$\begin{aligned} \mathbf{A}_{r,t}^N[f] &= -\eta \int_r^t \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \mu_s^N \otimes \gamma \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \rangle \pi(dx, dy) \\ &\quad + \frac{\eta}{N} \int_r^t \int_{\mathbf{X} \times \mathbf{Y}} \left\langle \langle \phi(\cdot, \cdot, x) - y, \gamma \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \gamma \rangle, \mu_s^N \right\rangle \pi(dx, dy) \\ &\quad - \frac{\eta}{N} \int_r^t \int_{\mathbf{X} \times \mathbf{Y}} \left\langle (\phi(\cdot, \cdot, x) - y) \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \right\rangle \pi(dx, dy). \end{aligned}$$

By (45) and (29),  $|\mathbf{A}_{r,t}^N[f]| \leq C\|f\|_{\mathcal{C}^1(\Theta)}[|t-r| + \frac{|t-r|}{N}]$ . In addition, since  $\theta \mapsto \mathcal{D}_{\text{KL}}(q_\theta^1|P_0^1)$  is bounded over  $\Theta$  (since it is smooth and  $\Theta$  is compact),

$$\left| \int_r^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q_s^1|P_0^1), \mu_s^N \rangle ds \right| \leq C\|f\|_{\mathcal{C}^1(\Theta)}|t-r|.$$

Furthermore, using (31),

$$|\mathbf{M}_t^N[f] - \mathbf{M}_r^N[f]| = \left| \sum_{k=\lfloor Nr \rfloor}^{\lfloor Nt \rfloor - 1} \mathbf{M}_k^N[f] \right| \leq (\lfloor Nt \rfloor - \lfloor Nr \rfloor)C\|f\|_{\mathcal{C}^1(\Theta)}/N.$$

Next, we have, by Item 3 in Lemma 6,  $|\mathbf{W}_t^N[f] - \mathbf{W}_r^N[f]| \leq |\mathbf{W}_t^N[f]| + |\mathbf{W}_r^N[f]| \leq C\|f\|_{\mathcal{C}^2(\Theta)}/N$ . Finally, by (30),

$$|\mathbf{R}_t^N[f] - \mathbf{R}_r^N[f]| = \left| \sum_{k=\lfloor Nr \rfloor}^{\lfloor Nt \rfloor - 1} \mathbf{R}_k^N[f] \right| \leq (\lfloor Nt \rfloor - \lfloor Nr \rfloor)C\|f\|_{\mathcal{C}^2(\Theta)}/N^2.$$

The proof of Proposition 7 is complete plugging all the previous estimates in (34).  $\blacksquare$

**Proposition 8 (Relative compactness)** *Assume  $\mathbf{A1} \rightarrow \mathbf{A4}$ . Then, the sequence  $(\mu^N)_{N \geq 1}$  is relatively compact in  $\mathcal{D}([0, 1], \mathcal{P}(\Theta))$ .*

**Proof** The proof consists in applying (Jakubowski, 1986, Theorem 3.1) with  $E = \mathcal{P}(\Theta)$  endowed with the weak convergence topology. Set  $\mathbb{F} = \{\mathfrak{L}_f, f \in \mathcal{C}^\infty(\Theta)\}$  where

$$\mathfrak{L}_f : \nu \in \mathcal{P}(\Theta) \mapsto \langle f, \nu \rangle.$$

The class of continuous functions  $\mathbb{F}$  on  $\mathcal{P}(\Theta)$  satisfies Conditions (Jakubowski, 1986, (3.1) and (3.2) in Theorem 3.1).

On the other hand, the condition (Jakubowski, 1986, (3.3) in Theorem 3.1) is satisfied since  $\mathcal{P}(\Theta)$  is compact because  $\Theta$  is compact (see e.g. (Panaretos and Zemel, 2020, Corollary 2.2.5) together with (Villani, 2009, Corollary 6.13)).

It remains to verify Condition (3.4) of (Jakubowski, 1986, Theorem 3.1), i.e. that for all  $f \in \mathcal{C}^\infty(\Theta)$ ,  $(\langle f, \mu^N \rangle)_{N \geq 1}$  is relatively compact in  $\mathcal{D}([0, 1], \mathbf{R})$ . To this end, we apply (Billingsley, 1999, Theorem 13.2). Condition (i) in (Billingsley, 1999, Theorem 13.2) is satisfied because  $|\langle f, \mu_t^N \rangle| \leq \|f\|_{\infty, \Theta}$  for all  $t \in [0, 1]$  and  $N \geq 1$ . Let us now show that Condition (ii) in (Billingsley, 1999, Theorem 13.2) holds. For this purpose, we use Lemma 7. For  $\delta, \beta > 0$  sufficiently small, it is possible to construct a subdivision  $\{t_i\}_{i=0}^v$  of  $[0, 1]$  such that  $t_0 = 0, t_v = 1, t_{i+1} - t_i = \delta + \beta$  for  $i \in \{0, \dots, v-2\}$  and  $\delta + \beta \leq t_v - t_{v-1} \leq 2(\delta + \beta)$ . According to the terminology introduced in (Billingsley, 1999, Section 12),  $\{t_i\}_{i=0}^v$  is  $\delta$ -sparse. Then, by Lemma 7, there exists  $C > 0$  such that a.s. for all  $\delta, \beta > 0$ , all such subdivision  $\{t_i\}_{i=0}^v, i \in \{0, \dots, v-1\}$ , and  $N \geq 1$ ,

$$\sup_{t, r \in [t_i, t_{i+1}]} |\langle f, \mu_t^N \rangle - \langle f, \mu_r^N \rangle| \leq C \left( |t_{i+1} - t_i| + \frac{|t_{i+1} - t_i|}{N} + \frac{1}{N} \right) \leq C \left( 2(\delta + \beta) + \frac{2(\delta + \beta)}{N} + \frac{1}{N} \right).$$



Thus, one has:

$$\inf_{\beta > 0} \max_i \sup_{t, r \in [t_i, t_{i+1}]} |\langle f, \mu_t^N \rangle - \langle f, \mu_r^N \rangle| \leq C \left( 2\delta + \frac{2\delta}{N} + \frac{1}{N} \right).$$

Consequently, there exists  $C > 0$  such that a.s. for all  $\delta > 0$  small enough and  $N \geq 1$ ,

$$w'_{\langle f, \mu^N \rangle}(\delta) := \inf_{\substack{\{t_i\} \\ \delta\text{-sparse}}} \max_i \sup_{t, r \in [t_i, t_{i+1}]} |\langle f, \mu_t^N \rangle - \langle f, \mu_r^N \rangle| \leq C \left( 2\delta + \frac{2\delta}{N} + \frac{1}{N} \right).$$

This implies  $\lim_{\delta \rightarrow 0} \limsup_{N \rightarrow +\infty} \mathbf{E}[w'_{\langle f, \mu^N \rangle}(\delta)] = 0$ . By Markov's inequality, this proves Condition (ii) of (Billingsley, 1999, Theorem 13.2). Therefore, for all  $f \in C^\infty(\Theta)$ , using also Prokhorov theorem, the sequence  $(\langle f, \mu^N \rangle)_{N \geq 1} \subset \mathcal{D}([0, 1], \mathbf{R})$  is relatively compact. In conclusion, according to (Jakubowski, 1986, Theorem 3.1),  $(\mu^N)_{N \geq 1} \subset \mathcal{D}([0, 1], \mathcal{P}(\Theta))$  is tight. ■

### A.2.3. LIMIT POINTS SATISFY THE LIMIT EQUATION (11)

In this section we prove that every limit point of  $(\mu^N)_{N \geq 1}$  in  $\mathcal{D}([0, 1], \mathcal{P}(\Theta))$  satisfies (11).

**Lemma 9** *Let  $\mathfrak{m}, (\mathfrak{m}^N)_{N \geq 1} \subset \mathcal{D}([0, 1], \mathcal{P}(\Theta))$  be such that  $\mathfrak{m}^N \rightarrow \mathfrak{m}$  in  $\mathcal{D}([0, 1], \mathcal{P}(\Theta))$ . Then, for all Lipschitz continuous function  $f : \Theta \rightarrow \mathbf{R}$ , we have  $\langle f, \mathfrak{m}^N \rangle \rightarrow \langle f, \mathfrak{m} \rangle$  in  $\mathcal{D}([0, 1], \mathbf{R})$ .*

**Proof** Let  $f$  be such a function. By (Billingsley, 1999, p.124),  $\mathfrak{m}^N \rightarrow \mathfrak{m}$  in  $\mathcal{D}([0, 1], \mathcal{P}(\Theta))$  iff there exist functions  $\lambda_N : [0, 1] \rightarrow [0, 1]$  continuous, increasing onto itself such that  $\sup_{t \in [0, 1]} |\lambda_N(t) - t| \rightarrow_{N \rightarrow \infty} 0$  and  $\sup_{t \in [0, 1]} \mathbb{W}_1(\mathfrak{m}_{\lambda_N(t)}^N, \mathfrak{m}_t) \rightarrow_{N \rightarrow \infty} 0$ . Then  $\langle f, \mathfrak{m}^N \rangle \rightarrow \langle f, \mathfrak{m} \rangle$  in  $\mathcal{D}([0, 1], \mathbf{R})$  since by (32),  $\sup_{t \in [0, 1]} |\langle f, \mathfrak{m}_{\lambda_N(t)}^N \rangle - \langle f, \mathfrak{m}_t \rangle| \leq \|f\|_{\text{Lip}} \sup_{t \in [0, 1]} \mathbb{W}_1(\mathfrak{m}_{\lambda_N(t)}^N, \mathfrak{m}_t) \rightarrow_{N \rightarrow \infty} 0$ . ■

**Proposition 10 (Continuity of the limit points of  $\langle f, \mu^N \rangle$ )** *Let  $f \in C^\infty(\Theta)$ . Then, any limit point of  $(\langle f, \mu^N \rangle)_{N \geq 1} \subset \mathcal{D}([0, 1], \mathbf{R})$  belong a.s. to  $\mathcal{C}([0, 1], \mathbf{R})$ .*

**Proof** Fix  $t \in (0, 1]$ . Letting  $r \rightarrow t$  in (33), we obtain  $|\langle f, \mu_t^N \rangle - \langle f, \mu_{t-}^N \rangle| \leq C/N$ . Therefore  $\sup_{t \in (0, 1]} |\langle f, \mu_t^N \rangle - \langle f, \mu_{t-}^N \rangle| \xrightarrow{\mathcal{D}} 0$  as  $N \rightarrow +\infty$ . The result follows from (Billingsley, 1999, Theorem 13.4). ■

