

# FastPart: Over-Parameterized Stochastic Gradient Descent for Sparse optimisation on Measures

Yohann De Castro<sup>1,4</sup>, Sébastien Gadat<sup>2,4</sup> and Clément Marteau<sup>3</sup>

<sup>1</sup>*École Centrale de Lyon, CNRS UMR 5208, Institut Camille Jordan, Écully, France.*

<sup>2</sup>*Université Toulouse 1 Capitole, Toulouse School of Economics, France.*

<sup>3</sup>*Univ. Claude Bernard, CNRS UMR 5208, Institut Camille Jordan, Villeurbanne, France.*

<sup>4</sup>*Institut Universitaire de France (IUF)*

version as of December 9, 2023

## Abstract

This paper presents a novel algorithm that leverages Stochastic Gradient Descent strategies in conjunction with Random Features to augment the scalability of Conic Particle Gradient Descent (CPGD) specifically tailored for solving sparse optimisation problems on measures. By formulating the CPGD steps within a variational framework, we provide rigorous mathematical proofs demonstrating the following key findings: (i) The total variation norms of the solution measures along the descent trajectory remain bounded, ensuring stability and preventing undesirable divergence; (ii) We establish a global convergence guarantee with a convergence rate of  $\mathcal{O}(\log(K)/\sqrt{K})$  over  $K$  iterations, showcasing the efficiency and effectiveness of our algorithm, (iii) Additionally, we analyze and establish local control over the first-order condition discrepancy, contributing to a deeper understanding of the algorithm's behavior and reliability in practical applications.

## 1 Introduction

### 1.1 Convex programming for sparse optimisation on measures

Convex optimisation on the space of measures has gained attention during the past decade, *e.g.*, Bach and Chizat [2021], Chizat and Bach [2018], Chizat [2022], De Castro et al. [2021], Poon et al. [2021], De Castro and Gamboa [2012], Candès and Fernandez-Granda [2014] and references therein. It is also a popular field of investigation to derive global optimisation methods (a.k.a. simulated annealing) through the embedding of  $\mathbb{R}^d$  into the space of measures [Bolte et al., 2023, Miclo, 2023].

At his core, it can be viewed as a fruitful way of expressing many non-convex signal processing and machine learning tasks into a convex one, where one searches for an element of a Hilbert space  $\mathbb{H}$  that can be described as a linear combination of a few, say  $\bar{s}$ , elements:

$$\bar{\mathbf{y}} = \sum_{j=1}^{\bar{s}} \bar{\omega}_j \varphi_{\bar{\mathbf{t}}_j}, \quad (1.1)$$

from a given parameterized set  $\{\varphi_{\mathbf{t}} : \mathbf{t} \in \mathcal{X}\}$  where  $\bar{\omega}_j \in \mathbb{R} \setminus \{0\}$  and  $\mathcal{X}$  a compact set of  $\mathbb{R}^d$ .

Given an empirical observation  $\mathbf{y}$ , we would like to find a sparse representation, akin to (1.1), that explains  $\mathbf{y}$  and for which the learnt parameters  $(\omega_j, \mathbf{t}_j)_{j=1}^s$  encode an output solution for which generalization properties can be proven. A common practice is to minimize:

$$(\omega_j, \mathbf{t}_j)_{j=1}^s \mapsto \left\| \mathbf{y} - \sum_{j=1}^s \omega_j \varphi_{\mathbf{t}_j} \right\|_{\mathbb{H}}^2, \quad (1.2)$$

which is a non-convex program. In the expression (1.2) above,  $s \geq 1$  is a tuning parameter quantifying the so called *sparsity* of the solution.

A substantial body of literature pertains to the minimization of Mean Squared Error (MSE) and aligns with the framework outlined herein. In this paper, we will expound upon our methodology in a generalized format that can effectively encompass a majority of fields. Notably, certain specific instances will be elaborated upon within this paper, including but not limited to sparse deconvolution [De Castro and Gamboa, 2012, Candès and Fernandez-Granda, 2014], infinitely wide neural networks [Bach and Chizat, 2021, Chizat and Bach, 2018], or Mixture Models [De Castro et al., 2021]. Our approach is deployed according to the following steps.

**Lifting on the space of signed measures** First, we lift Program (1.2) onto the space  $\mathcal{M}(\mathcal{X})$  of Radon measures with finite total variation norm on  $\mathcal{X}$ . Consider the positive definite kernel  $\mathbb{K}$  defined as the dot product  $\mathbb{K}(\mathbf{t}, \mathbf{t}') := \langle \varphi_{\mathbf{t}}, \varphi_{\mathbf{t}'} \rangle_{\mathbb{H}}$  for all  $\mathbf{t}, \mathbf{t}' \in \mathcal{X}$ , and assume that:

$$\text{the function } \mathbf{t} \in \mathcal{X} \mapsto \varphi_{\mathbf{t}} \in \mathbb{H} \text{ is continuous.} \quad (\mathbf{A}_0)$$

Consider the *kernel measure embedding* (we refer to Appendix A.2 for further details),

$$\Phi : \nu \in \mathcal{M}(\mathcal{X}) \mapsto \int_{\mathcal{X}} \varphi_{\mathbf{t}} d\nu(\mathbf{t}) \in \mathbb{H}. \quad (1.3)$$

It is proven in the appendix (see Lemma A.2) that  $\Phi$  is a bounded linear map under  $(\mathbf{A}_0)$ . We deduce that:

$$\nu \in \mathcal{M}(\mathcal{X}) \mapsto \|\mathbf{y} - \Phi(\nu)\|_{\mathbb{H}}^2, \quad (1.4)$$

is a convex function on  $\mathcal{M}(\mathcal{X})$ . Taking the set of discrete measures given by  $\nu = \sum_{j=1}^s \omega_j \delta_{\mathbf{t}_j}$ , where  $\delta_{\mathbf{t}}$  is the Dirac mass at point  $\mathbf{t}$ , we uncover the parametrisation (1.2). Hence, we have lifted a non-convex program on  $(\omega_j, \mathbf{t}_j)_{j=1}^s$  onto a convex program over a much larger space, the space of signed measures.

**Total variation norm regularization** The second step is regularization. One key parameter is  $s$ , the number of learnt parameters, that should be considered to estimate (an approximation of) the true function (1.1). In practice, it can be cumbersome to tune this parameter and it might be better to resort to regularization. One benefit of the lifting on the space of measures is that this can be simply done by the TV-norm. Inspired by  $L^1$ -regularization in (high-dimensional) inverse problems, we study the so-called *Beurling LASSO* [De Castro and Gamboa, 2012, Candès and Fernandez-Granda, 2014] referred to as BLASSO below, whose convex objective function is given by:

$$J(\nu) := \frac{1}{2} \|\mathbf{y} - \Phi(\nu)\|_{\mathbb{H}}^2 + \lambda \|\nu\|_{\text{TV}}, \quad (1.5)$$

where  $\lambda > 0$  is a tuning parameter. We denote by  $\mu^* \in \mathcal{M}(\mathcal{X})$  a solution to BLASSO

$$J(\mu^*) = \min_{\mu \in \mathcal{M}(\mathcal{X})} J(\mu). \quad (\mathcal{B})$$

The existence of a solution to the problem at hand is not manifestly evident. Seminal contributions in the field, as articulated in [Bredies and Pikkarainen, 2013, Proposition 3.1] and [Hofmann et al., 2007, Theorem 3.1], have shown the existence of solutions upon continuity prerequisites imposed on the operator  $\Phi$  or its pre-dual counterpart. Nevertheless, the arduousness associated with ascertaining the continuity and well-defined attributes of the operator  $\Phi$  as expounded in Equation (1.3), within a given framework, underscores the complexity involved. In contrast, Condition  $(\mathbf{A}_0)$  affords a more tractable means of validation. Remarkably, our contribution lies in presenting a - to the best of our knowledge - novel result displayed in Theorem 1.1 below (established in Appendix A.3) that demonstrates the existence of solutions for any convex optimisation problem formulated in the manner of Equation (1.6), subject to the condition of continuity stipulated in Equation  $(\mathbf{A}_0)$ .

**Theorem 1.1.** Let  $\mathbb{H}$  be separable Hilbert space and let  $\mathcal{X}$  be compact metric space. Consider the problem

$$\inf_{\mu \in \mathcal{M}(\mathcal{X})} \left\{ L(\Phi\mu) + \lambda \|\mu\|_{\text{TV}} \right\} \quad (1.6)$$

where  $L : \mathbf{h} \in \mathbb{H} \rightarrow L(\mathbf{h}) \in [0, \infty]$  is convex and lower semi-continuous. If  $(\mathbf{A}_0)$  holds then there exists a measure  $\mu^* \in \mathcal{M}(\mathcal{X})$  solution to (1.6). Furthermore, if  $L$  is strictly convex, then the vector  $\Phi(\mu^*) \in \mathbb{H}$  is unique (it does not depend on the choice of the solution  $\mu^*$ ).

**Remark 1.1.** By choosing  $L(\mathbf{h}) = (1/2)\|\mathbf{y} - \mathbf{h}\|_{\mathbb{H}}^2$  in Theorem 1.1, it is established that there exists a signed measure  $\mu^* \in \mathcal{M}(\mathcal{X})$  solution to BLASSO  $(\mathcal{B})$ .

Over the last decade, several investigations on the performances of the estimator associated to the solution of (1.6) have been proposed in several specific situations. This solution can be proven to be close, for some partial Wasserstein 2 distance [Poon et al., 2021], to the target measure  $\bar{\mu} = \sum_{j=1}^S \bar{\omega}_j \delta_{\bar{t}_j}$  involved in Equation (1.1) in some cases of interest (e.g., Mixture Models [De Castro et al., 2021] or sparse deconvolution [Poon et al., 2021, De Castro and Gamboa, 2012]) as soon as the support points  $\bar{t}_j$  of the target  $\bar{\mu}$  are sufficiently separated. Moreover, if the bounded linear map  $\Phi$  has finite rank  $m \geq 1$ , then there exists a solution to  $(\mathcal{B})$  with at most  $m$  atoms, as proven by [Boyer et al., 2019, Section 4].

In this contribution, we focus our attention on the practical implementation of the optimisation problem  $(\mathcal{B})$ . In particular, our aim is to provide a runnable and efficient algorithm, and to display associated theoretical guarantees.

## 1.2 Learning with over-parametrised non-convex objective functions

Solving  $(\mathcal{B})$  from a practical point of view is not an immediate task, due to the infinite dimensional nature of the target. In this context, the Sliding Frank-Wolfe algorithm (see, e.g., Denoyelle et al. [2019]) provides an answer to this question. In this paper, we focus instead on the convergence of a Stochastic and Random Feature version of the Conic Particle Gradient Descent (CPGD) [Chizat, 2022] towards a minimum of Program  $(\mathcal{B})$ .

Writing the weights  $\mathbb{W} := (\omega_1, \dots, \omega_p)$  and the positions  $\mathbb{T} := (\mathbf{t}_1, \dots, \mathbf{t}_p) \in \mathcal{X}^p$ , we consider a generic measure with  $p$  weighted particles by:

$$\nu(\mathbb{W}, \mathbb{T}) := \sum_{j=1}^p \varepsilon_j \omega_j \delta_{\mathbf{t}_j}, \quad (1.7)$$

where  $\omega_j > 0$  (resp.  $\varepsilon_j = \pm 1$ ) refers to the weight (resp. the sign) of the particle  $j$ . The signs are fixed along the descent while the positions  $\mathbb{T}$  and weights  $\mathbb{W}$  are updated each gradient step. By a symmetrization argument, see for instance [Chizat, 2022, Appendix A], we consider, without loss of generality, that  $\varepsilon = 1$ . It holds that minimizing  $(\mathcal{B})$  or minimizing  $J$  defined by:

$$J(\mu^*) = \min_{\mu \in \mathcal{M}(\mathcal{X})_+} J(\mu). \quad (\mathcal{B}_+)$$

are equivalent, in a sense made precise by [Chizat, 2022, Proposition A.1] for instance; where  $\mathcal{M}(\mathcal{X})_+$  is the set of nonnegative measures with finite TV-norm. The attentive reader can uncover the next results on  $\mathcal{M}(\mathcal{X})$  by replacing  $\omega_j$  by  $\varepsilon_j \omega_j$ . The gradient descent dynamics are the same and our results also holds in this latter case.

Our algorithm makes use of particles measures as a proxy for solving the problem  $(\mathcal{B})$ . To this end, we adapt the notation of objective functions and related quantities accordingly. Denoting  $\boldsymbol{\lambda} := (\lambda, \dots, \lambda)$ , the definition of  $\nu(\mathbb{W}, \mathbb{T})$  in (1.7) then yields

$$J(\nu(\mathbb{W}, \mathbb{T})) = \frac{1}{2} \left\| \mathbf{y} - \sum_{j=1}^p \omega_j \varphi_{\mathbf{t}_j} \right\|_{\mathbb{H}}^2 + \lambda \sum_{j=1}^p \omega_j := F(\mathbb{W}, \mathbb{T}) + \frac{1}{2} \|\mathbf{y}\|_{\mathbb{H}}^2,$$

where  $\mathbf{k}_{\mathbb{T}} := (\langle \mathbf{y}, \varphi_{\mathbf{t}_1} \rangle_{\mathbb{H}}, \dots, \langle \mathbf{y}, \varphi_{\mathbf{t}_p} \rangle_{\mathbb{H}}) \in \mathbb{R}^p$ ,  $\mathbb{K}_{\mathbb{T}}$  is a  $(p \times p)$  matrix with entries  $\mathbb{K}(\mathbf{t}_i, \mathbf{t}_j)$  defined by:

$$\mathbb{K}(\mathbf{t}_i, \mathbf{t}_j) = \langle \varphi_{\mathbf{t}_i}, \varphi_{\mathbf{t}_j} \rangle_{\mathbb{H}}, \quad (1.8)$$

and

$$F(\mathbf{W}, \mathbb{T}) := \langle \boldsymbol{\lambda} - \mathbf{k}_{\mathbb{T}}, \mathbf{W} \rangle + \frac{1}{2} \mathbf{W}^T \mathbf{K}_{\mathbb{T}} \mathbf{W}, \quad (1.9)$$

is equal to  $J(\nu(\mathbf{W}, \mathbb{T}))$  up to the additive constant term  $(1/2)\|\mathbf{y}\|_{\mathbb{H}}^2$  (that only depends on the observations and not on the parameters of the measure we are optimising on).

In contrast to the original problem presented in (B), the optimisation process now operates within a different domain. Instead of working within the space of measures  $\mathcal{M}(\mathcal{X})$  (or  $\mathcal{M}(\mathcal{X})_+$ ), it focuses on particle measures with a set of fixed size  $p$ . This shift in perspective involves optimising over both the positions  $\mathbb{T}$  and weights  $\mathbf{W}$ . Although this adjustment serves to simplify the model's complexity to some extent, it introduces certain computational challenges. Firstly, for each pair of parameters  $(\mathbf{W}, \mathbb{T})$ , the computation of  $F(\mathbf{W}, \mathbb{T})$  necessitates the evaluation of  $\mathbf{k}_{\mathbb{T}}$  and  $\mathbf{K}_{\mathbb{T}}$ . Depending on the structure of the Hilbert space  $\mathbb{H}$  and the associated scalar products, this computation can be time-consuming. The need to calculate these quantities at each iteration of a gradient descent algorithm can become problematic, especially when dealing with high-dimensional spaces ( $d$  being significant). Furthermore, considering a large number of particles during the optimisation process, which is often essential, can substantially escalate the computational burden. This computational overhead needs to be carefully managed and optimised to ensure the efficiency of the optimisation procedure. In this paper, we address these issues by presenting a novel algorithm that incorporates Stochastic Gradient Descent (SGD) iterations. To the best of our knowledge, this approach has not been previously explored in the context of sparse optimisation within the domain of measures. We rigorously examine the properties of this algorithm, conducting a comprehensive investigation from both theoretical and practical perspectives.