**Proposition 11 (Continuity of the limit points of  $\mu^N$ )** *Let  $\mu^* \in \mathcal{D}([0, 1], \mathcal{P}(\Theta))$  be a limit point of  $(\mu^N)_{N \geq 1} \subset \mathcal{D}([0, 1], \mathcal{P}(\Theta))$ . Then, a.s.  $\mu^* \in \mathcal{C}([0, 1], \mathcal{P}(\Theta))$ .*

**Proof** Up to extracting a subsequence, we assume that  $\mu^N \xrightarrow{\mathcal{D}} \mu^*$ . By Skorohod representation theorem, there exists another probability space  $(\hat{\Omega}, \hat{\mathcal{F}}, \hat{\mathbf{P}})$  on which are defined random elements  $(\hat{\mu}^N)_{N \geq 1}$  and  $\hat{\mu}^*$ , where,

$$\hat{\mu}^* \stackrel{\mathcal{D}}{=} \mu^*, \quad \text{and for all } N \geq 1, \hat{\mu}^N \stackrel{\mathcal{D}}{=} \mu^N,$$

and such that  $\hat{\mathbf{P}}$ -a.s.,  $\hat{\mu}^N \rightarrow \hat{\mu}^*$  in  $\mathcal{D}([0, 1], \mathcal{P}(\Theta))$  as  $N \rightarrow +\infty$ . Fix  $f \in \mathcal{C}^\infty(\Theta)$ . We have, by Lemma 9,

$$\hat{\mathbf{P}}\text{-a.s.}, \langle f, \hat{\mu}^N \rangle \rightarrow_{N \rightarrow +\infty} \langle f, \hat{\mu}^* \rangle \text{ in } \mathcal{D}([0, 1], \mathbf{R}).$$

In particular,  $\langle f, \hat{\mu}^N \rangle \rightarrow_{N \rightarrow +\infty} \langle f, \hat{\mu}^* \rangle$  in distribution. By Proposition 10, there exists  $\hat{\Omega}_f \subset \hat{\Omega}$  of  $\hat{\mathbf{P}}$ -mass 1 such that for all  $\omega \in \hat{\Omega}_f$ ,  $\langle f, \hat{\mu}^*(\omega) \rangle \in \mathcal{C}([0, 1], \mathbf{R})$ . Denote by  $\mathcal{F}$  the class polynomial functions with rational coefficients. Since this class is countable, the set  $\hat{\Omega}_{\mathcal{F}} := \bigcap_{f \in \mathcal{F}} \hat{\Omega}_f$  is of  $\hat{\mathbf{P}}$ -mass 1. Consider now an arbitrary  $f \in \mathcal{C}(\Theta)$  and let us show that for all  $\omega \in \hat{\Omega}_{\mathcal{F}}$ ,  $\langle f, \hat{\mu}^*(\omega) \rangle \in \mathcal{C}([0, 1], \mathbf{R})$ . By the Stone-Weierstrass theorem, there exist  $(f_n)_{n \geq 1} \subset \mathcal{F}$  such that  $\|f_n - f\|_{\infty, \Theta} \rightarrow_{n \rightarrow +\infty} 0$ . On  $\hat{\Omega}_{\mathcal{F}}$ , for all  $n, t \in [0, 1] \mapsto \langle f_n, \hat{\mu}_t^* \rangle$  is continuous and converges uniformly to  $t \in [0, 1] \mapsto \langle f, \hat{\mu}_t^* \rangle$ . Hence, for all  $\omega \in \hat{\Omega}_{\mathcal{F}}$  and  $f \in \mathcal{C}(\Theta)$ ,  $\langle f, \hat{\mu}^*(\omega) \rangle \in \mathcal{C}([0, 1], \mathbf{R})$ , i.e. for all  $\omega \in \hat{\Omega}_{\mathcal{F}}$ ,  $\hat{\mu}^*(\omega) \in \mathcal{C}([0, 1], \mathcal{P}(\Theta))$ . This concludes the proof.  $\blacksquare$

Now, we introduce, for  $t \in [0, 1]$  and  $f \in \mathcal{C}^\infty(\Theta)$ , the function  $\Lambda_t[f] : \mathcal{D}([0, 1], \mathcal{P}(\Theta)) \rightarrow \mathbf{R}_+$  defined by:

$$\begin{aligned} \Lambda_t[f] : \mathbf{m} \mapsto & \left| \langle f, \mathbf{m}_t \rangle - \langle f, \mu_0 \rangle \right. \\ & + \eta \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \mathbf{m}_s \otimes \gamma \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \mathbf{m}_s \otimes \gamma \rangle \pi(dx, dy) ds \\ & \left. + \eta \int_0^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \mathbf{m}_s \rangle ds \right|. \end{aligned} \quad (35)$$

We now study the continuity of  $\Lambda_t[f]$ .

**Lemma 12** *Let  $(\mathbf{m}^N)_{N \geq 1} \subset \mathcal{D}([0, 1], \mathcal{P}(\Theta))$  converge to  $\mathbf{m} \in \mathcal{D}([0, 1], \mathcal{P}(\Theta))$ . Then, for all continuity point  $t \in [0, 1]$  of  $\mathbf{m}$  and all  $f \in \mathcal{C}^\infty(\Theta)$ , we have  $\Lambda_t[f](\mathbf{m}^N) \rightarrow \Lambda_t[f](\mathbf{m})$ .*

**Proof** Let  $f \in \mathcal{C}^\infty(\Theta)$  and denote by  $\mathcal{C}(\mathbf{m}) \subset [0, 1]$  the set of continuity points of  $\mathbf{m}$ . Let  $t \in \mathcal{C}(\mathbf{m})$ . From (Billingsley, 1999, p. 124), we have, for all  $s \in \mathcal{C}(\mathbf{m})$ ,

$$\mathbf{m}_s^N \rightarrow \mathbf{m}_s \text{ in } \mathcal{P}(\Theta). \quad (36)$$

Thus,  $\langle f, \mathbf{m}_t^N \rangle \rightarrow_{N \rightarrow \infty} \langle f, \mathbf{m}_t \rangle$ . For all  $z \in \mathbf{R}^d$  and  $(x, y) \in \mathbf{X} \times \mathbf{Y}$ , **A1** and **A3** ensure that the functions  $\theta \in \Theta \mapsto \phi(\theta, z, x) - y$  and  $\theta \in \Theta \mapsto \nabla_\theta f(\theta) \cdot \nabla_\theta \phi(\theta, z, x)$  are continuous and also bounded because  $\Theta$  is compact. Hence, for all  $s \in [0, t] \cap \mathcal{C}(\mathbf{m})$ , using (36),

$$\langle \phi(\cdot, z, x) - y, \mathbf{m}_s^N \rangle \rightarrow \langle \phi(\cdot, z, x) - y, \mathbf{m}_s \rangle \text{ and } \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, z, x), \mathbf{m}_s^N \rangle \rightarrow \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, z, x), \mathbf{m}_s \rangle$$

Since  $[0, 1] \setminus \mathcal{C}(\mathbf{m})$  is at most countable (see (Billingsley, 1999, p. 124)) we have that for a.e.  $(s, z', z, x, y) \in [0, t] \times \mathbf{R}^d \times \mathbf{R}^d \times \mathbf{X} \times \mathbf{Y}$ ,

$$\langle \phi(\cdot, z', x) - y, \mathbf{m}_s^N \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, z, x), \mathbf{m}_s^N \rangle \rightarrow \langle \phi(\cdot, z', x) - y, \mathbf{m}_s \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, z, x), \mathbf{m}_s \rangle.$$

Since  $\phi(\theta, z', x) - y$  is bounded and by (46), there exists  $C > 0$  such that for all  $(s, z', z, x, y) \in [0, t] \times \mathbf{R}^d \times \mathbf{R}^d \times \mathbf{X} \times \mathbf{Y}$ ,  $\langle |\phi(\cdot, z', x) - y|, \mathbf{m}_s^N \rangle \langle |\nabla_\theta f \cdot \nabla_\theta \phi(\cdot, z, x)|, \mathbf{m}_s^N \rangle \leq C \|\nabla_\theta f\|_{\infty, \Theta} \mathbf{b}(z)$ .

By the dominated convergence theorem, we then have:

$$\begin{aligned} & \int_0^t \int_{\mathbb{X} \times \mathbb{Y}} \langle \phi(\cdot, \cdot, x) - y, \mathfrak{m}_s^N \otimes \gamma \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \mathfrak{m}_s^N \otimes \gamma \rangle \pi(dx, dy) ds \\ & \xrightarrow{N \rightarrow +\infty} \int_0^t \int_{\mathbb{X} \times \mathbb{Y}} \langle \phi(\cdot, \cdot, x) - y, \mathfrak{m}_s \otimes \gamma \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), \mathfrak{m}_s \otimes \gamma \rangle \pi(dx, dy) ds. \end{aligned}$$

With the same arguments as above, one shows that  $\int_0^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \mathfrak{m}_s^N \rangle ds \rightarrow \int_0^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \mathfrak{m}_s \rangle ds$ . The proof of the lemma is complete.  $\blacksquare$

**Proposition 13 (Convergence to the limit equation)** *Let  $\mu^* \in \mathcal{D}([0, 1], \mathcal{P}(\Theta))$  be a limit point of  $(\mu^N)_{N \geq 1} \subset \mathcal{D}([0, 1], \mathcal{P}(\Theta))$ . Then, a.s.  $\mu^*$  satisfies (11).*

**Proof** Up to extracting a subsequence, we can assume that  $\mu^N \xrightarrow{\mathcal{D}} \mu^*$  as  $N \rightarrow +\infty$ . Let  $f \in \mathcal{C}^\infty(\Theta)$ . The pre-limit equation (28) and Lemma 6 imply that a.s. for all  $N \geq 1$  and  $t \in [0, 1]$ ,  $\Lambda_t[f](\mu^N) \leq C/N + \mathbf{M}_t^N[f]$ . Hence, using the last statement in Lemma 6, it holds for all  $t \in [0, 1]$ ,

$$\lim_{N \rightarrow \infty} \mathbf{E}[\Lambda_t[f](\mu^N)] = 0.$$

In particular,  $\Lambda_t[f](\mu^N) \xrightarrow{\mathcal{D}} 0$ . Let us now show that  $\Lambda_t[f](\mu^N) \xrightarrow{\mathcal{D}} \Lambda_t[f](\mu^*)$ . Denoting by  $\mathbf{D}(\Lambda_t[f])$  the set of discontinuity points of  $\Lambda_t[f]$ , we have, from Proposition 11 and Lemma 12, for all  $t \in [0, 1]$  and  $f \in \mathcal{C}^\infty(\Theta)$ ,

$$\mathbf{P}(\mu^* \in \mathbf{D}(\Lambda_t[f])) = 0.$$

By the continuous mapping theorem,  $\Lambda_t[f](\mu^N) \xrightarrow{\mathcal{D}} \Lambda_t[f](\mu^*)$ . By uniqueness of the limit in distribution, we have that for all  $t \in [0, 1]$  and  $f \in \mathcal{C}^\infty(\Theta)$ , a.s.  $\Lambda_t[f](\mu^*) = 0$ . Let us now prove that a.s. for all  $t \in [0, 1]$  and  $f \in \mathcal{C}^\infty(\Theta)$ ,  $\Lambda_t[f](\mu^*) = 0$ .

On the one hand, for all  $f \in \mathcal{C}^\infty(\Theta)$  and  $\mathfrak{m} \in \mathcal{D}([0, 1], \mathcal{P}(\Theta))$ , the function  $t \mapsto \Lambda_t[f](\mathfrak{m})$  is right-continuous. Since  $[0, 1]$  is separable, we have that for all  $f \in \mathcal{C}^\infty(\Theta)$ , a.s. for all  $t \in [0, 1]$ ,  $\Lambda_t[f](\mu^*) = 0$ .

One the other hand  $\mathcal{C}^\infty(\Theta)$  is separable (when endowed with the norm  $\|f\|_{\mathcal{C}^\infty(\Theta)} = \sum_{k \geq 0} 2^{-k} \min(1, \sum_{|j|=k} \|\partial_j f\|_{\infty, \Theta})$ ) and the function  $f \in \mathcal{C}^\infty(\Theta) \mapsto \Lambda_t[f](\mathfrak{m})$  is continuous (for fixed  $t \in [0, 1]$  and  $\mathfrak{m} \in \mathcal{D}([0, 1], \mathcal{P}(\Theta))$ ) relatively to the topology induced by  $\|f\|_{\mathcal{C}^\infty(\Theta)}$ .

Hence, we obtain that a.s. for all  $t \in [0, 1]$  and  $f \in \mathcal{C}^\infty(\Theta)$ ,  $\Lambda_t[f](\mu^*) = 0$ . The proof of the proposition is thus complete.  $\blacksquare$

#### A.2.4. UNIQUENESS AND END OF THE PROOF OF THEOREM 2

**Proposition 14** *There exists a unique solution to (11) in  $\mathcal{C}([0, 1], \mathcal{P}(\Theta))$ .*

**Proof** First of all, the fact that there is a solution to (11) is provided by Propositions 8, 11 and 13. The proof of the fact that there is a unique solution to (11) relies on the same arguments as those used in the proof of (Descours et al., 2022, Proposition 2.14).

For  $\mu \in \mathcal{P}(\mathbf{R}^{d+1})$ , we introduce  $\mathbf{v}[\mu] : \mathbf{R}^{d+1} \rightarrow \mathbf{R}^{d+1}$  defined, for  $\theta = (m, \rho) \in \mathbf{R}^{d+1}$ , by

$$\mathbf{v}[\mu](\theta) = -\eta \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \mu \otimes \gamma \rangle \langle \nabla_{\theta} \phi(\theta, \cdot, x), \gamma \rangle \pi(dx, dy) - \eta \nabla_{\theta} \mathcal{D}_{\text{KL}}(q_{\theta}^1 | P_0^1). \quad (37)$$

In addition, if  $\bar{\mu} \in \mathcal{C}([0, 1], \mathcal{P}(\Theta))$  is solution to (11), it satisfies also (11) with test functions  $f \in \mathcal{C}_c^{\infty}(\mathbf{R}^{d+1})$ . Then, adopting the terminology of (Santambrogio, 2015, Section 4.1.2), any solution  $\bar{\mu}$  to (11) is a *weak solution*<sup>1</sup> on  $[0, T]$  of the measure-valued equation

$$\begin{cases} \partial_t \bar{\mu}_t = \text{div}(\mathbf{v}[\bar{\mu}_t] \bar{\mu}_t) \\ \bar{\mu}_0 = \mu_0. \end{cases} \quad (38)$$

Let us now prove that:

1. There exists  $C > 0$  such that for all  $\mu \in \mathcal{P}(\mathbf{R}^{d+1})$  and  $\theta \in \mathbf{R}^{d+1}$ ,

$$|\mathbf{J}_{\theta} \mathbf{v}[\mu](\theta)| \leq C.$$

2. There exists  $C > 0$  such that for all  $\bar{\mu} \in \mathcal{C}([0, 1], \mathcal{P}(\Theta))$  solution to (11),  $0 \leq s, t \leq 1$ , and  $\theta \in \mathbf{R}^{d+1}$ ,

$$|\mathbf{v}[\bar{\mu}_t](\theta) - \mathbf{v}[\bar{\mu}_s](\theta)| \leq C|t - s|.$$

3. There exists  $L' > 0$  such that for all  $\mu, \nu \in \mathcal{P}_1(\mathbf{R}^{d+1})$ ,

$$\sup_{\theta \in \mathbf{R}^d} |\mathbf{v}[\mu](\theta) - \mathbf{v}[\nu](\theta)| \leq L' W_1(\mu, \nu).$$

Before proving the three items above, we quickly conclude the proof of the proposition. Items 1 and 2 above imply that  $v(t, \theta) = \mathbf{v}[\bar{\mu}_t](\theta)$  is globally Lipschitz continuous over  $[0, 1] \times \mathbf{R}^{d+1}$  when  $\bar{\mu} \in \mathcal{C}([0, 1], \mathcal{P}(\Theta))$  is a solution to (11). Since  $\bar{\mu} \in \mathcal{C}([0, 1], \mathcal{P}(\Theta)) \subset \mathcal{C}([0, 1], \mathcal{P}(\mathbf{R}^{d+1}))$ , this allows to use the representation theorem (Villani, 2003, Theorem 5.34) for the solution of (38) in  $\mathcal{C}([0, 1], \mathcal{P}(\mathbf{R}^{d+1}))$ , i.e. it holds:

$$\forall t \in [0, 1], \bar{\mu}_t = \phi_t \# \mu_0, \quad (39)$$

where  $\phi_t$  is the flow generated by the vector field  $\mathbf{v}[\bar{\mu}_t](\theta)$  over  $\mathbf{R}^{d+1}$ . Equation (39) and the fact that  $\mathcal{C}([0, 1], \mathcal{P}(\Theta)) \subset \mathcal{C}([0, 1], \mathcal{P}_1(\mathbf{R}^{d+1}))$  together with Item 3 above and the same arguments as those used in the proof of (Descours et al., 2022, Proposition 2.14) (which we recall is based estimates in Wasserstein distances between two solutions of (11) derived in Piccoli and Rossi (2016)), one deduces that there is a unique solution to (11).

Let us prove Item 1. Recall  $g(\rho) = \ln(1 + e^{\rho})$ . The functions

$$\rho \mapsto g''(\rho)g(\rho), \quad \rho \mapsto g'(\rho), \quad \rho \mapsto \frac{g'(\rho)}{g(\rho)}, \quad \text{and} \quad \rho \mapsto \frac{g''(\rho)}{g(\rho)}$$

1. We mention that according to (Santambrogio, 2015, Proposition 4.2), the two notions of solutions of (38) (namely the weak solution and the *distributional* solution) are equivalent.

are bounded on  $\mathbf{R}$ . Thus, in view of (4),  $\|\text{Hess}_\theta \mathcal{D}_{\text{KL}}(q_\theta^1 | P_0^1)\|_{\infty, \mathbf{R}^{d+1}} < +\infty$ . On the other hand, by **A1** and **A3**, for  $x \in \mathbf{X}$ ,  $z \in \mathbf{R}^d$ ,  $\theta \in \Theta \mapsto \phi(\theta, z, x)$  is smooth and there exists  $C > 0$ , for all  $x \in \mathbf{X}$ ,  $\theta \in \mathbf{R}^{d+1}$ ,  $z \in \mathbf{R}^d$ :

$$|\text{Hess}_\theta \phi(\theta, z, x)| \leq C(\mathfrak{b}(z)^2 + \mathfrak{b}(z)).$$

This bound allows us to differentiate under the integral signs in (37) and proves that  $|\mathbf{J}_\theta \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \mu \otimes \gamma \rangle \langle \nabla_\theta \phi(\theta, \cdot, x), \gamma \rangle \pi(dx, dy)| \leq C$ , where  $C > 0$  is independent of  $\mu \in \mathcal{P}(\Theta)$  and  $\theta \in \Theta$ . The proof of Item 1 is complete.

Let us prove Item 2. Let  $\bar{\mu} \in \mathcal{C}([0, 1], \mathcal{P}(\Theta))$  be a solution to (11),  $0 \leq s \leq t \leq 1$ , and  $\theta \in \mathbf{R}^{d+1}$ . We have

$$\mathbf{v}[\bar{\mu}_t](\theta) - \mathbf{v}[\bar{\mu}_s](\theta) = -\eta \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x), (\bar{\mu}_t - \bar{\mu}_s) \otimes \gamma \rangle \langle \nabla_\theta \phi(\theta, \cdot, x), \gamma \rangle \pi(dx, dy). \quad (40)$$

Let  $z \in \mathbf{R}^d$  and  $x \in \mathbf{X}$ . By **A1** and **A3**,  $\phi(\cdot, z, x) \in \mathcal{C}^\infty(\Theta)$ . Therefore, by (11),

$$\begin{aligned} \langle \phi(\cdot, z, x), \bar{\mu}_t - \bar{\mu}_s \rangle &= -\eta \int_s^t \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x') - y, \bar{\mu}_r \otimes \gamma \rangle \langle \nabla_\theta \phi(\cdot, z, x) \cdot \nabla_\theta \phi(\cdot, \cdot, x'), \bar{\mu}_r \otimes \gamma \rangle \pi(dx', dy) dr \\ &\quad - \eta \int_s^t \langle \nabla_\theta \phi(\cdot, z, x) \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q_r^1 | P_0^1), \bar{\mu}_r \rangle dr \end{aligned}$$

We have  $\|\nabla_\theta \mathcal{D}_{\text{KL}}(q_\theta^1 | P_0^1)\|_{\infty, \Theta} < +\infty$ . Using also (46) and the fact that  $\mathbf{X} \times \mathbf{Y}$  is a compact (see **A2**), it holds:

$$|\langle \phi(\cdot, z, x), \bar{\mu}_t - \bar{\mu}_s \rangle| \leq C\mathfrak{b}(z)|t - s|.$$

Hence, for all  $x' \in \mathbf{X}$ ,

$$|\langle \phi(\cdot, \cdot, x'), (\bar{\mu}_t - \bar{\mu}_s) \otimes \gamma \rangle| \leq \langle |\langle \phi(\cdot, \cdot, x'), \bar{\mu}_t - \bar{\mu}_s \rangle|, \gamma \rangle \leq C|t - s|.$$

Thus, by (40) and (47),  $|\mathbf{v}[\bar{\mu}_t](\theta) - \mathbf{v}[\bar{\mu}_s](\theta)| \leq C|t - s|$ . This ends the proof of Item 2.