The paper is organised as follows. First, we describe in Section 2 the construction of the algorithm and related necessary assumptions. Section 3 is illustrated by the example of deconvolution for mixture models, an important topic in unsupervised learning. Some theoretical results describing the behaviour of the algorithm in terms of the number of iterations are presented in Section 4, while a numerical illustration on some toy examples are discussed in Section 5. Proofs and related technical results are gathered in Section 6 and in Appendices A, B and C.

## 2 Construction of a Stochastic Gradient Descent algorithm.

The primary objective of this contribution is to introduce a stochastic algorithm designed for tackling the optimisation Program (B), followed by an exploration of its underlying theoretical properties. Our approach initially stems from a deterministic algorithm, which is subsequently adapted into a stochastic variant to enhance computational efficiency. Considering Program (B) as an optimisation challenge, it can be effectively addressed through the application of gradient descent (GD) techniques within the domain of measures. A straightforward algorithm would entail performing a discretized gradient descent on the space of non-negative measures  $\mathcal{M}(\mathcal{X})_+$ . However, in the absence of a significant conceptual breakthrough pertaining to an efficient parameterization of the preceding iteration within  $\mathcal{M}(\mathcal{X})_+$  (or its dual spaces and Hilbert basis), we resort to emulating this gradient descent process through a collection of measures encoded with particles. This method of approximating optimisation over measures through particles finds its conceptual foundation in swarm optimisation approaches, as exemplified in recent works such as those of Bolte et al. [2023], Miclo [2023].

### 2.1 Mirror principled conic gradient descent

#### 2.1.1 The Mirror descent principle

We first introduce a basic ingredient related to optimisation problems on geometric spaces, that permits to adapt the evolution of an algorithm to some constrained sets where the problem is embedded, without using some non-smooth additional projection steps. Mirror Descent (MD below) originates from the pioneering work of Nemirovskij and Yudin [1983] and permits to naturally handle optimisation problems especially when the mirror/proximal mapping is explicit, which is indeed the case for a convex problem constrained on measures as  $\mathcal{M}(\mathcal{X})_+$  (see e.g. Lan et al. [2012], Bubeck et al. [2015]).

The MD approach has the nice feature to define a smooth evolution that lives inside the constrained set without adding some supplementary projection step and “pushes” the frontiers of  $\mathcal{M}(\mathcal{X})_+$  at an infinite distance from any point that lies strictly inside.

Consider a strongly convex function  $h$  on  $\mathbb{R}_+^p \times \mathcal{X}^p$ , we define the Bregman divergence associated to  $h$  as follows: for two pairs  $(\mathbb{W}_1, \mathbb{T}_1)$  and  $(\mathbb{W}_2, \mathbb{T}_2)$  in  $\mathbb{R}_+^p \times \mathcal{X}^p$ , we denote:

$$D_h((\mathbb{W}_1, \mathbb{T}_1), (\mathbb{W}_2, \mathbb{T}_2)) = h(\mathbb{W}_1, \mathbb{T}_1) - h(\mathbb{W}_2, \mathbb{T}_2) - \langle \nabla h(\mathbb{W}_2, \mathbb{T}_2), (\mathbb{W}_1, \mathbb{T}_1) - (\mathbb{W}_2, \mathbb{T}_2) \rangle. \quad (2.1)$$

The Bregman divergence  $D_h$  is then used to define the MD with the following variational characterisation

$$(\mathbb{W}^{k+1}, \mathbb{T}^{k+1}) = \arg \min_{(\mathbb{W}, \mathbb{T})} \left\{ \langle \nabla F(\mathbb{W}^k, \mathbb{T}^k), (\mathbb{W}, \mathbb{T}) - (\mathbb{W}^k, \mathbb{T}^k) \rangle + \frac{1}{\kappa} D_h((\mathbb{W}, \mathbb{T}), (\mathbb{W}^{k+1}, \mathbb{T}^{k+1})) \right\}, \quad (2.2)$$

where  $\kappa > 0$  is the gradient step size.

### 2.1.2 A conic descent

In this contribution, we will consider the entropy function Ent on  $\{\mathbb{R}_+\}^p$  as:

$$\text{Ent}(\mathbb{W}) := \sum_{j=1}^p \omega_j \log(\omega_j) \quad (2.3)$$

It induces the Bregman divergence defined on the set of positive weights as

$$D_{\text{Ent}}(\mathbb{W}^1, \mathbb{W}^2) = \sum_{j=1}^p \omega_j^1 \left( \frac{\omega_j^2}{\omega_j^1} - 1 - \log \frac{\omega_j^2}{\omega_j^1} \right). \quad (2.4)$$

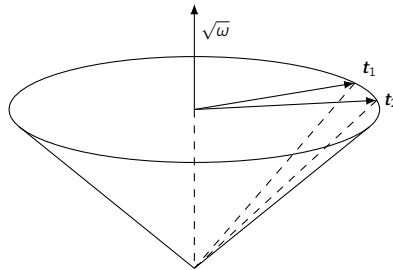
We then define the global divergence on the set  $\mathbb{R}_+^p \times \mathcal{X}^p$ , as

$$\Delta_{\alpha, \eta}((\mathbb{W}^1, \mathbb{T}^1), (\mathbb{W}^2, \mathbb{T}^2)) := \frac{1}{\alpha} D_{\text{Ent}}(\mathbb{W}^1, \mathbb{W}^2) + \frac{1}{\eta} D_{\text{Conic}}((\mathbb{W}^1, \mathbb{T}^1), (\mathbb{W}^2, \mathbb{T}^2)), \quad (2.5)$$

$$\text{where } D_{\text{Conic}}((\mathbb{W}^1, \mathbb{T}^1), (\mathbb{W}^2, \mathbb{T}^2)) := \sum_{j=1}^p \omega_j^1 \|\mathbf{t}_j^1 - \mathbf{t}_j^2\|^2. \quad (2.6)$$

In Equation (2.5) above, the parameter  $\alpha > 0$  is the gradient step size for the weights update and  $\eta > 0$  for the positions updates. This global divergence  $\Delta_{\alpha, \eta}$  is used to compute the gradient descent updates thanks to the variational formulation (2.2). We incorporate the term  $D_{\text{Conic}}$  with the specific intention of aligning it with the gradient updates associated with Conic Particle Gradient Descent (CPGD), as presented for instance by [Chizat, 2022, Section 2.2].

**Remark 2.1** (The conic metric...). We recall that in the CPGD framework [Chizat, 2022, Section 2.2], the set of particles  $(\omega, \mathbf{t}) \in \mathbb{R}_+ \times \mathcal{X}$  is equipped with the Riemannian metric defined by  $(1/2)\nabla_\omega^2 + (\omega/2)\|\nabla_{\mathbf{t}}\|^2$  where  $(\nabla_\omega, \nabla_{\mathbf{t}}) \in \mathbb{R} \times T_{(\omega, \mathbf{t})}(\mathcal{X})$  is a tangent vector at point  $(\omega, \mathbf{t})$ . As displayed below, we recognize a “conic” metric where two given fixed points  $\mathbf{t}_1, \mathbf{t}_2$  of  $\mathcal{X}$  gets linearly closer as  $\sqrt{\omega}$  goes to zero.



Then a mirror retraction is applied to  $\omega$  (see Chizat [2022][Definition 2.3]), corresponding to the Bregman divergence term  $D_{\text{Ent}}$  in our variational formulation. We notice that the term  $D_{\text{Conic}}$  comes from the latter conic metric.

**Remark 2.2** (...is not a Bregman divergence). *It is noteworthy that  $D_{\text{Conic}}$  does not conform to the definition of a Bregman divergence, as outlined in Definition (2.1). Specifically, there exists no global function  $h$  such that  $D_{\text{Conic}}$  can be expressed in the form of  $D_h$ . This is evident, for instance, by considering the boundedness of Bregman balls, a property that is conspicuously absent in the case of  $D_{\text{Conic}}$ . Furthermore, a direct proof of this assertion can be established through a contradiction argument. If such a function  $h$  were to exist for  $p = 1$ , it would imply the following relationship*

$$h(\omega_1, \mathbf{t}_1) - h(\omega_2, \mathbf{t}_2) - \langle \nabla h(\omega_2, \mathbf{t}_2), (\omega_1 - \omega_2, \mathbf{t}_1 - \mathbf{t}_2) \rangle = \omega_1 \|\mathbf{t}_1 - \mathbf{t}_2\|^2.$$

Then, consider  $\mathbf{t}_2 = 0$  and observe that necessarily  $h(\omega_1, \mathbf{t}_1) = h(0, 0) + \omega_1 \|\mathbf{t}_1\|^2$ , by removing the linear part that is necessarily vanishing. Subsequent straightforward calculations would lead to a contradiction, thereby confirming the incompatibility of  $D_{\text{Conic}}$  with the Bregman divergence framework.

According to the above remark, our descent algorithm should rather be understood as a Riemannian (stochastic) gradient descent instead of a purely mirror descent. We will keep the MD abuse of terms in what follows as it refers essentially to the evolution of the weights and since it is the commonly used term in machine learning with exponentially parameterized weights. Nevertheless, the difference between  $\Delta_{\alpha, \eta}$  and a Bregman divergence prevents the use of standard arguments of convergence for mirror descent algorithm.

## 2.2 Evaluation of the gradient and construction of a stochastic approximation

### 2.2.1 Algorithmic issues

Iterative algorithms based on (2.2) and (2.5) have already been at the core of theoretical investigations. We refer for instance to Chizat [2022] among others. The latter investigates an algorithm that requires some frequent calls to the gradient  $\nabla F$  of the objective  $F$  defined in (1.9). This gradient can be related to the Fréchet differential function  $\mathbf{t} \mapsto J'_\nu(\mathbf{t})$  of  $J(\cdot)$  at point  $\nu \in \mathcal{M}(\mathcal{X})_+$  and its gradient  $\nabla_{\mathbf{t}} J'_\nu$ . The Fréchet differential  $J'_\nu$  is defined through the following first order Taylor expansion:

$$\forall \nu \in \mathcal{M}(\mathcal{X})_+, \quad \nu + \sigma \in \mathcal{M}(\mathcal{X})_+ \quad J(\nu + \sigma) - J(\nu) = \langle J'_\nu, \sigma \rangle_{\mathcal{M}(\mathcal{X})^*, \mathcal{M}(\mathcal{X})} + q(\sigma), \quad (2.7)$$

where  $J'_\nu$  is the Fréchet gradient and  $q$  is a second order term. We refer to Proposition B.1 for details. According to Proposition B.3, for any  $\mathbf{t} \in \mathcal{X}$ , we have by Equation (B.3) that, given  $\nu = \sum_{j=1}^p \omega_j \delta_{\mathbf{t}_j}$ ,

$$J'_\nu(\mathbf{t}) = \sum_{j=1}^p \omega_j \langle \varphi_{\mathbf{t}}, \varphi_{\mathbf{t}_j} \rangle_{\mathbb{H}} - \langle \varphi_{\mathbf{t}}, \mathbf{y} \rangle_{\mathbb{H}} + \lambda. \quad (2.8)$$

The computation of these functions is time-consuming for the three reasons listed below. For each of these challenges, we outline our strategy to address them, and we introduce three notation—namely, the three random variables  $U$ ,  $V$ , and  $T$ —which will be elaborated upon in the subsequent section. In Section 3.1, we will provide a concrete example illustrating this phenomenon.

**Kernel evaluation: random variable  $U$  for Random Fourier feature strategy** Firstly, it is important to note that these functions may necessitate integral approximations due to their lack of closed-form expressions. For instance, the computation of  $J'_\nu(\mathbf{t})$  within each iteration of a gradient descent algorithm involves multiple evaluations of the kernel  $\mathbb{K}(\mathbf{t}, \mathbf{t}')$ , where  $\mathbb{K}(\mathbf{t}, \mathbf{t}')$  is defined in Equation (1.8). In many cases, these evaluations rely on non-explicit integrals, posing computational challenges.

To circumvent this issue, we will employ a Random Fourier feature strategy. This strategy aims to approximate the kernel  $\mathbb{K}$  by using a low-rank random kernel, achieved by evaluating the integral defining the kernel through independent Monte Carlo sampling.

**Large samples: random variable  $V$  for picking a data sample at random** A second strategy is to employ stochastic gradient computation with batch sub-sampling, which entails selecting a single data point at random - a fundamental component of various stochastic algorithms. Within our framework, this can be realized by utilizing the observation vector  $\mathbf{y}$ . It is worth emphasizing that in specific scenarios, the

observed signal  $\mathbf{y}$  can itself be a random variable. For instance, consider the case where  $\mathbf{y} := \Psi_N(\mathbf{X})$ , with  $\Psi_N : \mathbb{R}^N \rightarrow \mathbb{H}$  representing a bounded linear mapping and  $\mathbf{X} = (x_1, \dots, x_N)$  consisting of  $N$  i.i.d. samples drawn from a probability measure  $\rho$ .

In such cases, it becomes feasible to generate an unbiased stochastic version of  $\mathbf{y}$  by randomly selecting a single or a mini-batch data sample. Moreover, even if  $\mathbf{y}$  is deterministic or does not match the aforementioned identity  $\mathbf{y} := \Psi_N(\mathbf{X})$ , it is always possible to employ a similar strategy as described above (Random Fourier feature) to approximate  $\langle \varphi_t, \mathbf{y} \rangle_{\mathbb{H}}$  itself. Certainly, note that the unbiased stochastic version of  $\mathbf{y}$  and the low-rank random kernel approximation of  $\langle \varphi_t, \mathbf{y} \rangle_{\mathbb{H}}$  should be jointly considered whenever possible.

**Many particles: random variable  $T$  for picking a particle at random** The necessity for a large number of particles to attain convergence guarantees increases exponentially with the dimension  $d$  of  $\mathcal{X}$ . Consequently, it may become impractical to update all particles during each gradient step.

As a final ingredient, we opt to update the weight and position of a particle or a mini-batch of particles, selected randomly, as a mean to mitigate this challenge.

## 2.2.2 Stochastic approximation of the gradient

We provide in this paragraph a general framework allowing the construction of the stochastic conic gradient particle algorithm we are studying. Specific examples are discussed in Section 3 below.

The formulation of a stochastic approximation for the gradient  $\nabla F$  necessitates certain general assumptions, which are outlined below. Of particular significance is the introduction of a random variable  $Z = (T, U, V)$ , which serves as a way to alleviate the computational complexity associated with evaluating the variational formulation of the gradient descent step (2.2).

**Assumption (A<sub>1</sub>).** *There exists a pair of random variables  $(U, V)$  (not necessarily independent) such that, for any  $\mathbf{t}, \mathbf{t}' \in \mathcal{X}$ ,*

$$\langle \varphi_{\mathbf{t}}, \varphi_{\mathbf{t}'} \rangle_{\mathbb{H}} = \mathbb{E}_U \mathbf{g}_{\mathbf{t}, \mathbf{t}'}(U) \quad \text{and} \quad \langle \varphi_{\mathbf{t}}, \mathbf{y} \rangle_{\mathbb{H}} = \mathbb{E}_V \mathbf{h}_{\mathbf{t}}(V), \quad (\mathbf{A}_1)$$

for some explicit **bounded** functions  $\mathbf{g}$  and  $\mathbf{h}$ .

The latter assumption allows for a stochastic approximation of the functional  $J'_\nu$ . It exactly corresponds to what happens in Equation (3.11) below for the mixture model. Indeed, if we define a random variable  $T$  with distribution  $\nu/\nu(\mathcal{X})$  (assumed non-negative w.l.o.g as discussed in (1.7)), sampled independently from  $(U, V)$ , we then introduce:

$$J'_\nu(\mathbf{t}, Z) := \|\nu\|_{\text{TV}} \mathbf{g}_{\mathbf{t}, T}(U) - \mathbf{h}_{\mathbf{t}}(V) + \lambda \quad \text{where} \quad Z := (T, U, V). \quad (2.9)$$

According to (A<sub>1</sub>), we can write that

$$J'_\nu(\mathbf{t}, Z) := J'_\nu(\mathbf{t}) + \xi_\nu(\mathbf{t}, Z) \quad \text{with} \quad \mathbb{E}_Z \xi_\nu(\mathbf{t}, Z) = 0 \quad \forall \mathbf{t} \in \mathcal{X}. \quad (2.10)$$

We can construct a similar stochastic approximation for the gradient (w.r.t. the position parameter) of the functional  $J'_\nu$ . First remark that

$$\nabla_{\mathbf{t}} J'_\nu(\mathbf{t}) = \sum_{j=1}^p \omega_j \nabla_{\mathbf{t}} \langle \varphi_{\mathbf{t}}, \varphi_{\mathbf{t}_j} \rangle_{\mathbb{H}} - \nabla_{\mathbf{t}} \langle \varphi_{\mathbf{t}}, \mathbf{y} \rangle_{\mathbb{H}}.$$

The following assumption allows the commutativity between derivation and expectation and is satisfied in many situations, including batch or mini-batch strategies and smooth integral computations. The associated term in the mixture model is (3.12) and its stochastic counterpart is (3.15).

**Assumption (A<sub>2</sub>).** *The couple  $(U, V)$  and the functions  $\mathbf{g}, \mathbf{h}$  introduced in Assumption (A<sub>1</sub>) satisfy*

$$\nabla_{\mathbf{t}} \mathbb{E}_U \mathbf{g}_{\mathbf{t}, \mathbf{t}'}(U) = \mathbb{E}_U \nabla_{\mathbf{t}} \mathbf{g}_{\mathbf{t}, \mathbf{t}'}(U) \quad \text{and} \quad \nabla_{\mathbf{t}} \mathbb{E}_V \mathbf{h}_{\mathbf{t}}(V) = \mathbb{E}_V \nabla_{\mathbf{t}} \mathbf{h}_{\mathbf{t}}(V), \quad (\mathbf{A}_2)$$

for any  $\mathbf{t}, \mathbf{t}' \in \mathcal{X}$  and  $\mathbf{g}$  and  $\mathbf{h}$  have **bounded** derivatives.

---

**Algorithm 1** Stochastic & Random Feature Conic Particle Gradient Descent (FastPart)

---

**Require:** Constant learning rates  $(\alpha, \eta)$ , Initialization  $(\mathbf{W}^0, \mathbb{T}^0)$ ; ▷ Weights:  $\mathbf{W}^k$  and Positions:  $\mathbb{T}^k$   
1: **for**  $k = 1, \dots, K$  **do** ▷  $K$  gradient steps  
2:   Set  $\nu_k \leftarrow \nu(\mathbf{W}^k, \mathbb{T}^k)$ ; ▷ Particles  
3:   Sample  $Z^{k+1} \leftarrow (T^{k+1}, U^{k+1}, V^{k+1})$ ; ▷ Stochastic variables  
4:   Compute  $J'_{\nu_k}(\mathbb{T}^k, Z^{k+1})$  as defined by (2.9); ▷ Stochastic weights gradient  
5:   Compute  $D_{\nu_k}(\mathbb{T}^k, Z^{k+1})$  as defined by (2.11); ▷ Stochastic positions gradient  
6:   Compute 
$$d^{k+1} \leftarrow (J'_{\nu_k}(\mathbb{T}^k, Z^{k+1}), \mathbf{W}^k \odot D_{\nu_k}(\mathbb{T}^k, Z^{k+1}))$$
  
    where  $u \odot v$  denotes the vector  $(u_i v_i)_i$ ; ▷ Stochastic Conic Particle Gradient (see Proposition B.3)  
7:   Update the weights and positions with 
$$(\mathbf{W}^{k+1}, \mathbb{T}^{k+1}) := \arg \min_{(\mathbf{W}, \mathbb{T})} \langle d^{k+1}, (\mathbf{W}, \mathbb{T}) - (\mathbf{W}^k, \mathbb{T}^k) \rangle + \Delta_{\alpha, \eta}((\mathbf{W}, \mathbb{T}), (\mathbf{W}^k, \mathbb{T}^k)), \quad (2.13)$$
  
    with  $\Delta_{\alpha, \eta}$  given by (2.5). ▷ Mirror descent  
8: **end for**

---

Then, introduce

$$D_{\nu}(\mathbf{t}, Z) := \|\nu\|_{\text{TV}} \nabla_{\mathbf{t}} \mathbf{g}_{\mathbf{t}, T}(U) - \nabla_{\mathbf{t}} \mathbf{h}_{\mathbf{t}}(V), \quad (2.11)$$

where  $Z = (T, U, V)$  still denotes the same variable as above. We can then observe that

$$D_{\nu}(\mathbf{t}, Z) := \nabla_{\mathbf{t}} J'_{\nu}(\mathbf{t}) + \zeta_{\nu}(\mathbf{t}, Z) \quad \text{with} \quad \mathbb{E}_Z \zeta_{\nu}(\mathbf{t}, Z) = 0 \quad \forall \mathbf{t} \in \mathcal{X}, \quad (2.12)$$

where we have used Assumptions **(A<sub>1</sub>)** and **(A<sub>2</sub>)**.

Lemma B.1 provides a concrete and uniform upper bound on  $J'_{\nu}$  for any  $\nu \in \mathcal{M}(\mathcal{X})_+$ . From this remark, the boundedness of the functions  $\mathbf{g}$  and  $\mathbf{f}$  is indeed a reasonable assumption within this context.

Furthermore, in addition to the conclusions presented in Equations (2.10) and (2.12), that establish unbiased estimations of  $J'_{\nu}$  and  $\nabla_{\mathbf{t}} J'_{\nu}$ , it is imperative to highlight that our assumptions **(A<sub>1</sub>)** and **(A<sub>2</sub>)** also result in almost sure upper bounds on  $|J'_{\nu}(\cdot, Z)|$  and  $\|D_{\nu}(\cdot, Z)\|$  (for more comprehensive details, we refer to Lemma B.1 and Proposition C.1).

### 2.2.3 Stochastic conic particle gradient descent (FastPart)

With all these essential components in place, we are now prepared to construct our algorithm. The fundamental concept behind this approach is to substitute the deterministic gradient  $\nabla F$  of  $F$  in the mirror descent (2.2) with its stochastic counterpart, as derived from the stochastic gradients on weights (2.9) and positions (2.11) under Assumptions **(A<sub>1</sub>)** and **(A<sub>2</sub>)**. In Algorithm 1, we introduce our stochastic optimisation method. This algorithm is defined through the utilization of stochastic, unbiased realizations, as outlined in the preceding assumptions.

**Mirror descent updates** An important remark for the tractability of the Stochastic CPGD is that Equation (2.13) may be made explicit and simply corresponds to the following updates of the weights:

$$\forall j \in \{1, \dots, p\}, \quad \omega_j^{k+1} = \omega_j^k e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})}. \quad (2.14)$$

In a similar way, it may be verified that the positions are updated as follows:

$$\mathbb{T}^{k+1} = \mathbb{T}^k - \eta D_{\nu_k}(\mathbb{T}^k, Z^{k+1}) \quad (2.15)$$

**Properties** In what follows, we will study the properties of Algorithm 1 when the number  $K$  of iterations becomes large, and will establish some convergence and reconstruction properties when the several parameters involved in the method are correctly tuned (learning rates  $\alpha$  and  $\eta$ , number of particles  $p$ , and number of iterations  $K$ ).



## 3 Some examples from Unsupervised learning and Signal processing

### 3.1 Mixture Models (M.M.)

#### 3.1.1 Introduction

For the sake of clarity, we discuss briefly in this section the specific case of statistical mixture models. They are a class of statistical models that can be used for various purposes such as inference, testing, and modeling, have garnered significant attention in recent years due to their versatility and simplicity. However, the estimation of mixture models remains a complex task, with many aspects of the process not yet fully understood. The Expectation-Maximization (E.M.) algorithm, introduced by [Dempster et al. \[1977\]](#), and its subsequent generalization to stochastic variants in [Delyon et al. \[1999\]](#), have played a crucial role in the development of M.M.. Notably, the E.M. algorithm has been reinterpreted on exponential families as a descent algorithm with a surrogate in [Kunstner et al. \[2021\]](#), which has led to a renewed interest in M.M. within the machine learning and optimisation communities. Our work is also related to preliminary experiments conducted in [De Castro et al. \[2021\]](#), which employed a deterministic version of particle gradient descent. These experiments demonstrated the potential of using such methods in the context of M.M., and have inspired further research in this area. While M.M. may appear straightforward at first glance, their estimation poses significant challenges. However, recent advances in the field, including the reinterpretation of the EM algorithm and the use of descent algorithms with surrogates, have reignited interest in these models and their potential applications.

In this setting, the data  $\mathbf{X} = (x_1, \dots, x_N)$  are i.i.d. random variables having a density  $\rho$  in  $\mathbb{R}^d$  (w.r.t. the Lebesgue measure) verifying:

$$\rho = \theta \star \bar{\mu} = \sum_{j=1}^{\bar{s}} \bar{\omega}_j \theta_{\bar{\mathbf{t}}_j} \text{ with } \bar{\mu} := \sum_{j=1}^{\bar{s}} \bar{\omega}_j \delta_{\bar{\mathbf{t}}_j},$$

where  $\bar{\mu}$  is an *unknown* mixing distribution,  $\theta$  is a *known* even density on  $\mathbb{R}^d$  and  $\theta_{\mathbf{t}}(\cdot) = \theta(\mathbf{t} - \cdot)$ . In the expression above, the symbol  $\star$  denotes the convolution product. The goal in this context is to recover the target  $\bar{\mu}$  and/or the corresponding weights  $\bar{\mathbf{W}} := (\bar{\omega}_1, \dots, \bar{\omega}_{\bar{s}})$  and positions  $\bar{\mathbf{T}} := (\bar{\mathbf{t}}_1, \dots, \bar{\mathbf{t}}_{\bar{s}})$ .

#### 3.1.2 Model specification

Following [De Castro et al. \[2021\]](#), we consider the Hilbert space  $\mathbb{H}$  defined as the RKHS associated to the  $\text{sinc}_m$  kernel, denoted by  $\gamma_m$ , for some bandwidth parameter  $m$  and leading to

$$\mathbb{H} := \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R} : \|f\|_{\mathbb{H}}^2 = \int_{\mathbb{R}^d} \frac{|\mathcal{F}[f](\mathbf{t})|^2}{\mathcal{F}[\gamma_m](\mathbf{t})} d\mathbf{t} < +\infty \right\} \text{ with } \mathcal{F}[\gamma_m] = \mathbf{1}_{[-m, m]^d},$$

where  $\mathcal{F}$  is the Fourier transform defined as

$$\forall \mathbf{s} \in \mathbb{R}^d, \quad \forall f \in \mathbb{H}, \quad \mathcal{F}[f](\mathbf{s}) := \int_{\mathbb{R}^d} f(\mathbf{t}) e^{-i\langle \mathbf{s}, \mathbf{t} \rangle} d\mathbf{t},$$

and  $i$  refers to the complex number. The inner product associated to the space  $\mathbb{H}$  hence verifies

$$\forall f, g \in \mathbb{H}, \quad \langle f, g \rangle_{\mathbb{H}} = \mathcal{R} \left( \int_{[-m, m]^d} \mathcal{F}[f](\mathbf{t}) \times \overline{\mathcal{F}[g](\mathbf{t})} d\mathbf{t} \right), \quad (3.1)$$

where  $\mathcal{R}(z)$  denotes the real part of a complex number  $z$ . As in [De Castro et al. \[2021\]](#), we embed the sample  $\mathbf{X}$  and the density  $\rho$  in the same Hilbert space  $\mathbb{H}$  taking a convolution by the  $\text{sinc}_m$  kernel. It yields

$$\mathbf{y} = \Psi_N(\mathbf{X}) := \frac{1}{N} \sum_{i=1}^N \gamma_m(x_i - \cdot) \text{ and } \bar{\mathbf{y}} := \mathbb{E}_{\mathbf{X}}[\mathbf{y}] = (\gamma_m \star \theta) \star \bar{\mu} = \sum_{j=1}^{\bar{s}} \bar{\omega}_j (\gamma_m \star \theta)(\bar{\mathbf{t}}_j - \cdot).$$

We uncover that the feature map (1.1) and the measure embedding (1.3) are given by

$$\varphi_{\mathbf{t}}(\cdot) := (\gamma_m \star \theta)(\mathbf{t} - \cdot) \text{ and } \Phi(\mu) := (\gamma_m \star \theta) \star \mu.$$

Note that the observation  $\mathbf{y} \in \mathbb{H}$  corresponds to the non-parametric kernel estimator of  $\bar{\mathbf{y}} = \Phi(\bar{\mu})$  based on the sample  $\mathbf{X}$ . For the sake of simplicity we omit the dependency with the bandwidth  $m$  as we are only interested in the optimisation problem.