Let us now prove Item 3. Fix  $\mu, \nu \in \mathcal{P}_1(\mathbf{R}^{d+1})$  and  $\theta \in \mathbf{R}^{d+1}$ . We have

$$\mathbf{v}[\mu](\theta) - \mathbf{v}[\nu](\theta) = -\eta \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x), (\mu - \nu) \otimes \gamma \rangle \langle \nabla_\theta \phi(\theta, \cdot, x), \gamma \rangle \pi(dx, dy) \quad (41)$$

For all  $x \in \mathbf{X}$ , using (32) and (46), it holds:

$$\begin{aligned} |\langle \phi(\cdot, \cdot, x), (\mu - \nu) \otimes \gamma \rangle| &\leq \int_{\mathbf{R}^d} |\langle \phi(\cdot, z, x), \mu \rangle - \langle \phi(\cdot, z, x), \nu \rangle| \gamma(z) dz \\ &\leq C \int_{\mathbf{R}^d} W_1(\mu, \nu) \mathfrak{b}(z) \gamma(z) dz \leq CW_1(\mu, \nu). \end{aligned}$$

Finally, using in addition (47) and (41), we deduce Item 3.

This ends the proof of the proposition. ■

We are now ready to prove Theorem 2.

**Proof** [Proof of Theorem 2] Recall Lemma 1 ensures that a.s.  $(\mu^N)_{N \geq 1} \subset \mathcal{D}([0, 1], \mathcal{P}(\Theta))$ . By Proposition 8, this sequence is relatively compact. Let  $\mu^* \in \mathcal{D}([0, 1], \mathcal{P}(\Theta))$  be a limit point. Along some subsequence  $N'$ , it holds:

$$\mu^{N'} \xrightarrow{\mathcal{D}} \mu^*.$$

In addition, a.s.  $\mu^* \in \mathcal{C}([0, 1], \mathcal{P}(\Theta))$  (by Proposition 11) and  $\mu^*$  satisfies (11) (by Proposition 13). By Proposition 14, (11) admits a unique solution  $\bar{\mu} \in \mathcal{C}([0, 1], \mathcal{P}(\Theta))$ . Hence, a.s.  $\mu^* = \bar{\mu}$ . Therefore,

$$\mu^{N'} \xrightarrow{\mathcal{D}} \bar{\mu}.$$

Since the sequence  $(\mu^N)_{N \geq 1}$  admits a unique limit point, the whole sequence converges in distribution to  $\bar{\mu}$ . The convergence also holds in probability since  $\bar{\mu}$  is deterministic. The proof of Theorem 2 is complete.  $\blacksquare$

### A.3. Proof of Lemma 1

In this section we prove Lemma 1. We start with the following simple result.

**Lemma 15** *Let  $T > 0$ ,  $N \geq 1$ , and  $c_1 > 0$ . Consider a sequence  $(u_k)_{0 \leq k \leq \lfloor NT \rfloor} \subset \mathbf{R}_+$  for which there exists  $v_0$  such that  $u_0 \leq v_0$  and for all  $1 \leq k \leq \lfloor NT \rfloor$ ,  $u_k \leq c_1(1 + \frac{1}{N} \sum_{\ell=0}^{k-1} u_\ell)$ . Then, for all  $0 \leq k \leq \lfloor NT \rfloor$ ,  $u_k \leq v_0 e^{c_1 T}$ .*

**Proof** Define  $v_k = c_1(1 + \frac{1}{N} \sum_{\ell=0}^{k-1} v_\ell)$ . For all  $0 \leq k \leq \lfloor NT \rfloor$ ,  $u_k \leq v_k$  and  $v_k = v_{k-1}(1 + c_1/N)$ . Hence  $v_k = v_0(1 + c_1/N)^k \leq v_0(1 + c_1/N)^{\lfloor NT \rfloor} \leq v_0 e^{c_1 T}$ . This ends the proof of the Lemma.  $\blacksquare$

**Proof** [Proof of Lemma 1] Since  $\rho \mapsto g'(\rho)$  and  $\rho \mapsto g'(\rho)/g(\rho)$  are bounded continuous functions over  $\mathbf{R}$ , and since  $|g(\rho)| \leq C(1 + |\rho|)$ , according to (4), there exists  $c > 0$ , for all  $\theta \in \mathbf{R}^{d+1}$ ,

$$|\nabla_\theta \mathcal{D}_{\text{KL}}(q_\theta^1 | P_0^1)| \leq c(1 + |\theta|). \quad (42)$$

All along the proof,  $C > 0$  is a constant independent of  $N \geq 1$ ,  $T > 0$ ,  $i \in \{1, \dots, N\}$ ,  $1 \leq k \leq \lfloor NT \rfloor$ ,  $(x, y) \in \mathbf{X} \times \mathbf{Y}$ ,  $\theta \in \mathbf{R}^{d+1}$ , and  $z \in \mathbf{R}^d$ , which can change from one occurrence to another. It holds:

$$|\theta_k^i| \leq |\theta_0^i| + \sum_{\ell=0}^{k-1} |\theta_{\ell+1}^i - \theta_\ell^i|. \quad (43)$$

Using (5), we have, for  $0 \leq \ell \leq k-1$ ,

$$\begin{aligned} |\theta_{\ell+1}^i - \theta_\ell^i| &\leq \frac{\eta}{N^2} \sum_{j=1, j \neq i}^N \left| \langle \phi(\theta_\ell^j, \cdot, x_\ell), \gamma \rangle - y_\ell \langle \nabla_\theta \phi(\theta_\ell^i, \cdot, x_\ell), \gamma \rangle \right| \\ &\quad + \frac{\eta}{N^2} \left| \langle \phi(\theta_\ell^i, \cdot, x_\ell) - y_\ell \nabla_\theta \phi(\theta_\ell^i, \cdot, x_\ell), \gamma \rangle \right| + \frac{\eta}{N} |\nabla_\theta \mathcal{D}_{\text{KL}}(q_{\theta_\ell^i}^1 | P_0^1)|. \end{aligned} \quad (44)$$

For all  $\theta \in \mathbf{R}^{d+1}$ ,  $z \in \mathbf{R}^d$ ,  $(x, y) \in \mathbf{X} \times \mathbf{Y}$ , we have, by A2 and A3, since  $\phi(\theta, z, x) = s(\Psi_\theta(z), x)$ ,

$$|\phi(\theta, z, x) - y| \leq C. \quad (45)$$

Moreover, we have  $\nabla_\theta \phi(\theta, z, x) = \nabla_1 s(\Psi_\theta(z), x) J_\theta \Psi_\theta(z)$  (here  $\nabla_1 s$  refers to the gradient of  $s$  w.r.t. its first variable). By **A3**,  $|\nabla_1 s(\Psi_\theta(z), x)| \leq C$  and, hence, denoting by  $J_\theta$  the Jacobian w.r.t.  $\theta$ , using (10),

$$|\nabla_\theta \phi(\theta, z, x)| \leq C |J_\theta \Psi_\theta(z)| \leq C \mathfrak{b}(z). \quad (46)$$

Therefore, by (10),

$$\langle |\nabla_\theta \phi(\theta, \cdot, x)|, \gamma \rangle \leq C. \quad (47)$$

Hence, we obtain, using (44) and (42),

$$|\theta_{\ell+1}^i - \theta_\ell^i| \leq \frac{\eta}{N^2} \sum_{j=1, j \neq i}^N C + \frac{\eta}{N^2} C + \frac{c\eta}{N} (1 + |\theta_\ell^i|) \leq \frac{C}{N} (1 + |\theta_\ell^i|). \quad (48)$$

Using **A4**, there exists  $K_0 > 0$  such that a.s. for all  $i$ ,  $|\theta_0^i| \leq K_0$ . Then, from (43) and (48), for  $1 \leq k \leq \lfloor NT \rfloor$ , it holds:

$$|\theta_k^i| \leq K_0 + \frac{C}{N} \sum_{\ell=0}^{k-1} (1 + |\theta_\ell^i|) \leq K_0 + CT + \frac{C}{N} \sum_{\ell=0}^{k-1} |\theta_\ell^i| \leq C_{0,T} (1 + \frac{1}{N} \sum_{\ell=0}^{k-1} |\theta_\ell^i|),$$

with  $C_{0,T} = \max(K_0 + CT, C) \leq K_0 + C(1+T)$ . Then, by Lemma 15 and **A4**, we have that for all  $N \geq 1$ ,  $i \in \{1, \dots, N\}$  and  $0 \leq k \leq \lfloor NT \rfloor$ ,  $|\theta_k^i| \leq K_0 e^{[K_0 + C(1+T)]T}$ . The proof of Lemma 1 is thus complete.  $\blacksquare$

## Appendix B. Proof of Theorem 3

In this section, we assume **A1**  $\rightarrow$  **A5** (where in **A2**, when  $k \geq 1$ ,  $\mathcal{F}_k^N$  is now the one defined in (12)) and the  $\theta_k^i$ 's (resp.  $\mu^N$ ) are those defined by (7) for  $i \in \{1, \dots, N\}$  and  $k \geq 0$  (resp. by (13) for  $N \geq 1$ ).