### 3.1.3 Gradients

For any  $\mu \in \mathcal{M}(\mathbb{R}^d)$ , the observation  $\mathbf{y}$  is compared to  $\Phi(\mu) = (\gamma \star \theta) \star \mu$  inside the criterion  $J(\mu)$ . A direct application of Proposition B.1 leads to

$$J'_\nu = \Phi^*(\Phi(\nu) - \mathbf{y}) + \lambda = (\gamma \star \theta) \star ((\gamma \star \theta) \star \nu - \mathbf{y}) + \lambda,$$

One can check that, for any  $\mathbf{t} \in \mathbb{R}^d$ ,

$$J'_\nu(\mathbf{t}) = \int_{\mathbb{R}^d} \theta_t(\mathbf{s})(\Phi(\nu) - \mathbf{y})(\mathbf{s})d\mathbf{s} + \lambda = \langle \varphi_t, \Phi(\nu) - \mathbf{y} \rangle_{\mathbb{H}} + \lambda.$$

The latter formula can be re-written as

$$J'_\nu(\mathbf{t}) = \langle \varphi_t, \Phi(\nu) \rangle_{\mathbb{H}} - \langle \varphi_t, \mathbf{y} \rangle_{\mathbb{H}} + \lambda = \sum_{j=1}^p \omega_j \langle \varphi_t, \varphi_{t_j} \rangle_{\mathbb{H}} - \langle \varphi_t, \mathbf{y} \rangle_{\mathbb{H}} + \lambda, \quad (3.2)$$

In particular, if we denote  $\sigma = \mathcal{F}[\theta]$ , using (3.1) and the definition of  $\theta_t$  as a convolution, we obtain that

$$\forall \mathbf{t}, \mathbf{s} \in \mathbb{R}^d, \quad \mathcal{F}[\theta_t](\mathbf{s}) = \mathcal{F}[\theta](\mathbf{s})e^{-i\langle \mathbf{s}, \mathbf{t} \rangle} = \sigma(\mathbf{s})e^{-i\langle \mathbf{s}, \mathbf{t} \rangle},$$

where  $\sigma(\cdot)$  is a real valued function since  $\theta$  is even. A consequence is therefore that

$$\langle \varphi_t, \mathbf{y} \rangle_{\mathbb{H}} = \frac{1}{N} \sum_{i=1}^N \int_{[-m, m]^d} \cos(\langle \mathbf{u}, \mathbf{t} - \mathbf{x}_i \rangle) \sigma(\mathbf{u}) d\mathbf{u} \quad \text{and} \quad \langle \varphi_t, \varphi_{t_j} \rangle_{\mathbb{H}} = \int_{[-m, m]^d} \cos(\langle \mathbf{u}, \mathbf{t} - \mathbf{t}_j \rangle) \sigma^2(\mathbf{u}) d\mathbf{u}.$$

The previous computation then shows that:

$$J'_\nu(\mathbf{t}) = \sum_{j=1}^p \omega_j \int_{[-m, m]^d} \cos(\langle \mathbf{u}, \mathbf{t} - \mathbf{t}_j \rangle) \sigma^2(\mathbf{u}) d\mathbf{u} - \frac{1}{N} \sum_{i=1}^N \int_{[-m, m]^d} \cos(\langle \mathbf{u}, \mathbf{t} - \mathbf{x}_i \rangle) \sigma(\mathbf{u}) d\mathbf{u} + \lambda. \quad (3.3)$$

### 3.1.4 Assumptions $(\mathbf{A}_0)$ , $(\mathbf{A}_1)$ and $(\mathbf{A}_2)$

We can identify the  $\mathbf{g}$  and  $\mathbf{h}$  functions appearing in Assumptions  $(\mathbf{A}_1)$  and  $(\mathbf{A}_2)$ . It holds

$$\mathbf{g}_{\mathbf{t}, \mathbf{t}'}(\mathbf{u}) := \cos(\langle \mathbf{u}, \mathbf{t} - \mathbf{t}' \rangle) \sigma(\mathbf{u}) \mathbf{1}_{\{u \in [-m, m]^d\}} \quad (3.4)$$

$$\mathbf{h}_{\mathbf{t}}(\mathbf{u}, \mathbf{v}') := \cos(\langle \mathbf{u}, \mathbf{t} - \mathbf{v}' \rangle) \mathbf{1}_{\{u \in [-m, m]^d\}}. \quad (3.5)$$

Now, observe that  $\|\sigma\|_\infty \leq \|\theta\|_1 = 1$ , hence both functions are bounded and Assumption  $(\mathbf{A}_1)$  is satisfied. Using the dominated convergence theorem for the bounded gradients

$$\nabla_{\mathbf{t}} \mathbf{g}_{\mathbf{t}, \mathbf{t}'}(\mathbf{u}) = \mathbf{u} \sin(\langle \mathbf{u}, \mathbf{t} - \mathbf{t}' \rangle) \sigma(\mathbf{u}) \mathbf{1}_{\{u \in [-m, m]^d\}} \quad (3.6)$$

$$\nabla_{\mathbf{t}} \mathbf{h}_{\mathbf{t}}(\mathbf{u}, \mathbf{v}') = \mathbf{u} \sin(\langle \mathbf{u}, \mathbf{t} - \mathbf{v}' \rangle) \mathbf{1}_{\{u \in [-m, m]^d\}}, \quad (3.7)$$

one can check that Assumption  $(\mathbf{A}_2)$  is satisfied.

Using the dominated convergence theorem and  $\|\varphi_t - \varphi_s\|_{\mathbb{H}}^2 = \int_{[-m, m]^d} |e^{-i\langle \mathbf{u}, \mathbf{t} \rangle} - e^{-i\langle \mathbf{u}, \mathbf{s} \rangle}|^2 \sigma(\mathbf{u}) d\mathbf{u}$ , one can prove that Assumption  $(\mathbf{A}_0)$  is satisfied.

### 3.1.5 Stochastic counterpart

Remark from (3.3) that  $J'_\nu$  is obtained through some integral computations. Given  $\nu = \nu(\mathbb{W}, \mathbb{T})$ , we introduce a random variable  $Z = (T, U, V)$  built as follows. Consider three independent random variables given by

$$T \sim \frac{\nu}{\|\nu\|_{\text{TV}}}, \quad V' \sim \frac{1}{N} \sum_{i=1}^N \delta_{x_i}, \quad \text{and } U \sim \sigma, \quad (3.8)$$

where  $\nu$  is assumed nonnegative wlog as discussed in (1.7), hence  $\nu/\|\nu\|_{\text{TV}}$  is a discrete probability measure. We then define  $V := (U, V')$  and

$$J'_{\nu, \mathbf{t}}(Z) = \|\nu\|_{\text{TV}} \cos(\langle U, \mathbf{t} - T \rangle) \sigma(U) \mathbf{1}_{\{U \in [-m, m]^d\}} - \cos(\langle U, \mathbf{t} - V' \rangle) \mathbf{1}_{\{U \in [-m, m]^d\}} + \lambda, \quad (3.9)$$

$$= \|\nu\|_{\text{TV}} \mathbf{g}_{\mathbf{t}, T}(U) - \mathbf{h}_{\mathbf{t}}(V) + \lambda. \quad (3.10)$$

and we uncover (2.9). According to (3.2) and (3.8), we can verify that once  $\nu$  is fixed,  $J'_{\nu, \mathbf{t}}(Z)$  is an unbiased estimation of  $J'_\nu$  and we can observe that there exists a random variable  $\xi_\nu(\mathbf{t}, Z)$  such that

$$J'_{\nu, \mathbf{t}}(Z) = J'_\nu(\mathbf{t}) + \xi_\nu(\mathbf{t}, Z) \quad \text{with} \quad \mathbb{E}_Z[\xi_\nu(\mathbf{t}, Z)] = 0. \quad (3.11)$$

A same remark occurs for the gradient of  $J'_\nu$  since for any  $\mathbf{t} \in \mathbb{R}^d$ ,

$$\nabla_{\mathbf{t}} J'_\nu(\mathbf{t}) = \sum_{j=1}^p \omega_j \int_{[-m, m]^d} \mathbf{u} \sin(\langle \mathbf{u}, \mathbf{t} - \mathbf{t}_j \rangle) \sigma^2(\mathbf{u}) d\mathbf{u} - \frac{1}{N} \sum_{i=1}^N \int_{[-m, m]^d} \mathbf{u} \sin(\langle \mathbf{u}, \mathbf{t} - x_i \rangle) \sigma(\mathbf{u}) d\mathbf{u}. \quad (3.12)$$

We shall then define

$$\mathbf{D}_{\nu, \mathbf{t}}(Z) = \|\nu\|_{\text{TV}} U \sin(\langle U, \mathbf{t} - T \rangle) \sigma(U) \mathbf{1}_{\{U \in [-m, m]^d\}} - U \sin(\langle U, \mathbf{t} - V' \rangle) \mathbf{1}_{\{U \in [-m, m]^d\}}, \quad (3.13)$$

$$= \|\nu\|_{\text{TV}} \nabla_{\mathbf{t}} \mathbf{g}_{\mathbf{t}, T}(U) - \nabla_{\mathbf{t}} \mathbf{h}_{\mathbf{t}}(V). \quad (3.14)$$

and we uncover (2.11). We get that

$$\mathbf{D}_{\nu, \mathbf{t}}(Z) = \nabla_{\mathbf{t}} J'_\nu(\mathbf{t}) + \zeta_\nu(\mathbf{t}, Z) \quad \text{with} \quad \mathbb{E}_Z[\zeta_\nu(\mathbf{t}, Z)] = 0, \quad (3.15)$$

for some centered random vector  $\zeta_\nu(\mathbf{t}, Z)$ .

Therefore, in this example of mixture deconvolution, our situation perfectly fits the stochastic gradient setting where we can easily access some unbiased random realization of the true gradient of  $F$  for any position of a current algorithm  $\nu = \nu(\mathbb{W}, \mathbb{T})$ .

## 3.2 Interlude, a comparison with sketching

In the field of machine learning, replacing the computation of the complex integrals present in Equation (3.3) is commonly referred to as sketching [Keriven et al., 2018]. The key distinction between our approach and sketching-type methods lies in how the frequencies ( $U_k$ ) are sampled. In the sketching approach, the frequencies are initially sampled only once at the beginning of the algorithm. Consequently, the gradients  $J'_\nu(\mathbf{t})$  and  $\nabla_{\mathbf{t}} J'_\nu(\mathbf{t})$  are approximated using a Monte-Carlo version of (3.3) that employs the same Monte-Carlo sample ( $U_k$ ) along the descent and the size of  $U_k$  needs to be large to guarantee a good Monte-Carlo approximation of the several integrals involved in the iterations of the method. Therefore, the sketching strategy of [Keriven et al., 2018] leads to costly iterations as the size of ( $U_k$ ) is not negligible. In contrast, our method involves sampling a unique Monte-Carlo sample  $U_k$  at each step  $k$ , and we replace the integrals of (3.3) with an evaluation at sample  $U_k$  (or we could also adopt a mini-batch strategy). It is important to note that in our case, the number of frequencies ( $U_k$ ) is equal to the number of steps  $K$  (or times the size of mini-batch), whereas in sketching approaches, the number of frequencies is directly proportional to the number of parameters to be estimated, up to logarithmic factors. For more information on this topic, refer to Keriven et al. [2018].

### 3.3 Sparse deconvolution with positive definite kernel

#### 3.3.1 Introduction

The statistical analysis of the  $\ell_1$ -regularization in the space of measures was initiated by Donoho [Donoho, 1992] and then investigated by Gamboa and Gassiat [1996]. Recently, this problem has attracted a lot of attention in the “*Super-Resolution*” community and its companion formulation in “*Line spectral estimation*”. In the Super-Resolution frame, one aims at recovering fine scale details of an image from few low frequency measurements, ideally the observation is given by a low-pass filter. The novelty of this body of work lies in new theoretical guarantees of the  $\ell_1$ -minimization over the space of discrete measures in a gridless manner, referred to as “off-the-grid” methods. Some recent work on this topic can be found in Bredies and Pikkarainen [2013], Tang et al. [2013], Candès and Fernandez-Granda [2014, 2013], Fernandez-Granda [2013], Duval and Peyré [2015], De Castro and Gamboa [2012], Azais et al. [2015]. More precisely, pioneering work can be found in Bredies and Pikkarainen [2013], which treats of inverse problems on the space of Radon measures and Candès and Fernandez-Granda [2014], which investigates the Super-Resolution problem via Semi-Definite Programming and the ground breaking construction of a “*dual certificate*”. Exact Reconstruction property (in the noiseless case), minimax prediction and localization (in the noisy case) have been performed using the “*Beurling Lasso*” estimator ( $\mathcal{B}$ ) introduced in Azais et al. [2015] and also studied in Tang et al. [2013], Fernandez-Granda [2013], Tang et al. [2015] which minimizes the total variation norm over complex Borel measures. Noise robustness (as the noise level tends to zero) has been investigated in the captivating paper Duval and Peyré [2015]. A sketching formulation and a construction of dual certificates with respect to the Fisher metric has been pioneering studied in Poon et al. [2021].

#### 3.3.2 Model specification

In sparse deconvolution, the convolution  $\bar{\mathbf{y}}$  of some measure  $\bar{\mu}$  with a  $\mathcal{C}^1$ -continuous **positive definite** function  $\varphi$  is given by

$$\forall \mathbf{t} \in \mathcal{X}, \quad \bar{\mathbf{y}}(\mathbf{t}) = \sum_{j=1}^{\bar{s}} \bar{\omega}_j \varphi(\mathbf{t} - \bar{\mathbf{t}}_j),$$

and we observe  $\mathbf{y}$ , a noisy version of  $\bar{\mathbf{y}}$ , given by  $\mathbf{y} = \bar{\mathbf{y}} + \mathbf{e}$ , where  $\mathbf{e} : \mathcal{X} \rightarrow \mathbb{R}$  is some function. The feature map (1.1) is the convolution kernel and the measure embedding (1.3) are given by

$$\varphi_{\mathbf{t}}(\cdot) := \varphi(\mathbf{t} - \cdot) \text{ and } \Phi(\mu) := \varphi \star \mu.$$

Recall that  $\varphi$  is a positive definite function. We assume that  $\varphi$  is defined on the  $d$ -Torus  $\mathbb{T}^d$  or on  $\mathbb{R}^d$ , and we further assume that  $\varphi(0) = 1$ , without loss of generality on this latter point. By Bochner’s theorem, there exists a probability measure  $\sigma$ , referred to as the spectral measure of  $\varphi$ , such that

$$\forall \mathbf{t}, \mathbf{s} \in \mathcal{X}, \quad \mathcal{F}[\varphi_{\mathbf{t}}](\mathbf{s}) = \sigma(\mathbf{s}) e^{-i\langle \mathbf{s}, \mathbf{t} \rangle}. \quad (3.16)$$

We consider the following assumption.

**Assumption ( $\mathbf{A}_{\text{conv}}$ ).** *The spectral measure of  $\varphi$  has compact support.*

**Super Resolution** In this case,  $\mathcal{X} = \mathbb{T}^d$  and  $\sigma = \sum_{\mathbf{u} \in \mathbb{Z}^d} \sigma_{\mathbf{u}} \delta_{\mathbf{u}}$  is a probability measure on  $\mathbb{Z}^d$ . When  $\sigma$  is the uniform measure on  $[-f_c, f_c]^d$ , with  $f_c \geq 1$ , we uncover the Super-Resolution framework and  $\varphi$  is the Dirichlet kernel. For latter use, we denote

$$\int e^{-i\langle \mathbf{u}, \mathbf{t} - \mathbf{s} \rangle} d\sigma(\mathbf{u}) := \sum_{\mathbf{u} \in \mathbb{Z}^d} \sigma_{\mathbf{u}} e^{-i\langle \mathbf{u}, \mathbf{t} - \mathbf{s} \rangle}.$$

**Continuous sampling Fourier transform** In this case,  $\mathcal{X}$  is a compact set of  $\mathbb{R}^d$  and  $\sigma$  is a probability measure on  $\mathbb{R}^d$ . For latter use, we denote

$$\int e^{-i\langle \mathbf{u}, \mathbf{t} - \mathbf{s} \rangle} d\sigma(\mathbf{u}) := \int_{\mathbb{R}^d} e^{-i\langle \mathbf{u}, \mathbf{t} - \mathbf{s} \rangle} d\sigma(\mathbf{u}).$$

By Lemma A.1, we can assume that  $\mathbb{H}$  is the RKHS associated with the kernel defined by  $\varphi$ . In particular,

$$\forall \mathbf{t} \in \mathcal{X}, \quad \langle \varphi_{\mathbf{t}}, \varphi_{\mathbf{t}_j} \rangle_{\mathbb{H}} = \varphi_{\mathbf{t}_j}(\mathbf{t}) \text{ and } \langle \varphi_{\mathbf{t}}, \mathbf{y} \rangle_{\mathbb{H}} = \mathbf{y}(\mathbf{t}). \quad (3.17)$$

**About the noise term and the regularity of the observation** Using the projection  $\Pi$  and the isometry  $\square$  of Lemma A.1, we can assume without loss of generality that the noise term  $\mathbf{e}$  belongs to  $\mathbb{H}$ , the RKHS associated with  $\varphi$ . By Assumption ( $\mathbf{A}_{\text{conv}}$ ), it can be shown that  $\varphi$  is smooth, hence all the elements of  $\mathbb{H}$  are smooth and so is  $\mathbf{e}$ . Since  $\mathcal{X}$  is compact, we deduce that  $\mathbf{y}$  and its gradient are bounded.

### 3.3.3 Gradients

For any  $\mu \in \mathcal{M}(\mathbb{R}^d)$ , the observation  $\mathbf{y}$  is compared to  $\Phi(\mu) = \varphi \star \mu$  inside the criterion  $J(\mu)$ . A direct application of Proposition B.1 leads to  $J'_{\nu} = \Phi^*(\Phi(\nu) - \mathbf{y}) + \lambda$ , and one can check that, for any  $\mathbf{t} \in \mathbb{R}^d$ ,

$$J'_{\nu}(\mathbf{t}) = \langle \varphi_{\mathbf{t}}, \Phi(\nu) - \mathbf{y} \rangle_{\mathbb{H}} + \lambda.$$

The latter formula can be re-written as

$$J'_{\nu}(\mathbf{t}) = \langle \varphi_{\mathbf{t}}, \Phi(\nu) \rangle_{\mathbb{H}} - \langle \varphi_{\mathbf{t}}, \mathbf{y} \rangle_{\mathbb{H}} + \lambda = \sum_{j=1}^p \omega_j \langle \varphi_{\mathbf{t}}, \varphi_{\mathbf{t}_j} \rangle_{\mathbb{H}} - \langle \varphi_{\mathbf{t}}, \mathbf{y} \rangle_{\mathbb{H}} + \lambda. \quad (3.18)$$

By (3.16) and (3.17), it yields that

$$J'_{\nu}(\mathbf{t}) = \sum_{j=1}^p \omega_j \int e^{-i\langle \mathbf{u}, \mathbf{t} - \mathbf{t}_j \rangle} d\sigma(\mathbf{u}) - \mathbf{y}(\mathbf{t}) + \lambda. \quad (3.19)$$

### 3.3.4 Assumptions

We can identify the  $\mathbf{g}$  and  $\mathbf{h}$  functions appearing in Assumptions ( $\mathbf{A}_1$ ) and ( $\mathbf{A}_2$ ). It holds

$$\mathbf{g}_{\mathbf{t}, \mathbf{t}'}(\mathbf{u}) := e^{-i\langle \mathbf{u}, \mathbf{t} - \mathbf{t}' \rangle} \quad (3.20)$$

$$\mathbf{h}_{\mathbf{t}} := \mathbf{y}(\mathbf{t}), \quad (3.21)$$

which are bounded functions and Assumption ( $\mathbf{A}_1$ ) is satisfied. Using the dominated convergence theorem for the bounded gradients

$$\nabla_{\mathbf{t}} \mathbf{g}_{\mathbf{t}, \mathbf{t}'}(\mathbf{u}) := -i\mathbf{u} e^{-i\langle \mathbf{u}, \mathbf{t} - \mathbf{t}' \rangle} \mathbf{1}_{\{\mathbf{u} \in \text{Supp}(\sigma)\}} \quad (3.22)$$

$$\nabla_{\mathbf{t}} \mathbf{h}_{\mathbf{t}} := \nabla \mathbf{y}(\mathbf{t}), \quad (3.23)$$

where  $\text{Supp}(\sigma)$  denotes the support of the spectral measure  $\sigma$ . One can check that Assumption ( $\mathbf{A}_2$ ) is satisfied under Assumption ( $\mathbf{A}_{\text{conv}}$ ).

### 3.3.5 Stochastic counterpart

Given  $\nu = \nu(\mathbb{W}, \mathbb{T})$ , we introduce a random variable  $Z = (T, U)$  built as follows (there is no  $V$  random variable in this case). Consider two independent random variables given by

$$T \sim \frac{\nu}{\|\nu\|_{\text{TV}}} \text{ and } U \sim \sigma, \quad (3.24)$$

where  $\nu$  assumed nonnegative wlog as discussed in (1.7), hence  $\nu/\|\nu\|_{\text{TV}}$  is a discrete probability measure. We then define

$$J'_{\nu, \mathbf{t}}(Z) = \|\nu\|_{\text{TV}} e^{-i\langle U, \mathbf{t} - T \rangle} - \mathbf{y}(\mathbf{t}) + \lambda, \quad (3.25)$$

$$= \|\nu\|_{\text{TV}} \mathbf{g}_{\mathbf{t}, T}(U) - \mathbf{h}_{\mathbf{t}}(V) + \lambda, \quad (3.26)$$

writing, with a slight abuse of notation,  $\mathbf{h}_t(V) = \mathbf{y}(t)$ . We hence uncover (2.9). According to (3.18) and (3.24), we can verify that once  $\nu$  is fixed,  $J'_{\nu,t}(Z)$  is an unbiased estimation of  $J'_\nu$  and we can observe that there exists a random variable  $\xi_\nu(\mathbf{t}, Z)$  such that

$$J'_{\nu,t}(Z) = J'_\nu(\mathbf{t}) + \xi_\nu(\mathbf{t}, Z) \quad \text{with} \quad \mathbb{E}_Z[\xi_\nu(\mathbf{t}, Z)] = 0. \quad (3.27)$$

A same remark occurs for the gradient of  $J'_\nu$  since for any  $\mathbf{t} \in \mathbb{R}^d$ ,

$$\nabla_{\mathbf{t}} J'_\nu(\mathbf{t}) = -i \sum_{j=1}^p \omega_j \int \mathbf{u} e^{-i\langle \mathbf{u}, \mathbf{t} - t_j \rangle} d\sigma(\mathbf{u}) - \nabla \mathbf{y}(\mathbf{t}). \quad (3.28)$$

We shall then define

$$\mathbf{D}_{\nu,t}(Z) = -i \|\nu\|_{\text{TV}} U e^{-i\langle \mathbf{u}, \mathbf{t} - T \rangle} - \nabla \mathbf{y}(\mathbf{t}), \quad (3.29)$$

$$= \|\nu\|_{\text{TV}} \nabla_{\mathbf{t}} \mathbf{g}_{\mathbf{t},T}(U) - \nabla_{\mathbf{t}} \mathbf{h}_{\mathbf{t}}(V). \quad (3.30)$$

and we uncover (2.11). We get that

$$\mathbf{D}_{\nu,t}(Z) = \nabla_{\mathbf{t}} J'_\nu(\mathbf{t}) + \zeta_\nu(\mathbf{t}, Z) \quad \text{with} \quad \mathbb{E}_Z[\zeta_\nu(\mathbf{t}, Z)] = 0, \quad (3.31)$$

for some centered random vector  $\zeta_\nu(\mathbf{t}, Z)$ .

## 4 Main results

In this section, we state our main results related to the behaviour of Algorithm 1 when the number of iterations become large. The starting point of our analysis is Proposition B.1. Below, we denote by  $\nu_k$  the measure of particles produced at step  $k$  by Algorithm 1. Our main contributions are threefold:

- A control on the total-variation norm of the measures  $\nu_k$  along the different iterations.
- A global minimisation result.
- A local investigation on the evolution of  $J'_{\nu_k}$  and its gradient.

We emphasize that these results are stated in a finite horizon setting and are therefore non-asymptotic in terms of  $K$ .

Here and below, we will require additional notation. We introduce

$$\begin{aligned} \|\underline{\varphi}\|_{\mathbb{H}} &= \inf_{s,t \in \mathcal{X}} \langle \varphi_t, \varphi_s \rangle_{\mathbb{H}}, & \|\varphi\|_{\infty, \mathbb{H}} &= \sup_{t \in \mathcal{X}} \|\varphi_t\|_{\mathbb{H}}, & \|\mathbf{g}\|_{\text{Inf}} &= \inf_{s,t,u} \mathbf{g}_{t,s}(u), \\ \|\mathbf{g}\|_{\infty} &:= \sup_{t,s,u} |\mathbf{g}_{t,s}(u)| & \text{and} & & \|\mathbf{h}\|_{\infty} &:= \sup_{t,v} |\mathbf{h}_t(v)|, \end{aligned}$$

where the functions  $\mathbf{g}$  and  $\mathbf{h}$  have been introduced in Section 2.2.2.

### 4.1 Boundedness of the sequence

The next result establishes a preliminary upper bound of the total variation norm of  $(\nu_k)_{k \geq 1}$  that holds *uniformly* over the iterations. It will be the key to obtain the convergence towards minimizers.

**Proposition 4.1.** *Assume that  $\|\mathbf{g}\|_{\text{Inf}} > 0$ . The algorithm is initialized with a measure  $\nu_0$  such that  $\|\nu_0\|_{\text{TV}} \leq R_0$  where*

$$R_0 = \frac{\|\mathbf{h}\|_{\infty} + 2\mathcal{C}_1}{\|\mathbf{g}\|_{\text{Inf}}} (1 + \alpha(\|\mathbf{h}\|_{\infty} + 2\mathcal{C}_1)) \quad \text{and} \quad \mathcal{C}_1 = \max(\|\mathbf{g}\|_{\infty}, \|\mathbf{h}\|_{\infty} + \lambda). \quad (4.1)$$

Assume that  $\alpha$  is chosen such that:

$$\alpha \leq \frac{1}{2\mathcal{C}_1(R_0 + 1)^2}. \quad (4.2)$$

Then, for any  $k \in \mathbb{N}$ , we have

$$\|\nu_k\|_{\text{TV}} \leq R_0.$$

Provided the mass of the measure  $\nu_0$  at the initialization step is not too large,  $\|\nu_k\|_{TV}$  remains bounded along the iterations of the algorithm. We stress that the control is deterministic although we consider a stochastic algorithm. The main ingredient of the proof is to show that  $J'_{\nu_k}$  (together with its stochastic counterpart) can be related, up to some constants, to  $\|\nu_k\|_{TV}$ . The update displayed in Algorithm 1 then allows to conclude. The complete proof is postponed to Section 6.1.

The assumption  $\|\mathbf{g}\|_{Inf} > 0$  appears to be quite reasonable as soon as  $\|\varphi\|_{\mathbb{H}} > 0$ . For any  $s, t \in \mathcal{X}$ , the variable  $\mathbf{g}_{s,t}(U)$  is indeed a stochastic approximation of  $\langle \varphi_s, \varphi_t \rangle_{\mathbb{H}}$ . We get that  $\|\mathbf{g}\|_{Inf} > 0$  with common assumptions on the construction of this approximation.

**Remark 4.1.** *The inequality displayed in (4.2) implicitly provides a condition on the parameter  $\alpha$ . Indeed, considering for instance the specific case of Gaussian mixtures, we can notice that, provided  $\mathcal{X} = [0, 1]^d$ ,*

$$\|\underline{\theta}\|_{\mathbb{H}} = \left[ \frac{1}{2\sqrt{\pi}} e^{-1/4} \right]^d.$$

Assuming for instance that  $\|\mathbf{g}\|_{Inf} \geq C\|\underline{\theta}\|_{\mathbb{H}}$  for some constant  $C \in ]0, 1[$ , Inequality (4.2) hence holds as soon as  $\alpha \leq \bar{\eta}^d$  for some positive  $\bar{\eta} \in ]0, 1[$ .

## 4.2 Global minimization with swarm stochastic optimisation

Consider  $\mu^*$  a measure that *globally* minimizes  $J$ , obtained from Theorem 1.1. The aim of this section is to show that our algorithm can produce a solution close to  $\mu^*$  (in a sense made precise in Theorem 4.1 below), under some specific conditions. Given a fixed number  $K$  of iterations of Algorithm 1, we define the Cesaro average of our sequence  $(\nu_k)_{k \geq 0}$  by:

$$\bar{\nu}_K = \frac{1}{K+1} \sum_{k=0}^K \nu_k. \quad (4.3)$$

We then obtain the following global minimization result whose proof is displayed in Section 6.2.

**Theorem 4.1.** *Consider an integer  $K$  and assume that the learning rates are chosen as:*

$$\alpha = \sqrt{\frac{d\|\mu^*\|_{TV}}{R_0^3 K}} \quad \text{and} \quad \eta = \sqrt{\frac{dR_0}{K^3\|\mu^*\|_{TV}}},$$

where  $R_0$  is introduced (4.1). Assume furthermore that  $K$  is picked large enough so that  $\alpha$  satisfies (4.2), and assume that the measure  $\nu_0$  is uniformly distributed over a uniform grid of step-size  $\delta = 2\sqrt{\frac{d}{\|\mu^*\|_{TV}KM}}$ , then:

$$\mathbb{E} [J(\bar{\nu}_K) - J(\mu^*)] \leq \mathfrak{C} \sqrt{\frac{d\|\mu^*\|_{TV}R_0^3}{K}} \left[ \log(d\|\mu^*\|_{TV}R_0^3K) + \frac{\log(|\mathcal{X}|)}{d} \right],$$

for some positive constant  $\mathfrak{C}$  depending only on  $\|\varphi\|_{\infty, \mathbb{H}}$ ,  $\|\mathbf{y}\|_{\mathbb{H}}$ ,  $\|\bar{\sigma}'\|_{\mathbb{H}}$ ,  $\|\mathbf{h}'\|_{\mathbb{H}}$ .

A careful inspection of the previous upper bound shows that to obtain an  $\epsilon$  approximation with our Cesaro averaged measure  $\bar{\nu}_K$ , (while removing the effect of the log term) we need to choose  $K$  as:

$$K_\epsilon = dR_0^3\|\mu^*\|_{TV}\epsilon^{-2}.$$

Then, the grid step-size is then of the order  $\delta_\epsilon$  given by:

$$\delta_\epsilon = \frac{\epsilon}{\|\mu^*\|_{TV}}.$$

We finally observe that the number of particles  $p$  needed to obtain an  $\epsilon$  approximation is then of the order:

$$p_\epsilon = |\mathcal{X}|\|\mu^*\|_{TV}^d \epsilon^{-d}.$$

Hence, if the number of iteration varies polynomially in terms of  $\epsilon^{-2}$ , we observe the degradation of the number of particles needed to well approximate any distribution over  $\mathcal{X}$  in terms of the dimension  $d$ .