### B.1. Preliminary analysis and pre-limit equation

#### B.1.1. NOTATION AND WEIGHTED SOBOLEV EMBEDDINGS

For  $J \in \mathbf{N}$  and  $\beta \geq 0$ , let  $\mathcal{H}^{J,\beta}(\mathbf{R}^{d+1})$  be the closure of the set  $\mathcal{C}_c^\infty(\mathbf{R}^{d+1})$  for the norm

$$\|f\|_{\mathcal{H}^{J,\beta}} := \left( \sum_{|k| \leq J} \int_{\mathbf{R}^{d+1}} \frac{|\partial_k f(\theta)|^2}{1 + |\theta|^{2\beta}} d\theta \right)^{1/2}.$$

The space  $\mathcal{H}^{J,\beta}(\mathbf{R}^{d+1})$  is a separable Hilbert space and we denote its dual space by  $\mathcal{H}^{-J,\beta}(\mathbf{R}^{d+1})$  (see e.g. Fernandez and Méléard (1997); Jourdain and Méléard (1998)). The associated scalar product on  $\mathcal{H}^{J,\beta}(\mathbf{R}^{d+1})$  will be denoted by  $\langle \cdot, \cdot \rangle_{\mathcal{H}^{J,\beta}}$ . For  $\Phi \in \mathcal{H}^{-J,\beta}(\mathbf{R}^{d+1})$ , we use the notation

$$\langle f, \Phi \rangle_{J,\beta} = \Phi[f], \quad f \in \mathcal{H}^{J,\beta}(\mathbf{R}^{d+1}).$$

For ease of notation, and if no confusion is possible, we simply denote  $\langle f, \Phi \rangle_{J,\beta}$  by  $\langle f, \Phi \rangle$ . The set  $\mathcal{C}_0^{J,\beta}(\mathbf{R}^{d+1})$  (resp.  $\mathcal{C}^{J,\beta}(\mathbf{R}^{d+1})$ ) is defined as the space of functions  $f : \mathbf{R}^{d+1} \rightarrow \mathbf{R}$  with continuous partial derivatives up to order  $J \in \mathbf{N}$  such that

$$\text{for all } |k| \leq J, \quad \lim_{|\theta| \rightarrow \infty} \frac{|\partial_k f(\theta)|}{1 + |\theta|^\beta} = 0 \quad (\text{resp. } \sum_{|k| \leq J} \sup_{\theta \in \mathbf{R}^{d+1}} \frac{|\partial_k f(\theta)|}{1 + |\theta|^\beta} < +\infty).$$

The spaces  $\mathcal{C}^{J,\beta}(\mathbf{R}^{d+1})$  and  $\mathcal{C}_0^{J,\beta}(\mathbf{R}^{d+1})$  is endowed with the norm

$$\|f\|_{\mathcal{C}^{J,\beta}} := \sum_{|k| \leq J} \sup_{\theta \in \mathbf{R}^{d+1}} \frac{|\partial_k f(\theta)|}{1 + |\theta|^\beta}.$$

We note that

$$\theta \in \mathbf{R}^{d+1} \mapsto (1 - \chi(\theta))|\theta|^\alpha \in \mathcal{H}^{J,\beta}(\mathbf{R}^{d+1}) \text{ if } \beta - \alpha > (d+1)/2, \quad (49)$$

where  $\chi \in \mathcal{C}_c^\infty(\mathbf{R}^{d+1})$  equals 1 near 0. We recall that from (Fernandez and Méléard, 1997, Section 2), for  $m' > (d+1)/2$  and  $\alpha, j \geq 0$ ,  $\mathcal{H}^{m'+j,\alpha}(\mathbf{R}^{d+1}) \hookrightarrow \mathcal{C}_0^{j,\alpha}(\mathbf{R}^{d+1})$ . In the following, we consider  $\gamma_0, \gamma_1 \in \mathbf{R}$  and  $L_0 \in \mathbf{N}$  such that

$$\gamma_1 > \gamma_0 > \frac{d+1}{2} + 1 \text{ and } L_0 > \frac{d+1}{2} + 1.$$

We finally recall the following standard result.

**Proposition 16** *Let  $q > p \geq 1$  and  $C > 0$ . The set  $\mathcal{K}_C^q := \{\mu \in \mathcal{P}_p(\mathbf{R}^{d+1}), \int_{\mathbf{R}^{d+1}} |x|^q \mu(dx) \leq C\}$  is compact.*

### B.1.2. BOUND ON THE MOMENTS OF THE $\theta_k^i$ 'S

We have the following uniform bound in  $N \geq 1$  on the moments of the sequence  $\{\theta_k^i, i \in \{1, \dots, N\}\}_{k=0, \dots, \lfloor NT \rfloor}$  defined by (7).

**Lemma 17** *Assume **A1**  $\rightarrow$  **A5**. For all  $T > 0$  and  $p \geq 1$ , there exists  $C > 0$  such that for all  $N \geq 1$ ,  $i \in \{1, \dots, N\}$  and  $0 \leq k \leq \lfloor NT \rfloor$ ,*

$$\mathbf{E}[|\theta_k^i|^p] \leq C.$$

**Proof** Let  $p \geq 1$ . By **A4**,  $\mathbf{E}[|\theta_0^i|^p] \leq C_p$  for all  $i \in \{1, \dots, N\}$ . Let  $T > 0$ . In the following  $C > 0$  is a constant independent of  $N \geq 1$ ,  $i \in \{1, \dots, N\}$ , and  $1 \leq k \leq \lfloor NT \rfloor$ . Using (7), the fact that  $\phi$  is bounded,  $\mathbf{Y}$  is bounded, and (46), we have, for  $0 \leq n \leq k-1$ ,

$$\begin{aligned} |\theta_{n+1}^i - \theta_n^i| &\leq \frac{C}{N^2 B} \sum_{j=1}^N \sum_{\ell=1}^B \mathfrak{b}(Z_n^{i,\ell}) + \frac{C}{N} |\nabla_{\theta} \mathcal{D}_{\text{KL}}(q_{\theta_n^i}^1 | P_0^1)| \\ &\leq \frac{C}{NB} \sum_{\ell=1}^B (1 + \mathfrak{b}(Z_n^{i,\ell})) + \frac{C}{N} (1 + |\theta_n^i|), \end{aligned} \quad (50)$$

where we have also used (42) for the last inequality. Let us recall the following convexity inequality: for  $m, p \geq 1$  and  $x_1, \dots, x_p \in \mathbf{R}_+$ ,

$$\left( \sum_{n=1}^m x_n \right)^p \leq m^{p-1} \sum_{n=1}^m x_n^p. \quad (51)$$

Using (43), **A1** with  $q = p$ , and the fact that  $1 \leq k \leq \lfloor NT \rfloor$ , one has setting  $u_k = \mathbf{E}[|\theta_k^i|^p]$ ,  $u_k \leq C(1 + \frac{1}{N} \sum_{n=0}^{k-1} u_n)$ . The result then follows from Lemma 15.  $\blacksquare$



## B.1.3. PRE-LIMIT EQUATION

In this section, we derive the pre-limit equation for  $\mu^N$  defined by (13). For simplicity we will keep the same notations as those introduced in Section A.1.1, though these objects will now be defined with  $\theta_k^j$  set by (7), and on  $\mathcal{C}^{2,\gamma_1}(\mathbf{R}^{d+1})$ , for all integer  $k \geq 0$ , and all time  $t \geq 0$ . Let  $f \in \mathcal{C}^{2,\gamma_1}(\mathbf{R}^{d+1})$ . Then, set for  $k \geq 0$ ,

$$\begin{aligned} \mathbf{D}_k^N[f] &= -\frac{\eta}{N^3} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \int_{\mathbf{X} \times \mathbf{Y}} (\langle \phi(\theta_k^j, \cdot, x), \gamma \rangle - y) \langle \nabla_{\theta} f(\theta_k^i) \cdot \nabla_{\theta} \phi(\theta_k^i, \cdot, x), \gamma \rangle \pi(dx, dy) \\ &\quad - \frac{\eta}{N^2} \int_{\mathbf{X} \times \mathbf{Y}} \langle (\phi(\cdot, \cdot, x) - y) \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \nu_k^N \otimes \gamma \rangle \pi(dx, dy). \end{aligned}$$

Note that  $\mathbf{D}_k^N$  above is the one defined in (21) but now on  $\mathcal{C}^{2,\gamma_1}(\mathbf{R}^{d+1})$  and with  $\theta_k^i$  defined by (7). For  $k \geq 0$ , we set

$$\mathbf{M}_k^N[f] = -\frac{\eta}{N^3 B} \sum_{i,j=1}^N \sum_{\ell=1}^B (\phi(\theta_k^j, Z_k^{j,\ell}, x_k) - y_k) \nabla_{\theta} f(\theta_k^i) \cdot \nabla_{\theta} \phi(\theta_k^i, Z_k^{i,\ell}, x_k) - \mathbf{D}_k^N[f]. \quad (52)$$

By Lemma 17 together with (45) and (46),  $\mathbf{M}_k^N[f]$  is integrable. Also, using A5 and the fact that  $\theta_k^j$  is  $\mathcal{F}_k^N$ -measurable (see (12)),

$$\mathbf{E}[\mathbf{M}_k^N[f] | \mathcal{F}_k^N] = 0.$$

Set  $\mathbf{M}_t^N[f] = \sum_{k=0}^{\lfloor Nt \rfloor - 1} \mathbf{M}_k^N[f]$ ,  $t \geq 0$ . We now extend the definition of  $\mathbf{W}_t^N[f]$  and  $\mathbf{R}_k^N[f]$  in (27) and (19) to any time  $t \geq 0$ ,  $k \geq 0$ , and  $f \in \mathcal{C}^{2,\gamma_1}(\mathbf{R}^{d+1})$ , and with  $\theta_k^i$  set by (7). We then set

$$\mathbf{R}_t^N[f] = \sum_{k=0}^{\lfloor Nt \rfloor - 1} \mathbf{R}_k^N[f], \quad t \geq 0.$$

With the same algebraic computations as those made in Section A.1.1, one obtains the following pre-limit equation: for  $N \geq 1$ ,  $t \geq 0$ , and  $f \in \mathcal{C}^{2,\gamma_1}(\mathbf{R}^{d+1})$ ,

$$\begin{aligned} \langle f, \mu_t^N \rangle - \langle f, \mu_0^N \rangle &= -\eta \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \langle \phi(\cdot, \cdot, x) - y, \mu_s^N \otimes \gamma \rangle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \rangle \pi(dx, dy) ds \\ &\quad - \eta \int_0^t \langle \nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \mu_s^N \rangle ds \\ &\quad + \frac{\eta}{N} \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \langle \langle \phi(\cdot, \cdot, x) - y, \gamma \rangle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \gamma \rangle, \mu_s^N \rangle \pi(dx, dy) ds \\ &\quad - \frac{\eta}{N} \int_0^t \int_{\mathbf{X} \times \mathbf{Y}} \langle (\phi(\cdot, \cdot, x) - y) \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x), \mu_s^N \otimes \gamma \rangle \pi(dx, dy) ds \\ &\quad + \mathbf{M}_t^N[f] + \mathbf{W}_t^N[f] + \mathbf{R}_t^N[f]. \end{aligned} \quad (53)$$

We will now show that the sequence  $(\mu^N)_{N \geq 1}$  is relatively compact in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ .

## B.2. Relative compactness and convergence to the limit equation

### B.2.1. RELATIVE COMPACTNESS IN $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$

In this section we prove the following result.

**Proposition 18** *Assume A1→A5. Recall  $\gamma_0 > \frac{d+1}{2} + 1$ . Then, the sequence  $(\mu^N)_{N \geq 1}$  is relatively compact in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ .*

We start with the following lemma.

**Lemma 19** *Assume A1→A5. Then,  $\forall T > 0$  and  $f \in \mathcal{C}^{2,\gamma_1}(\mathbf{R}^{d+1})$ ,*

$$\sup_{N \geq 1} \mathbf{E} \left[ \sup_{t \in [0, T]} \langle f, \mu_t^N \rangle^2 \right] < +\infty.$$

**Proof** Let  $T > 0$ . In what follows,  $C > 0$  is a constant independent of  $f \in \mathcal{C}^{2,\gamma_1}(\mathbf{R}^{d+1})$ ,  $(s, t) \in [0, T]^2$ , and  $z \in \mathbf{R}^d$  which can change from one occurrence to another. We have by A4,  $\mathbf{E}[\langle f, \mu_0^N \rangle^2] \leq C \|f\|_{\mathcal{C}^{2,\gamma_1}}^2$ . By (53) and (45), it holds:

$$\begin{aligned} \sup_{t \in [0, T]} \langle f, \mu_t^N \rangle^2 &\leq C \left[ \|f\|_{\mathcal{C}^{2,\gamma_1}}^2 + \int_0^T \int_{\mathbf{X} \times \mathbf{Y}} \left| \left\langle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x) \rangle, \gamma \right\rangle, \mu_s^N \right|^2 \pi(dx, dy) ds \right. \\ &\quad \left. + \int_0^T \left| \langle \nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q^1 | P_0^1) \rangle, \mu_s^N \right|^2 ds \right. \\ &\quad \left. + \frac{1}{N^2} \int_0^T \int_{\mathbf{X} \times \mathbf{Y}} \left| \left\langle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x) \rangle, \gamma \right\rangle, \mu_s^N \right|^2 \pi(dx, dy) ds \right. \\ &\quad \left. + \sup_{t \in [0, T]} |\mathbf{M}_t^N[f]|^2 + \sup_{t \in [0, T]} |\mathbf{W}_t^N[f]|^2 + \sup_{t \in [0, T]} |\mathbf{R}_t^N[f]|^2 \right]. \end{aligned} \quad (54)$$

We have using (46), for  $s \in [0, T]$  and  $z \in \mathbf{R}^d$ ,

$$|\nabla_{\theta} f(\theta_{[Ns]}^i) \cdot \nabla_{\theta} \phi(\theta_{[Ns]}^i, z, x)| \leq C \|f\|_{\mathcal{C}^{1,\gamma_1}} \mathfrak{b}(z) (1 + |\theta_{[Ns]}^i|^{\gamma_1}). \quad (55)$$

Thus, using Lemma 17,

$$\mathbf{E} \left[ \left\langle \langle \nabla_{\theta} f \cdot \nabla_{\theta} \phi(\cdot, \cdot, x) \rangle, \gamma \right\rangle, \mu_s^N \right]^2 \leq C \|f\|_{\mathcal{C}^{1,\gamma_1}}^2. \quad (56)$$

Using (42), for  $s \in [0, T]$ , it holds:

$$|\nabla_{\theta} f(\theta_{[Ns]}^i) \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q_{\theta_{[Ns]}^i}^1 | P_0^1)| \leq C \|f\|_{\mathcal{C}^{1,\gamma_1}} (1 + |\theta_{[Ns]}^i|^{\gamma_1+1}). \quad (57)$$

Thus, using Lemma 17,

$$\mathbf{E} \left[ \left| \langle \nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q^1 | P_0^1) \rangle, \mu_s^N \right|^2 \right] \leq C \|f\|_{\mathcal{C}^{1,\gamma_1}}^2. \quad (58)$$

On the other hand, we have using (51):

$$\sup_{t \in [0, T]} |\mathbf{M}_t^N[f]|^2 \leq \lfloor NT \rfloor \sum_{k=0}^{\lfloor NT \rfloor - 1} |\mathbf{M}_k^N[f]|^2. \quad (59)$$

Recall (52). By (21), (51), **A1**, and (55), it holds:

$$|\mathbf{D}_k^N[f]|^2 \leq C \|f\|_{\mathcal{C}^{1,\gamma_1}}^2 \left[ \frac{1}{N^4} \sum_{i \neq j=1}^N (1 + |\theta_k^i|^{2\gamma_1}) + \frac{1}{N^4} (1 + \langle |\cdot|^{2\gamma_1}, \nu_k^N \rangle) \right] \leq \frac{C}{N^2} \|f\|_{\mathcal{C}^{1,\gamma_1}}^2 (1 + |\theta_k^i|^{2\gamma_1})$$

and

$$|\mathbf{M}_k^N[f]|^2 \leq \frac{C}{N^4 B} \sum_{i,j=1}^N \sum_{\ell=1}^B \|f\|_{\mathcal{C}^{1,\gamma_1}}^2 |\mathbf{b}(Z_k^{i,\ell})|^2 (1 + |\theta_{[Ns]}^i|^{2\gamma_1}) + |\mathbf{D}_k^N[f]|^2.$$

By Lemma 17 and **A1**, one deduces that

$$\mathbf{E}[|\mathbf{M}_k^N[f]|^2] \leq C \|f\|_{\mathcal{C}^{1,\gamma_1}}^2 / N^2. \quad (60)$$

Going back to (59), we then have  $\mathbf{E}[\sup_{t \in [0, T]} |\mathbf{M}_t^N[f]|^2] \leq C \|f\|_{\mathcal{C}^{1,\gamma_1}}^2$ . Using the same arguments as those used so far, one also deduces that for  $t \in [0, T]$

$$\begin{aligned} \sup_{t \in [0, T]} |\mathbf{W}_t^N[f]|^2 &\leq \frac{C \|f\|_{\mathcal{C}^{1,\gamma_1}}^2}{N^2} \sup_{t \in [0, T]} (1 + \langle |\cdot|^{\gamma_1+1}, \nu_{[Nt]}^N \rangle)^2 \\ &= \frac{C \|f\|_{\mathcal{C}^{1,\gamma_1}}^2}{N^2} \max_{0 \leq k \leq [NT]} (1 + \langle |\cdot|^{\gamma_1+1}, \nu_k^N \rangle)^2 \\ &\leq \frac{C \|f\|_{\mathcal{C}^{1,\gamma_1}}^2}{N^2} \sum_{k=0}^{[NT]} (1 + \langle |\cdot|^{\gamma_1+1}, \nu_k^N \rangle)^2. \end{aligned}$$

and thus

$$\mathbf{E} \left[ \sup_{t \in [0, T]} |\mathbf{W}_t^N[f]|^2 \right] \leq C \|f\|_{\mathcal{C}^{1,\gamma_1}}^2 / N. \quad (61)$$

Let us finally deal with the term involving  $\mathbf{R}_t^N[f]$ . One has using (51):

$$\sup_{t \in [0, T]} |\mathbf{R}_t^N[f]|^2 \leq [NT] \sum_{k=0}^{[NT]-1} |\mathbf{R}_k[f]|^2.$$

For  $0 \leq k \leq [NT] - 1$ , we have, from (19),

$$\begin{aligned} |\mathbf{R}_k^N[f]|^2 &\leq \frac{C \|f\|_{\mathcal{C}^{2,\gamma_1}}^2}{N} \sum_{i=1}^N |\theta_{k+1}^i - \theta_k^i|^4 (1 + |\hat{\theta}_k^i|^{\gamma_1})^2 \\ &\leq \frac{C \|f\|_{\mathcal{C}^{2,\gamma_1}}^2}{N} \sum_{i=1}^N |\theta_{k+1}^i - \theta_k^i|^4 (1 + |\theta_{k+1}^i|^{2\gamma_1} + |\theta_k^i|^{2\gamma_1}). \end{aligned}$$

Using (50),

$$|\theta_{k+1}^i - \theta_k^i|^4 \leq C \left[ \frac{1}{N^4} + \frac{|\theta_k^i|^4}{N^4} + \frac{1}{N^4 B} \sum_{\ell=1}^B |\mathbf{b}(Z_k^{i,\ell})|^4 \right].$$

By Lemma 17 and A1, it then holds  $\mathbf{E}[|\theta_{k+1}^i - \theta_k^i|^4(1 + |\theta_{k+1}^i|^{2\gamma_1} + |\theta_k^i|^{2\gamma_1})] \leq C/N^4$ . Hence, one deduces that

$$\mathbf{E}[\sup_{t \in [0, T]} |\mathbf{R}_t^N[f]|^2] \leq C\|f\|_{\mathcal{C}^{2, \gamma_1}}^2/N^2. \quad (62)$$

This ends the proof of Lemma 19.  $\blacksquare$

**Lemma 20 (Compact containment for  $(\mu^N)_{N \geq 1}$ )** Assume A1  $\rightarrow$  A5. Let  $0 < \epsilon < \gamma_1 - \gamma_0$ . For every  $T > 0$ ,

$$\sup_{N \geq 1} \mathbf{E} \left[ \sup_{t \in [0, T]} \int_{\mathbf{R}^{d+1}} |x|^{\gamma_0 + \epsilon} \mu_t^N(dx) \right] < +\infty. \quad (63)$$

**Proof** Apply Lemma 19 with  $f : \theta \mapsto (1 - \chi)|\theta|^{\gamma_0 + \epsilon} \in \mathcal{C}^{2, \gamma_1}(\mathbf{R}^{d+1})$ .  $\blacksquare$

**Lemma 21** Assume A1  $\rightarrow$  A5. Let  $T > 0$  and  $f \in \mathcal{C}^{2, \gamma_1}(\mathbf{R}^{d+1})$ . Then, there exists  $C > 0$  such that for all  $\delta > 0$  and  $0 \leq r < t \leq T$  such that  $t - r \leq \delta$ , one has for all  $N \geq 1$ ,

$$\mathbf{E}[|\langle f, \mu_t^N \rangle - \langle f, \mu_r^N \rangle|^2] \leq C(\delta^2 + \delta/N + 1/N).$$

**Proof** Using (53), Jensen's inequality, (45), (56), and (58), one has for  $f \in \mathcal{C}^{2, \gamma_1}(\mathbf{R}^{d+1})$ ,

$$\begin{aligned} \mathbf{E}[|\langle f, \mu_t^N \rangle - \langle f, \mu_r^N \rangle|^2] &\leq C \left[ (t - r)^2 (1 + 1/N^2) \|f\|_{\mathcal{C}^{1, \gamma_1}}^2 + \mathbf{E} \left[ \left| \sum_{k=\lfloor Nr \rfloor}^{\lfloor Nt \rfloor - 1} \mathbf{M}_k^N[f] \right|^2 \right] \right. \\ &\quad \left. + \mathbf{E} [|\mathbf{W}_t^N[f] - \mathbf{W}_r^N[f]|^2] + \mathbf{E} [|\mathbf{R}_t^N[f] - \mathbf{R}_r^N[f]|^2] \right]. \end{aligned} \quad (64)$$

We also have with the same arguments as those used just before (31)

$$\mathbf{E} \left[ \left| \sum_{k=\lfloor Nr \rfloor}^{\lfloor Nt \rfloor - 1} \mathbf{M}_k^N[f] \right|^2 \right] = \sum_{k=\lfloor Nr \rfloor}^{\lfloor Nt \rfloor - 1} \mathbf{E} [|\mathbf{M}_k^N[f]|^2].$$

Using in addition (60), one has  $\mathbf{E} \left[ \left| \sum_{k=\lfloor Nr \rfloor}^{\lfloor Nt \rfloor - 1} \mathbf{M}_k^N[f] \right|^2 \right] \leq C(N\delta + 1) \|f\|_{\mathcal{C}^{1, \gamma_1}}^2/N^2$ . Note that with this argument, we also deduce that

$$\mathbf{E} [|\mathbf{M}_t^N[f]|^2] \leq C \|f\|_{\mathcal{C}^{1, \gamma_1}}^2/N. \quad (65)$$

On the other hand, by (61) and (62), one has

$$\mathbf{E} [|\mathbf{W}_t^N[f] - \mathbf{W}_r^N[f]|^2] \leq C \|f\|_{\mathcal{C}^{1, \gamma_1}}^2/N \text{ and } \mathbf{E} [|\mathbf{R}_t^N[f] - \mathbf{R}_r^N[f]|^2] \leq C \|f\|_{\mathcal{C}^{2, \gamma_1}}^2/N^2.$$

One then plugs all the previous estimates in (64) to deduce the result of Lemma 21.  $\blacksquare$

We are now in position to prove Proposition 18.