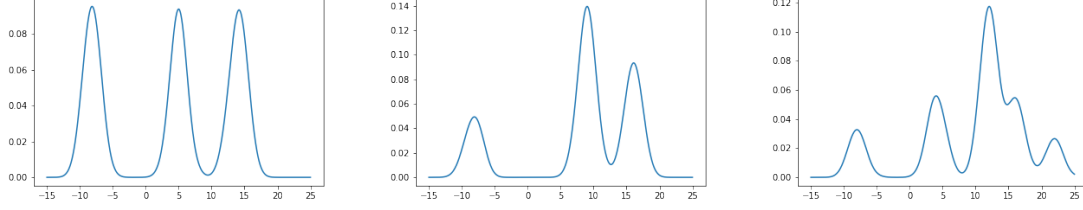


Figure 1: Three 1-D Gaussian mixture distributions to be learnt by Supermix and Stochastic Conic Particle Gradient Descent.

### 4.3 Local minimization with swarm stochastic optimisation

To conclude this contribution, we state a complementary result that quantifies the behaviour of Algorithm 1 when the number of particles used is not “as large” as the one indicated in Theorem 4.1.

**Theorem 4.2.** Assume that  $\alpha = \eta = 1/\sqrt{K}$  is chosen such that Proposition 4.1 holds. If  $\tau_K$  refers to a random variable uniformly distributed over  $\{1, \dots, K\}$ , independent from  $(\nu_k)_{k \geq 1}$ , then:

$$\mathbb{E} \left[ \|J'_{\nu_{\tau_K}}\|_{\nu_{\tau_K}}^2 + \|\nabla J'_{\nu_{\tau_K}}\|_{\nu_{\tau_K}}^2 \right] \leq \frac{J(\nu_0) + \mathfrak{C}(1 + R_0^4)}{\sqrt{K}},$$

for some constant  $\mathfrak{C}$  that depends on  $\|\varphi\|_{Lip}$ .

Even though weaker than Theorem 4.1, the previous “local” result deserves several comments as it raises some meaningful and challenging questions.

The previous result is a strong indicator of the convergence of our swarm stochastic particle algorithms towards a minimizer of  $J$ . Indeed, Proposition B.2 stated in Appendix B states that any minimizer  $\mu^*$  of  $J$  necessarily satisfies  $J'_{\mu^*} = 0$  on the support of  $\mu^*$  and  $J'_{\mu^*} \geq 0$  everywhere, which implicitly means that  $\nabla J'_{\mu^*} = 0$  on the support of  $\mu^*$  otherwise the previous positivity condition would not hold. In Theorem 4.2, we obtain that  $\|J'_{\nu_k}\|_{\nu_{\tau_K}}^2$  and  $\|\nabla J'_{\nu_k}\|_{\nu_{\tau_K}}^2$  become arbitrarily small when  $k$  becomes larger and larger, which is perfectly in accordance with the previous conditions. Nevertheless, the nature of our algorithm does not permit to push further our conclusions, especially about the positivity of  $J'_{\nu_k}$  *everywhere* (and not only on the support of  $\nu_k$ ). In particular, to extend our conclusion towards a global minimization result, we need to be able to address a kind of “density of the support”, which is absolutely unattainable in our present analysis without multiplying the number of particles all over the state space, which is in some sense what is done in Theorem 4.1.

## 5 Numerical experiments on FastPart

### 5.1 Experimental setup

In this short section, we develop a brief numerical study, that may be seen as a proof of concept, to assess the efficiency of our stochastic gradient descent, when using some sketched randomized evaluations of  $J'_\nu$ , with the help of sampling involved in Equations (3.3) and (3.8).<sup>1</sup>

For this purpose, we consider the Supermix problem introduced in De Castro et al. [2021], which is described in Section 3.1, when considering a mixture of Gaussian densities. We consider three toy situations in 1D. Figure 1 represents the mixture densities considered in this study, that contains for two of them 3 components, and for the last one 5 components.

<sup>1</sup>Our simulations greatly benefit from the previous work of Nicolas Jouvin <https://nicolasjouvin.github.io/>, while the original numerical Python code is made available here [https://forgemia.inra.fr/njouvin/particle\\_blasso](https://forgemia.inra.fr/njouvin/particle_blasso).



## 5.2 Benchmark

Our experimental setup is essentially built with the help of three versions of the Conic Particle Gradient Descent.

- The first method we will use is the deterministic CPGD introduced in Chizat [2022], implemented by N. Jouvin [https://forgemia.inra.fr/njouvin/particle\\_blasso](https://forgemia.inra.fr/njouvin/particle_blasso). This method depends on the number of particles we use, the learning rate that encodes the gain of the algorithm at each iteration, and the number of iterations.
- The second method is the Stochastic-CPGD introduced in this work. Our method depends on the same previous set of parameters (number of particles, learning rate, number of iterations) and of the batch size of the data we sample per iteration and the number of randomly sketched frequencies.
- The last method is simply the Cesaro averaged counterpart of our Stochastic-CPGD, but this method raises some technical computational difficulties since averaging a sequence of measures seriously complicates the final estimates  $\bar{\nu}_K = \frac{1}{K+1} \sum_{k=0}^K \nu_k$ . To overcome this difficulty, we have chosen to use instead the measures supported by the averaged means all along the trajectory of the S-CPGD, and weighted by the averaged weights of the S-CPGD. For this purpose, we introduce

$$\forall j \in \{1, \dots, p\} \quad \forall K \geq 0 \quad \bar{\mathbf{t}}_j^K = \frac{1}{K+1} \sum_{k=0}^K \mathbf{t}_j^k \quad \text{and} \quad \bar{\omega}_j^K = \frac{1}{K+1} \sum_{k=0}^K \omega_j^k$$

and we approximate  $\bar{\nu}_K$  with the help of the sequence  $\hat{\nu}_K$ , defined by:

$$\hat{\nu}_K = \sum_{j=1}^p \bar{\omega}_j^K \delta_{\bar{\mathbf{t}}_j^K}. \quad (5.1)$$

Again, this sequence  $\hat{\nu}_K$  depends on several parameters, the learning rate, the number of particles and iterations, and the size of batches and sketches as well.

## 5.3 Results

**Loss function: Averaging vs no averaging** We show in Figure 2 the evolution of the loss function  $J_\nu$  over the iterations of the algorithm. We emphasize that the complexity of the S-CPGD and of the approximated Cesaro average are almost the same, since the sequence  $\hat{\nu}_K$  introduced in (5.1) is a cheap approximation of the true Cesaro averaged sequence  $\bar{\nu}_K$ .

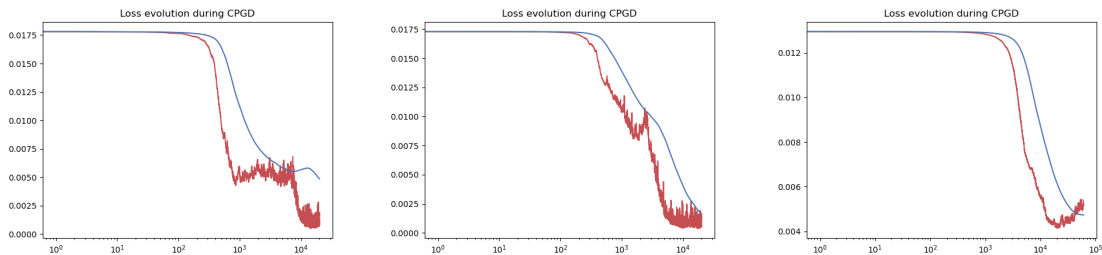


Figure 2: Evolution of the loss (log scale) with our S-CPGD algorithm (red) and its averaged counterpart (in blue) when using our methods on the mixture problems illustrated in Figure 1.

First, as indicated in Figure 2, we shall observe that the sequence  $\hat{\nu}_K$  always produces the desired smoothing effect all over the iterations of the algorithm, while slowing a bit the decrease of the loss function  $J_\nu$  over the iterations. As a consequence, it seems more appropriate to use several long-range parallelized S-CPGD instead of a unique average thread of S-CPGD. In the same time, it may be remarked that our sequence  $(\hat{\nu}_K)_{K \geq 1}$  is a rough approximation of the true Cesaro averaging that is studied in our

paper and the numerical approximation introduced in (5.1) may not be as good as the true Cesaro sequence  $(\bar{D}_K)_{K \geq 1}$ .

Second, as a classical phenomenon in machine learning when using stochastic approximation algorithm, or over-parameterized neural networks, our S-CPGD commonly generates some double-descent phenomena (see the 3 sub-figures of Figure 2) that translates some local minimizer escape of the swarm of particles.

**Loss function: Averaging vs no averaging vs Deterministic** Figure 3 represents the evolution of the cost function with respect to the numerical cost which is a far better indicator than the number of iterations of the algorithm in our case since the S-CPGD algorithm is designed to be much more cheaper than the deterministic CPGD. Figure 3 clearly illustrates the efficiency of our method with regards to the deterministic

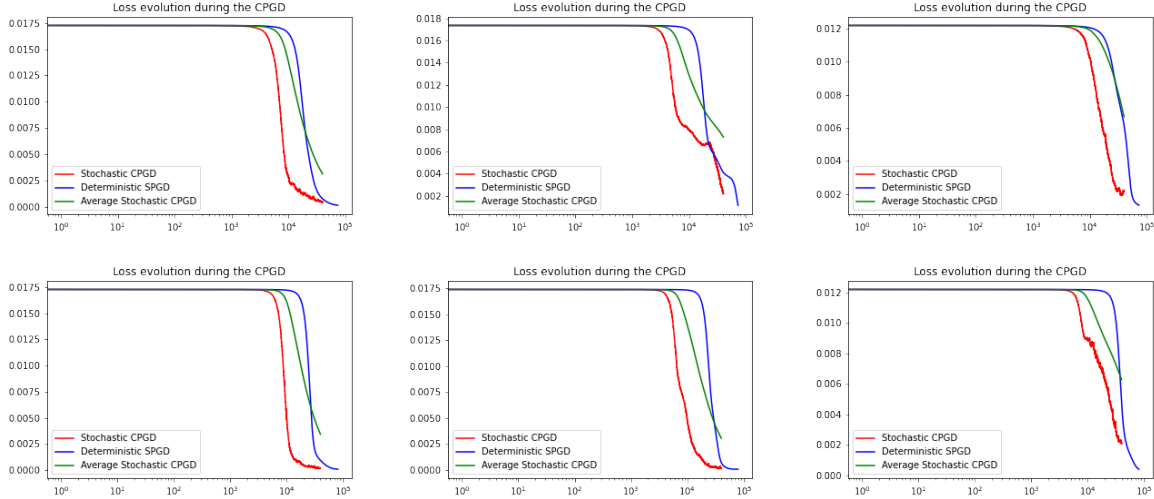


Figure 3: Evolution of the loss (log scale of the **computational time**) with our S-CPGD algorithm (red), its averaged counterpart (in green) and the deterministic CPGD (in blue) on the mixture problem of Figure 1 with 20 particles (top) and with 50 particles (bottom). As indicated in the text, and observed with the shift to the right of the blue curve when compared to the red one, the cost of our S-CPGD is much cheaper than the deterministic CPGD.

one as the red curve shows that the non-averaged S-CPGD produces comparable results as those obtained by CPGD with a significantly lower needs of computational cost: the red curve is clearly shifted on the left when compared to the blue one. It is furthermore possible to quantitatively assess the numerical gain produced by the S-CPGD when compared to the deterministic one: on our toy example, the deterministic CPGD requires approximately 4 more computations to attain the same decrease of the  $J_\nu$ , this effect being even amplified when the number of particles is increasing.

**Loss function: Effect of the number of particles** In the meantime, we observe that the loss function benefits from a large number of particles (see the comparison between top and bottom lines of Figure 3) but this should be tempered by the increasing number of simulations, which varies linearly with the number of particles. We should finally observe that using a large number of particles seems to be important especially in difficult situations (as illustrated in the right column of Figure 3 where using 50 particles instead of 20 significantly improves the loss function, which is not the case on the right column of Figure 3).

The effect of the number of particles can also be illustrated while looking at the trajectories themselves of the particles as shown in Figure 4. We observe that the number of particles is a clear key parameter that strongly influences the success of the method. In our example of 5 components GMM, (last example in Figure 1), we see that a too small number of particles completely miss some components of the mixture while using a strongly over-parameterized set of particles permit to fully recover the support of the mixing distribution, even in the situation where some components of the mixture overlap.

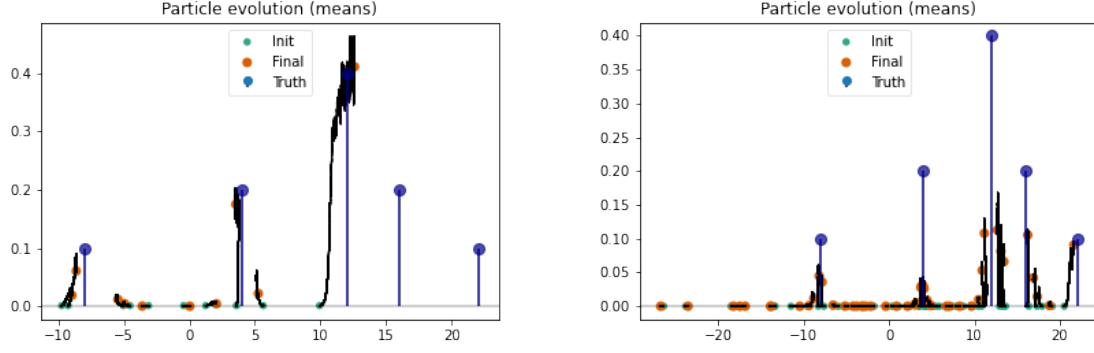


Figure 4: Trajectories of the particles of our  $S$ -CPGD algorithm in the third example of the mixture problem of Figure 1 with 5 modes. Left: trajectories using 10 particles. Right: Same with 50 particles. The l.h.s. shows the behaviour of  $S$ -CPGD when a too small number of particles is used. Particles concentrate around good positions but may miss some of the important locations of the mixture. The r.h.s. demonstrates that a sufficiently large number of particles is necessary to guarantee an exhaustive reconstruction of the mixing distribution.

From our brief numerical study, we can conclude that both sketching and batch subsampling with a stochastic gradient strategy appears to strongly improve the numerical cost of the Conic Particle Gradient Descent, which permits to increase the number of particles used in the mean field approximation. We also have shown that in some difficult inverse problem examples, a large number of particles seems necessary to perfectly recover the solution of the optimisation problem. It appears that the problem seriously benefits from a strong over-parametrisation, that may be handled with our cheap stochastic computing approach, which is not the case in a reasonable time with the deterministic CPGD.

## 6 Proof of the main results

### 6.1 Proof of Proposition 4.1

*Proof.* The principle of this proof is as follows. We identify two radii  $R$  and  $R_0$  such that  $R < R_0$  for which the two following properties hold:

- If  $\|\nu_k\|_{\text{TV}} \leq R < R_0$ , then  $\|\nu_{k+1}\|_{\text{TV}} \leq R_0$
- If  $R \leq \|\nu_k\|_{\text{TV}} \leq R_0$ , then  $\|\nu_{k+1}\|_{\text{TV}} \leq \|\nu_k\|_{\text{TV}}$ .

We will establish these properties with a suitable choice for  $R_0$  and a related condition on  $\alpha$  as displayed in Proposition 4.1. The proof is divided into two steps.