**Proof** [Proof of Proposition 18] The proof consists in applying (Jakubowski, 1986, Theorem 4.6) with  $E = \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1})$  and  $\mathbb{F} = \{H_f, f \in \mathcal{C}_c^\infty(\mathbf{R}^{d+1})\}$  where

$$H_f : \nu \in \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}) \mapsto \langle f, \nu \rangle.$$

The set  $\mathbb{F}$  on  $\mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1})$  satisfies Conditions (Jakubowski, 1986, (3.1) and (3.2) in Theorem 3.1). Condition (4.8) there follows from Proposition 16, Lemma 20, and Markov's inequality. Let us now show (Jakubowski, 1986, Condition (4.9)) is verified, i.e. that for all  $f \in \mathcal{C}_c^\infty(\mathbf{R}^{d+1})$ , the family  $(\langle f, \mu^N \rangle)_{N \geq 1}$  is relatively compact in  $\mathcal{D}(\mathbf{R}_+, \mathbf{R})$ . To do this, it suffices to use Lemma 21 and (Descours et al., 2022, Proposition A.1) (with  $\mathcal{H}_1 = \mathcal{H}_2 = \mathbf{R}$  there). In conclusion, according to (Jakubowski, 1986, Theorem 4.6), the sequence  $(\mu^N)_{N \geq 1} \subset \mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$  is relatively compact.  $\blacksquare$

### B.2.2. LIMIT POINTS SATISFY THE LIMIT EQUATION (15)

For  $f \in \mathcal{C}^{1, \gamma_0-1}(\mathbf{R}^{d+1})$  and  $t \geq 0$ , we introduce for  $m \in \mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ ,

$$\begin{aligned} \Phi_t[f] : m \mapsto & \langle f, m_t \rangle - \langle f, \mu_0 \rangle \\ & + \eta \int_0^t \int_{X \times Y} \langle \phi(\cdot, \cdot, x) - y, m_s \otimes \gamma \rangle \langle \nabla_\theta f \cdot \nabla_\theta \phi(\cdot, \cdot, x), m_s \otimes \gamma \rangle \pi(dx, dy) ds \\ & + \eta \int_0^t \langle \nabla_\theta f \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), m_s \rangle ds. \end{aligned} \quad (66)$$

Note that  $\Phi_t[f]$  is the function  $\Lambda_t[f]$  previously defined in (35) for test functions  $f \in \mathcal{C}^{1, \gamma_0-1}(\mathbf{R}^{d+1})$  and for  $m \in \mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ .

**Lemma 22** *Assume A1  $\rightarrow$  A5. Let  $f \in \mathcal{C}^{1, \gamma_0-1}(\mathbf{R}^{d+1})$ . Then  $\Phi_t[f]$  is well defined. In addition, if a sequence  $(m^N)_{N \geq 1}$  converges to  $m$  in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ , then, for all continuity point  $t \geq 0$  of  $m$ , we have  $\Phi_t[f](m^N) \rightarrow \Phi_t[f](m)$ .*

**Proof** Using A1, and because  $Y$  is bounded and the function  $\phi$  is bounded,  $\mathcal{G}_1^{x,y} : \theta \mapsto \langle \phi(\theta, \cdot, x) - y, \gamma \rangle \in \mathcal{C}_b^\infty(\mathbf{R}^{d+1})$ . In addition, for all multi-index  $\alpha \in \mathbf{N}^{d+1}$ , there exists  $C > 0$ , for all  $x, y \in X \times Y$  and all  $\theta \in \mathbf{R}^{d+1}$ ,  $|\partial_\alpha \mathcal{G}_1^{x,y}(\theta)| \leq C$ . The same holds for the function  $\mathcal{G}_2^x : \theta \in \mathbf{R}^{d+1} \mapsto \langle \nabla_\theta \phi(\theta, \cdot, x), \gamma \rangle$ . Consequently,  $\theta \mapsto \nabla_\theta f(\theta) \cdot \mathcal{G}_2^x(\theta) \in \mathcal{C}^{0, \gamma_0-1}(\mathbf{R}^{d+1}) \hookrightarrow \mathcal{C}^{0, \gamma_0}(\mathbf{R}^{d+1})$ . Then, there exists  $C > 0$  independent of  $(x, y) \in X \times Y$  and  $s \in [0, t]$  such that

$$|\langle \mathcal{G}_1^{x,y}, m_s \rangle| \leq C,$$

and

$$|\langle \nabla_\theta f \cdot \mathcal{G}_2^x, m_s \rangle| \leq C \|f\|_{\mathcal{C}^{1, \gamma_0-1}} \langle 1 + |\cdot|^{\gamma_0}, m_s \rangle.$$

Finally, the function  $\theta \mapsto \nabla_\theta \mathcal{D}_{\text{KL}}(q_\theta^1 | P_0^1)$  is smooth (see (4)) and (42) extends to all its derivatives, i.e. for all multi-index  $\alpha \in \mathbf{N}^{d+1}$ , there exists  $c > 0$ , for all  $\theta \in \mathbf{R}^{d+1}$ ,

$$|\partial_\alpha \nabla_\theta \mathcal{D}_{\text{KL}}(q_\theta^1 | P_0^1)| \leq c(1 + |\theta|).$$

Thus,  $\nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q_{\theta}^1 | P_0^1) \in \mathcal{C}^{0,\gamma_0}(\mathbf{R}^{d+1})$  and for some  $C > 0$  independent of  $s \in [0, t]$

$$|\langle \nabla_{\theta} f \cdot \nabla_{\theta} \mathcal{D}_{\text{KL}}(q_{\theta}^1 | P_0^1), \mathbf{m}_s \rangle| \leq C \|f\|_{\mathcal{C}^{1,\gamma_0-1}} \langle 1 + |\cdot|^{\gamma_0}, \mathbf{m}_s \rangle.$$

Since in addition  $\sup_{s \in [0,t]} \langle 1 + |\cdot|^{\gamma_0}, \mathbf{m}_s \rangle < +\infty$  (since  $s \mapsto \langle 1 + |\cdot|^{\gamma_0}, \mathbf{m}_s \rangle \in \mathcal{D}(\mathbf{R}_+, \mathbf{R})$ ),  $\Phi_t[f]$  is well defined. To prove the continuity property of  $\Phi_t[f]$  it then suffices to use the previous upper bounds together similar arguments as those used in the proof of Lemma 12 (see also Descours et al. (2022)).  $\blacksquare$

**Proposition 23** *Assume A1  $\rightarrow$  A5. Let  $\mu^*$  be a limit point of  $(\mu^N)_{N \geq 1}$  in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ . Then,  $\mu^*$  satisfies a.s. Equation (15).*

**Proof** Let us consider  $f \in \mathcal{C}_c^{\infty}(\mathbf{R}^{d+1})$  and  $\mu^*$  be a limit point of  $(\mu^N)_{N \geq 1}$  in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ . Recall that by (Ethier and Kurtz, 2009, lemma 7.7 in Chapter 3), the complementary of the set

$$\mathcal{C}(\mu^*) = \{t \geq 0, \mathbf{P}(\mu_{t-}^* = \mu_t^*) = 1\}$$

is at most countable. Let  $t_* \in \mathcal{C}(\mu^*)$ . Then, by Lemma 22, one has that  $\mathbf{P}(\mu^* \in \text{D}(\Phi_{t_*}[f])) = 0$ . Thus, by the continuous mapping theorem, it holds

$$\Phi_{t_*}[f](\mu^N) \xrightarrow{\mathcal{D}} \Phi_{t_*}[f](\mu^*).$$

On the other hand, using (53) and the estimates (62), (61), (65), (56), and (58), it holds

$$\lim_{N \rightarrow \infty} \mathbf{E}[\Phi_{t_*}[f](\mu^N)] = 0.$$

Consequently, for all  $f \in \mathcal{C}_c^{\infty}(\mathbf{R}^{d+1})$  and  $t_* \in \mathcal{C}(\mu^*)$ , it holds a.s.  $\Phi_{t_*}[f](\mu^*) = 0$ . On the other hand, for all  $\psi \in \mathcal{C}_c^{\infty}(\mathbf{R}^{d+1})$ ,  $\mathbf{m} \in \mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ , and  $s \geq 0$ , the mappings

$$t \geq 0 \mapsto \Phi_t[\psi](\mathbf{m})$$

is right continuous, and

$$f \in \mathcal{H}^{L_0, \gamma_0-1}(\mathbf{R}^{d+1}) \mapsto \Phi_s[f](\mathbf{m})$$

is continuous (because  $\mathcal{H}^{L_0, \gamma_0-1}(\mathbf{R}^{d+1}) \hookrightarrow \mathcal{C}_0^{1, \gamma_0-1}(\mathbf{R}^{d+1})$ ). In addition,  $\mathcal{H}^{L_0, \gamma_0-1}(\mathbf{R}^{d+1})$  admits a dense and countable subset of elements in  $\mathcal{C}_c^{\infty}(\mathbf{R}^{d+1})$ . Moreover, there exists a countable subset  $\mathcal{T}_{\mu^*}$  of  $\mathcal{C}(\mu^*)$  such that for all  $t \geq 0$  and  $\epsilon > 0$ , there exists  $s \in \mathcal{T}_{\mu^*}$ ,  $s \in [t, t + \epsilon]$ . We prove this claim. Since  $\mathbb{R}_+$  is a metric space,  $\mathcal{C}(\mu^*)$  is separable and thus admits a dense subset  $\mathcal{O}_{\mu^*}$ . Since  $[t + \epsilon/4, t + 3\epsilon/4] \cap \mathcal{C}(\mu^*) \neq \emptyset$ , there exists  $u \in [t + \epsilon/4, t + 3\epsilon/4] \cap \mathcal{C}(\mu^*)$ . Consider now  $s \in \mathcal{O}_{\mu^*}$  such that  $|s - u| \leq \epsilon/4$ . It then holds  $t \leq s \leq t + \epsilon$ , proving the claim with  $\mathcal{T}_{\mu^*} = \mathcal{O}_{\mu^*}$ .