Step 1: total variation norm recursion. Let  $k \in \mathbb{N}$  be fixed, using an appropriate Taylor expansion, we get

$$\begin{aligned}
\|\nu_{k+1}\|_{\text{TV}} &= \sum_{j=1}^p \omega_j^{k+1}, \\
&= \sum_{j=1}^p \omega_j^k e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})}, \\
&= \|\nu_k\|_{\text{TV}} + \sum_{j=1}^p \omega_j^k \left( e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right), \\
&= \|\nu_k\|_{\text{TV}} - \alpha \sum_{j=1}^p \omega_j^k J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) + \mathcal{R}(\alpha, \nu_k, Z^{k+1}),
\end{aligned}$$

where  $\mathcal{R}(\alpha, \nu_k, Z^{k+1})$  is a remaining term whose value will be made precise later on. The last equality can be re-written as

$$\frac{\|\nu_{k+1}\|_{\text{TV}}}{\|\nu_k\|_{\text{TV}}} - 1 = \underbrace{-\alpha \sum_{j=1}^p \tilde{\omega}_j^k \mathbf{J}'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})}_{:=A_1} + \underbrace{\|\nu_k\|_{\text{TV}}^{-1} \mathcal{R}(\alpha, \nu_k, Z^{k+1})}_{:=A_2}, \quad (6.1)$$

where for any  $k \in \mathbb{N}$ ,  $\tilde{\nu}_k := \|\nu_k\|_{\text{TV}}^{-1} \nu_k = \sum_{j=1}^p \tilde{\omega}_j^k \delta_{\mathbf{t}_j^k}$  and  $\tilde{\omega}_j^k = \omega_j^k / \|\nu_k\|_{\text{TV}}$  for any  $j \in \{1, \dots, p\}$ . First, we concentrate our attention on the term  $A_1$ . Using (2.9),

$$\begin{aligned} A_1 &:= -\alpha \sum_{j=1}^p \tilde{\omega}_j^k \mathbf{J}'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}), \\ &= -\alpha \lambda - \alpha \|\nu_k\|_{\text{TV}} \int \mathbf{g}_{t, T^{k+1}}(U^{k+1}) d\tilde{\nu}_k(t) + \alpha \int \mathbf{h}_t(V^{k+1}) d\tilde{\nu}_k(t). \end{aligned}$$

Using a rough bound, we get:

$$\int \mathbf{h}_t(V^{k+1}) d\tilde{\nu}_k(t) \leq \|\mathbf{h}\|_{\infty}. \quad (6.2)$$

where we have used  $\|\tilde{\nu}_k\|_{\text{TV}} = 1$ . Moreover,

$$\begin{aligned} \|\nu_k\|_{\text{TV}} \int \mathbf{g}_{t, T^{k+1}}(U^{k+1}) d\tilde{\nu}_k(t) &= \|\nu_k\|_{\text{TV}} \sum_{j=1}^p \tilde{\omega}_j^k \mathbf{g}_{\mathbf{t}_j^k, T^{k+1}}(U^{k+1}), \\ &\geq \|\nu_k\|_{\text{TV}} \|\mathbf{g}\|_{\text{Inf}}, \end{aligned}$$

the last term being strictly positive thanks to our assumption. We eventually get, gathering the previous bounds

$$\begin{aligned} A_1 &\leq -\alpha \lambda - \alpha \|\mathbf{g}\|_{\text{Inf}} \|\nu_k\|_{\text{TV}} + \alpha \|\mathbf{h}\|_{\infty} \\ &= -\alpha [\lambda + \|\mathbf{g}\|_{\text{Inf}} \|\nu_k\|_{\text{TV}} - \|\mathbf{h}\|_{\infty}]. \end{aligned} \quad (6.3)$$

Now, we turn our attention on the term  $A_2$ . Recall that

$$\begin{aligned} A_2 &:= \|\nu_k\|_{\text{TV}}^{-1} \mathcal{R}(\alpha, \nu_k, Z^{k+1}), \\ &= \sum_{j=1}^p \tilde{\omega}_j^k \left[ e^{-\alpha \mathbf{J}'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 + \alpha \mathbf{J}'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) \right]. \end{aligned}$$

Using the inequality

$$e^{-x} \leq 1 - x + \frac{x^2}{2} e^{|x|} \quad \forall x \in \mathbb{R},$$

we get

$$A_2 \leq \sum_{j=1}^p \tilde{\omega}_j^k (\alpha \mathbf{J}'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}))^2 \times e^{\alpha |\mathbf{J}'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})|}.$$

According to Lemma B.1, for any  $l \in \{1, \dots, p\}$  and  $k \in \mathbb{N}$ ,

$$\mathbf{J}'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) \leq \mathcal{C}_1 (\|\nu_k\|_{\text{TV}} + 1) \quad \text{with} \quad \mathcal{C}_1 = \max(\|\mathbf{g}\|_{\infty}; \|\mathbf{h}\|_{\infty} + \lambda).$$

Hence:

$$\begin{aligned} A_2 &\leq \sum_{j=1}^p \tilde{\omega}_j^k (\alpha \mathcal{C}_1 (\|\nu_k\|_{\text{TV}} + 1))^2 \times e^{\alpha \mathcal{C}_1 (\|\nu_k\|_{\text{TV}} + 1)}, \\ &\leq [\alpha \mathcal{C}_1 (\|\nu_k\|_{\text{TV}} + 1)]^2 e^{\alpha \mathcal{C}_1 (\|\nu_k\|_{\text{TV}} + 1)}. \end{aligned} \quad (6.4)$$

Assume that  $\|\nu_k\|_{\text{TV}} \leq R_0$  for some  $R_0$ . Gathering (6.1), (6.3) and (6.4), we obtain that:

$$\begin{aligned} & \frac{\|\nu_{k+1}\|_{\text{TV}}}{\|\nu_k\|_{\text{TV}}} - 1 \\ & \leq -\alpha \left[ \lambda + \|\mathbf{g}\|_{\text{Inf}} \|\nu_k\|_{\text{TV}} - \|\mathbf{h}\|_{\infty} \right] + \alpha^2 \mathcal{C}_1^2 (\|\nu_k\|_{\text{TV}} + 1)^2 e^{\alpha \mathcal{C}_1 (\|\nu_k\|_{\text{TV}} + 1)}, \\ & = -\alpha \left[ \lambda + \|\mathbf{g}\|_{\text{Inf}} \|\nu_k\|_{\text{TV}} - \left( \|\mathbf{h}\|_{\infty} + \alpha \mathcal{C}_1^2 (\|\nu_k\|_{\text{TV}} + 1)^2 e^{\alpha \mathcal{C}_1 (\|\nu_k\|_{\text{TV}} + 1)} \right) \right], \\ & \leq -\alpha \left[ \lambda + \|\mathbf{g}\|_{\text{Inf}} \|\nu_k\|_{\text{TV}} - \left( \|\mathbf{h}\|_{\infty} + \alpha \mathcal{C}_1^2 (R_0 + 1)^2 e^{\alpha \mathcal{C}_1 (R_0 + 1)} \right) \right]. \end{aligned}$$

Step 2: calibration of  $R, R_0$  and  $\alpha$ . We first make precise the value of the parameters involved in our bound and then verify that the stability property holds. We define  $R_0$  and  $R$  as:

$$R = \frac{\|\mathbf{h}\|_{\infty} + 2\mathcal{C}_1}{\|\mathbf{g}\|_{\text{Inf}}}, \quad R_0 = R(1 + \alpha(\|\mathbf{h}\|_{\infty} + 2\mathcal{C}_1)), \quad \text{and} \quad \alpha \leq \frac{1}{2\mathcal{C}_1(R_0 + 1)^2},$$

where the constant  $\mathcal{C}_1$  is introduced in (B.5). At this step, two different situations can occur.

- 1<sup>st</sup> case: We consider the case where  $\|\nu_k\|_{\text{TV}} \geq R$ . In such a situation, the choice of  $\alpha$  induces that:

$$\begin{aligned} \frac{\|\nu_{k+1}\|_{\text{TV}}}{\|\nu_k\|_{\text{TV}}} - 1 & \leq -\alpha \left[ \lambda + \|\mathbf{g}\|_{\text{Inf}} \|\nu_k\|_{\text{TV}} - \left( \|\mathbf{h}\|_{\infty} + \alpha \mathcal{C}_1^2 (R_0 + 1)^2 e^{\alpha \mathcal{C}_1 (R_0 + 1)} \right) \right], \\ & \leq -\alpha \left[ \lambda + \|\mathbf{g}\|_{\text{Inf}} \|\nu_k\|_{\text{TV}} - \left( \|\mathbf{h}\|_{\infty} + \mathcal{C}_1 e^{\frac{1}{2(R_0 + 1)}} \right) \right], \\ & \leq -\alpha \left[ \lambda + \|\mathbf{g}\|_{\text{Inf}} R - (\|\mathbf{h}\|_{\infty} + 2\mathcal{C}_1) \right], \\ & \leq 0. \end{aligned}$$

The last inequality induces

$$\|\nu_{k+1}\|_{\text{TV}} \leq \|\nu_k\|_{\text{TV}} \leq R_0.$$

- 2<sup>nd</sup> case: We consider the situation where  $\|\nu_k\|_{\text{TV}} \leq R$ . With the same condition on  $\alpha$ , we use the rough bound

$$\begin{aligned} \frac{\|\nu_{k+1}\|_{\text{TV}}}{\|\nu_k\|_{\text{TV}}} - 1 & \leq -\alpha \left[ \lambda + \|\mathbf{g}\|_{\text{Inf}} \|\nu_k\|_{\text{TV}} - \left( \|\mathbf{h}\|_{\infty} + \alpha \mathcal{C}_1^2 (R_0 + 1)^2 e^{\alpha \mathcal{C}_1 (R_0 + 1)} \right) \right], \\ & \leq \alpha (\|\mathbf{h}\|_{\infty} + 2\mathcal{C}_1). \end{aligned}$$

This implies that

$$\begin{aligned} \|\nu_{k+1}\|_{\text{TV}} & \leq \|\nu_k\|_{\text{TV}} (1 + \alpha(\|\mathbf{h}\|_{\infty} + 2\mathcal{C}_1)), \\ & \leq R(1 + \alpha(\|\mathbf{h}\|_{\infty} + 2\mathcal{C}_1)) = R_0. \end{aligned}$$

Finally, we end the proof with an induction argument. □

## 6.2 Proof of Theorem 4.1

### 6.2.1 The shadow sequence

Consider an integer  $k \geq 0$  and the map  $\mathcal{T}_{k+1}$  defined as:

$$\forall t \in \mathcal{X} \quad \mathcal{T}_{k+1}(t) = t - \eta \mathbf{D}_{\nu_k}(t, Z^{k+1}), \quad (6.5)$$

where  $(Z^{k+1})_{k \geq 0}$  is the sequence of random variables sampled in Algorithm 1 and used in Equations (2.14) and (2.15). The sequence of maps  $(\mathcal{T}_{k+1})_{k \geq 0}$  only acts on the positions of  $\mathcal{X}$  and is built with the random sequence  $(\nu_k)_{k \geq 1}$ .

Using  $(\mathcal{T}_{k+1})_{k \geq 0}$ , we then define the shadow sequence  $(\nu_k^\varepsilon)_{k \geq 1}$  obtained through an iterative push-forward from a given initialisation measure  $\nu_0^\varepsilon \in \mathcal{M}(\mathcal{X})_+$  with the sequence of maps  $(\mathcal{T}_k)_{k \geq 1}$ . More formally, we set in an iterative way

$$\nu_{k+1}^\varepsilon = \mathcal{T}_{k+1}^\#(\nu_k^\varepsilon) \quad \forall k \in \mathbb{N}^*, \quad (6.6)$$

where, for any continuous function  $\psi$ ,

$$\int_{\mathcal{X}} \psi d\mathcal{T}_{k+1}^\#(\nu) = \int \psi(\mathcal{T}_{k+1}(\cdot)) d\nu \quad \forall \nu \in \mathcal{M}(\mathcal{X}).$$

The measure  $\nu_0^\varepsilon$  will be defined carefully at the very end of our study. Roughly speaking, the shadow sequence  $(\nu_k^\varepsilon)_{k \geq 1}$  moves exactly like  $(\nu_k)_{k \geq 1}$  and will share the same support, but the weights on the particles for the sequence  $(\nu_k^\varepsilon)_{k \geq 1}$  will be optimised to allow for a good approximation of  $\mu^*$ . In particular, if we decompose the initial measure  $\nu_0$  as

$$\nu_0 = \sum_{j=1}^p \omega_j \delta_{\mathbf{t}_j^0},$$

then we can write  $\nu_0^\varepsilon = \nu(\mathbb{W}^\varepsilon, \mathbf{t}^0)$ , so that:

$$\nu_0^\varepsilon = \sum_{j=1}^p \omega_j^\varepsilon \delta_{\mathbf{t}_j^0},$$

for some weights  $(\omega_j^\varepsilon)_{j=1..p}$  that will be chosen in an appropriate way.

## 6.2.2 Excess risk decomposition

The starting point is Proposition B.1 that is used with  $\nu = \nu_k$  and  $\sigma = \mu^* - \nu_k$ . We write:

$$\begin{aligned} J(\nu_k) - J(\mu^*) &= \int_{\mathcal{X}} J'_{\nu_k} d[\nu_k - \mu^*] - \frac{1}{2} \|\Phi(\mu^* - \nu_k)\|_{\mathbb{H}}^2 \\ &= \underbrace{\int_{\mathcal{X}} J'_{\nu_k} d[\nu_k - \nu_k^\varepsilon]}_{:=\textcircled{1}} + \underbrace{\int_{\mathcal{X}} J'_{\nu_k} d[\nu_k^\varepsilon - \mu^*]}_{:=\textcircled{2}} - \frac{1}{2} \|\Phi(\mu^* - \nu_k)\|_{\mathbb{H}}^2 \end{aligned} \quad (6.7)$$

where  $(\nu_k^\varepsilon)_{k \geq 1}$  is the auxiliary shadow sequence of measures introduced in (6.6). First, we establish that the mirror descent adapts the weights of  $(\nu_k)_{k \geq 1}$  to those of the shadow sequence  $(\nu_k^\varepsilon)_{k \geq 1}$ . For any  $\mu_1, \mu_2 \in \mathcal{M}(\mathcal{X})_+$ , we introduce the following entropy:

$$\mathcal{H}(\mu_1, \mu_2) = - \int_{\mathcal{X}} \log \left( \frac{d\mu_1}{d\mu_2} \right) d\mu_2 - \|\mu_2\|_{\text{TV}} + \|\mu_1\|_{\text{TV}}. \quad (6.8)$$

The next proposition focuses on the first term of Equation (6.7).