Hence, we have with a classical argument that a.s. for all  $f \in \mathcal{H}^{L_0, \gamma_0-1}(\mathbf{R}^{d+1})$  and  $t \geq 0$ ,  $\Lambda_t[f](\mu^*) = 0$ . Note also that  $\mathcal{C}_b^{\infty}(\mathbf{R}^{d+1}) \subset \mathcal{H}^{L_0, \gamma_0-1}(\mathbf{R}^{d+1})$  since  $2\gamma_0 > d + 1$ . This ends the proof of the proposition.  $\blacksquare$

### B.3. Uniqueness of the limit equation and end of the proof of Theorem 3

In this section, we prove that there is a unique solution to (15) in  $\mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$ . To this end, we first need to prove that every limit points of  $(\mu^N)_{N \geq 1}$  a.s. belongs to  $\mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$ .

B.3.1. LIMIT POINTS BELONG TO  $\mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$ 

**Proposition 24** *Assume **A1**→**A5**. Let  $\mu^* \in \mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$  be a limit point of  $(\mu^N)_{N \geq 1}$  in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ . Then, a.s.  $\mu^* \in \mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$ .*

**Proof** Note that since  $W_1 \leq W_{\gamma_0}$ ,  $\mu^{N'} \xrightarrow{\mathcal{D}} \mu^*$  also in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$ , along some subsequence  $N'$ . According to (Jacod and Shiryaev, 1987, Proposition 3.26 in Chapter VI),  $\mu^* \in \mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$  a.s. if for all  $T > 0$ ,  $\lim_{N \rightarrow +\infty} \mathbf{E}[\sup_{t \in [0, T]} W_1(\mu_{t-}^N, \mu_t^N)] = 0$ . Using (32), this is equivalent to prove that

$$\lim_{N \rightarrow +\infty} \mathbf{E} \left[ \sup_{t \in [0, T]} \sup_{\|f\|_{\text{Lip}} \leq 1} |\langle f, \mu_{t-}^N \rangle - \langle f, \mu_t^N \rangle| \right] = 0. \quad (67)$$

Let us consider  $T > 0$  and a Lipschitz function  $f : \mathbf{R}^{d+1} \rightarrow \mathbf{R}$  such that  $\|f\|_{\text{Lip}} \leq 1$ . We have  $\langle f, \mu_t^N \rangle = \langle f, \mu_0^N \rangle + \sum_{k=0}^{\lfloor Nt \rfloor - 1} \langle f, \nu_{k+1}^N \rangle - \langle f, \nu_k^N \rangle$  (with usual convention  $\sum_0^{-1} = 0$ ). Thus the discontinuity points of  $t \in [0, T] \mapsto \langle f, \mu_t^N \rangle$  lies exactly at  $\{1/N, 2/N, \dots, \lfloor NT \rfloor / N\}$  and

$$|\langle f, \mu_{t-}^N \rangle - \langle f, \mu_t^N \rangle| \leq \max_{k=0, \dots, \lfloor NT \rfloor - 1} |\langle f, \nu_{k+1}^N \rangle - \langle f, \nu_k^N \rangle|, \quad \forall t \in [0, T], f \text{ Lipschitz.} \quad (68)$$

Pick  $k = 0, \dots, \lfloor NT \rfloor - 1$ . We have by (50),

$$|\langle f, \nu_{k+1}^N \rangle - \langle f, \nu_k^N \rangle| \leq \frac{1}{N} \sum_{i=1}^N |\theta_{k+1}^i - \theta_k^i| \leq \frac{C}{N} \sum_{i=1}^N \left[ \frac{1}{NB} \sum_{\ell=1}^B (1 + \mathfrak{b}(Z_k^{i, \ell})) + \frac{1}{N} (1 + |\theta_k^i|) \right] =: d_k^N \quad (69)$$

Hence, it holds:

$$|d_k^N|^2 \leq \frac{C}{N} \sum_{i=1}^N \left[ \frac{1}{N^2 B} \sum_{\ell=1}^B (1 + \mathfrak{b}^2(Z_k^{i, \ell})) + \frac{1}{N^2} (1 + |\theta_k^i|^2) \right],$$

where thanks to Lemma 17 and **A1**, for all  $k = 0, \dots, \lfloor NT \rfloor - 1$ ,  $\mathbf{E}[|d_k^N|^2] \leq C/N^2$  for some  $C > 0$  independent of  $N \geq 1$  and  $k = 0, \dots, \lfloor NT \rfloor - 1$ . Thus, using (68) and (69),

$$\begin{aligned} \mathbf{E} \left[ \sup_{t \in [0, T]} \sup_{\|f\|_{\text{Lip}} \leq 1} |\langle f, \mu_{t-}^N \rangle - \langle f, \mu_t^N \rangle| \right] &\leq \mathbf{E} \left[ \sup_{\|f\|_{\text{Lip}} \leq 1} \max_{k=0, \dots, \lfloor NT \rfloor - 1} |\langle f, \nu_{k+1}^N \rangle - \langle f, \nu_k^N \rangle| \right] \\ &\leq \mathbf{E} \left[ \max_{k=0, \dots, \lfloor NT \rfloor - 1} d_k^N \right] \\ &\leq \mathbf{E} \left[ \sqrt{\sum_{k=0}^{\lfloor NT \rfloor - 1} |d_k^N|^2} \right] \\ &\leq \sqrt{\mathbf{E} \left[ \sum_{k=0}^{\lfloor NT \rfloor - 1} |d_k^N|^2 \right]} \leq \frac{C}{\sqrt{N}}. \end{aligned}$$

This concludes the proof of Proposition 24. ■

## B.3.2. UNIQUENESS OF THE SOLUTION TO (15)

**Proposition 25** *There is a unique solution  $\bar{\mu} \in \mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$  to (15).*

**Proof** First of all, the existence of a solution is provided by Propositions 18, 24 and 23. Let us now prove that there is a unique solution to (15) in  $\mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$ .

Recall the definition of  $\mathbf{v}[\mu]$  in (37). We claim that for all  $T > 0$  and all solution  $\bar{\mu} \in \mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$  of (15), there exists  $C > 0$  such that

$$|\mathbf{v}[\bar{\mu}_t](\theta) - \mathbf{v}[\bar{\mu}_s](\theta)| \leq C|t - s|, \text{ for all } 0 \leq s \leq t \leq T \text{ and } \theta \in \mathbf{R}^{d+1}. \quad (70)$$

The proof of item (70) is the same as the one made for Item 2 in Proposition 14 since it holds using (42) and (46), for all  $0 \leq s \leq t \leq T$  and  $z \in \mathbf{R}^d$ ,

$$\begin{aligned} \left| \int_s^t \langle \nabla_\theta \phi(\cdot, z, x) \cdot \nabla_\theta \mathcal{D}_{\text{KL}}(q^1 | P_0^1), \bar{\mu}_r \rangle dr \right| &\leq C \mathfrak{b}(z) \int_s^t \langle (1 + |\cdot|), \bar{\mu}_r \rangle dr \\ &\leq C \mathfrak{b}(z) \max_{r \in [0, T]} \langle (1 + |\cdot|), \bar{\mu}_r \rangle |t - s|. \end{aligned}$$

We now conclude the proof of Proposition 25. Item 1 in the proof of Proposition 14 and (70) imply that  $v(t, \theta) = \mathbf{v}[\bar{\mu}_t](\theta)$  is globally Lipschitz on  $[0, T] \times \mathbf{R}^{d+1}$ , for all  $T > 0$ , when  $\bar{\mu} \in \mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$  is a solution of (15). Since in addition a solution  $\bar{\mu}$  to (15) is a weak solution on  $\mathbf{R}_+$  to (38) in  $\mathcal{C}(\mathbf{R}_+, \mathcal{P}(\mathbf{R}^{d+1}))$ , it holds by (Villani, 2003, Theorem 5.34):

$$\forall t \geq 0, \bar{\mu}_t = \phi_t \# \mu_0, \quad (71)$$

where  $\phi_t$  is the flow generated by the vector field  $\mathbf{v}[\bar{\mu}_t](\theta)$  over  $\mathbf{R}^{d+1}$ . Together with Item 3 in the proof of Proposition 14 and using the same arguments as those used in Step 3 of the proof of (Descours et al., 2022, Proposition 2.14), two solutions agrees on each  $[0, T]$  for all  $T > 0$ . One then deduces the uniqueness of the solution to (11). The proof of Proposition 25 is complete. ■

We are now in position to end the proof of Theorem 3.

**Proof [Proof of Theorem 3]** By Proposition 18,  $(\mu^N)_{N \geq 1}$  is relatively compact in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ . Let  $\mu^1, \mu^2 \in \mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$  be two limit points of this sequence. By Proposition 24, a.s.  $\bar{\mu}^1, \bar{\mu}^2 \in \mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1}))$ . In addition, according to Proposition 23,  $\mu^1$  and  $\mu^2$  are a.s. solutions of (15). Denoting by  $\bar{\mu} \in \mathcal{C}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$  the unique solution to (15) (see Proposition 25), we have a.s.

$$\bar{\mu}^1 = \bar{\mu} \text{ and } \bar{\mu}^2 = \bar{\mu} \text{ in } \mathcal{C}(\mathbf{R}_+, \mathcal{P}_1(\mathbf{R}^{d+1})).$$

In particular  $\bar{\mu} \in \mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$  and  $\bar{\mu}^j = \bar{\mu}$  in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ ,  $j \in \{1, 2\}$ . As a consequence,  $\bar{\mu}$  is the unique limit point of  $(\mu^N)_{N \geq 1}$  in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$  and the whole sequence  $(\mu^N)_{N \geq 1}$  converges to  $\bar{\mu}$  in  $\mathcal{D}(\mathbf{R}_+, \mathcal{P}_{\gamma_0}(\mathbf{R}^{d+1}))$ . Since  $\bar{\mu}$  is deterministic, the convergence also holds in probability. The proof of Theorem 3 is complete. ■

Let us now prove Proposition 4.

**Proof [Proof of Proposition 4]** Any solution to (11) in  $\mathcal{C}([0, T], \mathcal{P}(\Theta_T))$  is a solution to (15) in  $\mathcal{C}([0, T], \mathcal{P}_1(\mathbf{R}^{d+1}))$ . The result follows from Proposition 25. ■