**Proposition 6.1.** *Term  $\textcircled{1}$  of Equation (6.7) may be decomposed as:*

$$\begin{aligned} \textcircled{1} &= \int_{\mathcal{X}} J'_{\nu_k} d[\nu_k - \nu_k^\varepsilon] = \frac{1}{\alpha} [\mathcal{H}(\nu_k, \nu_k^\varepsilon) - \mathcal{H}(\nu_{k+1}, \nu_{k+1}^\varepsilon)] \\ &\quad + \frac{1}{\alpha} \sum_{j=1}^p \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right] + \sum_{j=1}^p \omega_j^\varepsilon \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) \end{aligned}$$

*Proof.* Since both measures  $\nu_k$  and  $\nu_k^\varepsilon$  share the same particle locations  $\mathbf{t}^k$ , we can remark that

$$\textcircled{1} = \int_{\mathcal{X}} J'_{\nu_k} d[\nu_k - \nu_k^\varepsilon] = \sum_{j=1}^p [\omega_j^k - \omega_j^\varepsilon] J'_{\nu_k}(\mathbf{t}_j^k) \quad (6.9)$$

We then observe from Equation (2.14) that:

$$\omega_j^{k+1} = \omega_j^k e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} \quad \Rightarrow \quad J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) = -\frac{1}{\alpha} \log \left( \frac{\omega_j^{k+1}}{\omega_j^k} \right).$$

Using now Equation (2.10), we observe that:

$$J'_{\nu_k}(\mathbf{t}_j^k) = -\frac{1}{\alpha} \log \left( \frac{\omega_j^{k+1}}{\omega_j^k} \right) - \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}).$$

We then use the previous equality in (6.9) and obtain that:

$$\begin{aligned} \textcircled{1} &= \sum_{j=1}^p \left( \omega_j^k J'_{\nu_k}(\mathbf{t}_j^k) - \omega_j^\varepsilon \left[ -\frac{1}{\alpha} \log \left( \frac{\omega_j^{k+1}}{\omega_j^k} \right) - \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) \right] \right), \\ &= \sum_{j=1}^p \left( \omega_j^k J'_{\nu_k}(\mathbf{t}_j^k) + \omega_j^\varepsilon \frac{1}{\alpha} \log \left( \frac{\omega_j^{k+1}}{\omega_j^k} \right) \right) + \sum_{j=1}^p \omega_j^\varepsilon \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) \\ &= \frac{1}{\alpha} \sum_{j=1}^p \left( \alpha \omega_j^k J'_{\nu_k}(\mathbf{t}_j^k) + \omega_j^\varepsilon \left[ \log \left( \frac{\omega_j^{k+1}}{\omega_j^\varepsilon} \right) - \log \left( \frac{\omega_j^k}{\omega_j^\varepsilon} \right) \right] \right) + \sum_{j=1}^p \omega_j^\varepsilon \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}). \end{aligned}$$

We use the entropy  $\mathcal{H}$  introduced in Equation (6.8) and deduce that:

$$\begin{aligned} \textcircled{1} &= \frac{1}{\alpha} \sum_{j=1}^p \left[ \alpha \omega_j^k J'_{\nu_k}(\mathbf{t}_j^k) + \omega_j^{k+1} - \omega_j^k \right] + \frac{\mathcal{H}(\nu_k, \nu_k^\varepsilon) - \mathcal{H}(\nu_{k+1}, \nu_{k+1}^\varepsilon)}{\alpha} + \sum_{j=1}^p \omega_j^\varepsilon \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) \\ &= \frac{1}{\alpha} \sum_{j=1}^p \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right] + \frac{\mathcal{H}(\nu_k, \nu_k^\varepsilon) - \mathcal{H}(\nu_{k+1}, \nu_{k+1}^\varepsilon)}{\alpha} + \sum_{j=1}^p \omega_j^\varepsilon \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}). \end{aligned}$$

We then obtain the conclusion of the proof.  $\square$

Now, we study the second term of Equation (6.7), which is an ‘‘approximation’’ term. We essentially follow the same methodology proposed in Chizat [2022] but we use the specificity of our model to properly analyze this term. For this purpose, we use the BL norm (over functions) and dual norm (over measures) introduced in Chizat [2022], defined as:

$$\forall f : \mathcal{X} \rightarrow \mathbb{R} \quad \|f\|_{BL} = \|f\|_\infty + \|f\|_{Lip}, \quad (6.10)$$

where  $\|\cdot\|_\infty$  refers to the supremum norm over  $\mathcal{X}$ ,  $\|\cdot\|_{Lip}$  to the Lipschitz constant for  $f$ , and

$$\forall \nu \in \mathcal{M}(\mathcal{X})_+ \quad \|\nu\|_{BL}^* = \sup_{\|f\|_{BL} \leq 1} \int f d\nu. \quad (6.11)$$

We also introduce the constant  $Lip(\varphi)$  defined as

$$Lip(\varphi) = \sup_{s, t \in \mathcal{X}} \frac{\|\varphi_t - \varphi_s\|_{\mathbb{H}}}{\|t - s\|_{\mathcal{X}}}.$$

Using these notation, we can propose a bound on the second term of Equation (6.7) as displayed in the following proposition.

**Proposition 6.2.** *The approximation term  $\textcircled{2}$  satisfies:*

$$\forall k \geq 1 \quad \textcircled{2} = \int_{\mathcal{X}} J'_{\nu_k} d[\nu_k^\varepsilon - \mu^*] \leq \mathfrak{A}_k \|\nu_k^\varepsilon - \mu^*\|_{BL}^* \quad a.s.$$

where  $\mathfrak{A}_k$  is given by:

$$\mathfrak{A}_k = \mathcal{C}_0 (\|\nu_k\|_{TV} + 1) + Lip(\varphi) [\|\nu_k\|_{TV} \|\varphi\|_{\infty, \mathbb{H}} + \|\mathcal{Y}\|_{\mathbb{H}}], \quad (6.12)$$

where

$$\mathcal{C}_0 = \max(\lambda + \|\varphi\|_{\infty, \mathbb{H}} \|\mathcal{Y}\|_{\mathbb{H}}; \|\varphi\|_{\infty, \mathbb{H}}^2).$$

*Proof.* We can immediately remark that

$$\textcircled{2} = \int_{\mathcal{X}} J'_{\nu_k} d[\nu_k^\varepsilon - \mu^*] \leq \|J'_{\nu_k}\|_{BL} \|\nu_k^\varepsilon - \mu^*\|_{BL}^*.$$

Then, according to Lemma B.1,

$$\|J'_{\nu_k}\|_{BL} = \|J'_{\nu_k}\|_\infty + \|J'_{\nu_k}\|_{Lip} \leq C_0(\|\nu_k\|_{TV} + 1) + \|J'_{\nu_k}\|_{Lip}.$$

To conclude the proof, we have to propose an upper bound on  $\|J'_{\nu_k}\|_{Lip}$ . For any  $s, t \in \mathcal{X}$  we have

$$\begin{aligned} |J'_{\nu_k}(s) - J'_{\nu_k}(t)| &= \left| \sum_{j=1}^p \omega_j^k \langle \varphi_t - \varphi_s, \varphi_{t_j^k} \rangle_{\mathbb{H}} - \langle \varphi_t - \varphi_s, \mathbf{y} \rangle_{\mathbb{H}} \right|, \\ &\leq \sum_{j=1}^p \omega_j^k \|\varphi_t - \varphi_s\|_{\mathbb{H}} \|\varphi_{t_j^k}\|_{\mathbb{H}} + \|\varphi_t - \varphi_s\|_{\mathbb{H}} \|\mathbf{y}\|_{\mathbb{H}}, \\ &\leq Lip(\varphi) [\|\nu_k\|_{TV} \|\varphi\|_{\infty, \mathbb{H}} + \|\mathbf{y}\|_{\mathbb{H}}] \times \|t - s\|_{\mathcal{X}}. \end{aligned}$$

The results is obtained by gathering the previous bounds.  $\square$

We finally introduce a key term that quantifies the way where  $\mu^*$  can be approximated by a discrete measure. This term, denoted by  $\mathcal{Q}$ , is defined as

$$\mathcal{Q}_{\mu^*, \nu_0}(\tau) := \inf_{\mu \in \mathcal{M}(\mathcal{X})_+} \left[ \|\mu^* - \mu\|_{BL}^* + \frac{1}{\tau} \mathcal{H}(\mu, \nu_0) \right] \quad \forall \tau > 0. \quad (6.13)$$

### 6.2.3 Proof of Theorem 4.1

*Proof.* Below,  $\mathfrak{C}$  will refer to a constant independent on  $K$  and  $p$ , whose value may change from line to line. The proof is decomposed into three steps.

Step 1: Decomposition of the excess risk with the convexity of  $J$ .

We denote by  $(\mathfrak{F}_k)_{k \geq 0}$  the natural canonical filtration associated to the sequence of random variables  $(Z^k)_{k \geq 0}$ . The next upper bound is a consequence of the relationship (6.7) and Propositions 6.1, 6.2. We have

$$\begin{aligned} J(\nu_k) - J(\mu^*) &\leq \frac{\mathcal{H}(\nu_k, \nu_k^\varepsilon) - \mathcal{H}(\nu_{k+1}, \nu_{k+1}^\varepsilon)}{\alpha} + \mathfrak{A}_k \|\nu_k^\varepsilon - \mu^*\|_{BL}^* \\ &\quad + \frac{1}{\alpha} \sum_{j=1}^p \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right] + \sum_{j=1}^p \omega_j^\varepsilon \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}). \end{aligned}$$

We then use a telescopic sum argument and obtain that:

$$\begin{aligned} \sum_{k=0}^K (J(\nu_k) - J(\mu^*)) &\leq \frac{\mathcal{H}(\nu_0, \nu_0^\varepsilon)}{\alpha} + \sum_{k=0}^K \mathfrak{A}_k \|\nu_k^\varepsilon - \mu^*\|_{BL}^* \\ &\quad + \frac{1}{\alpha} \sum_{k=0}^K \sum_{j=1}^p \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right] + \sum_{k=0}^K \sum_{j=1}^p \omega_j^\varepsilon \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}). \end{aligned}$$

Finally, using the convexity of  $J$ , the Cesaro average defined by:

$$\bar{\nu}_K = \frac{1}{K+1} \sum_{k=0}^K \nu_k, \quad (6.14)$$



satisfies:

$$\begin{aligned}
J(\bar{\nu}_K) - J(\mu^*) &\leq \frac{\mathcal{H}(\nu_0, \nu_0^\varepsilon)}{\alpha(K+1)} + \frac{\sum_{k=0}^K \mathfrak{A}_k \|\nu_k^\varepsilon - \mu^*\|_{BL}^*}{K+1} \\
&\quad + \frac{\sum_{k=0}^K \sum_{j=1}^p \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right]}{\alpha(K+1)} + \frac{\sum_{k=0}^K \sum_{j=1}^p \omega_j^\varepsilon \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})}{K+1} \\
&\leq \frac{\mathcal{H}(\nu_0, \nu_0^\varepsilon)}{\alpha K} + \frac{\sum_{k=0}^K \mathfrak{A}_k \left[ \|\nu_0^\varepsilon - \mu^*\|_{BL}^* + \sum_{\ell=0}^{k-1} \|\nu_{\ell+1}^\varepsilon - \nu_\ell^\varepsilon\|_{BL}^* \right]}{K} \\
&\quad + \frac{\sum_{k=0}^K \sum_{j=1}^p \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right]}{\alpha K} + \frac{\sum_{k=0}^K \sum_{j=1}^p \omega_j^\varepsilon \xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})}{K}
\end{aligned}$$

where we used the triangle inequality on the telescopic decomposition

$$\nu_k^\varepsilon - \mu^* = (\nu_k^\varepsilon - \nu_{k-1}^\varepsilon) + (\nu_{k-1}^\varepsilon - \nu_{k-2}^\varepsilon) + \dots + (\nu_0^\varepsilon - \mu^*).$$

We then take the expectation and use in particular a standard conditional expectation argument. Since  $\mathbb{E}[\xi_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) | \mathfrak{F}_k] = 0$ , we deduce that:

$$\begin{aligned}
\mathbb{E}[J(\bar{\nu}_K) - J(\mu^*)] &\leq \frac{\mathcal{H}(\nu_0, \nu_0^\varepsilon)}{\alpha(K+1)} + \underbrace{\|\nu_0^\varepsilon - \mu^*\|_{BL}^* \frac{\sum_{k=0}^K \mathbb{E}[\mathfrak{A}_k]}{K+1}}_{:=A_1} + \underbrace{\frac{\sum_{k=0}^K \mathbb{E} \left[ \mathfrak{A}_k \sum_{\ell=0}^{k-1} \|\nu_{\ell+1}^\varepsilon - \nu_\ell^\varepsilon\|_{BL}^* \right]}{K+1}}_{:=A_2} \\
&\quad + \underbrace{\frac{\sum_{k=0}^K \sum_{j=1}^p \mathbb{E} \left[ \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right] \right]}{\alpha(K+1)}}_{:=A_3}. \tag{6.15}
\end{aligned}$$

Step 2a: Study of  $A_1$ . We use the definition of  $\mathfrak{A}_k$  in Equation (6.12) and observe that  $\mathfrak{A}_k \leq \mathfrak{C}(1 + \|\nu_k\|_{TV})$ . We then use Proposition 4.1 and conclude that:

$$\mathbb{E}[A_1] = \|\nu_0^\varepsilon - \mu^*\|_{BL}^* \frac{\sum_{k=0}^K \mathbb{E}[\mathfrak{A}_k]}{K+1} \leq \mathfrak{C} R_0 \|\nu_0^\varepsilon - \mu^*\|_{BL}^*. \tag{6.16}$$

Step 2b: Study of  $A_2$ . We focus on the shadow sequence that involves  $\nu_{\ell+1}^\varepsilon - \nu_\ell^\varepsilon$  and observe that:

$$\begin{aligned}
\|\nu_{\ell+1}^\varepsilon - \nu_\ell^\varepsilon\|_{BL}^* &= \sup_{\|\psi\|_{BL} \leq 1} \int_{\mathcal{X}} \psi(t) d[\nu_{\ell+1}^\varepsilon - \nu_\ell^\varepsilon](t) \\
&= \sup_{\|\psi\|_{BL} \leq 1} \sum_{j=1}^p \omega_j^\varepsilon [\psi(\mathbf{t}_j^{\ell+1}) - \psi(\mathbf{t}_j^\ell)] \\
&\leq \sum_{j=1}^p \omega_j^\varepsilon \|\mathbf{t}_j^{\ell+1} - \mathbf{t}_j^\ell\|_{\mathcal{X}} \\
&= \eta \sum_{j=1}^p \omega_j^\varepsilon \|\mathbf{D}_{\nu_\ell}(\mathbb{T}^\ell, Z^{\ell+1})_j\|_{\mathcal{X}} \\
&\leq C\eta \sum_{j=1}^p \omega_j^\varepsilon (1 + \|\nu_\ell\|_{TV}) \\
&= C\eta \|\nu_0^\varepsilon\|_{TV} (1 + \|\nu_\ell\|_{TV}). \tag{6.17}
\end{aligned}$$

where we used the almost sure upper bound in Proposition C.1. A simple sum yields:

$$\begin{aligned}
\frac{\sum_{k=0}^K \mathfrak{A}_k \sum_{\ell=0}^{k-1} \|\nu_{\ell+1}^\varepsilon - \nu_\ell^\varepsilon\|_{BL}^*}{K+1} &\leq \mathfrak{C}\eta \|\nu_0^\varepsilon\|_{TV} \frac{\sum_{k=0}^K \mathfrak{A}_k \sum_{\ell=1}^k (1 + \|\nu_\ell\|_{TV})}{K+1} \\
&\leq \mathfrak{C}\eta \|\nu_0^\varepsilon\|_{TV} \frac{\sum_{k=0}^K (1 + \|\nu_k\|_{TV}) \sum_{\ell=1}^k (1 + \|\nu_\ell\|_{TV})}{K+1},
\end{aligned}$$

for  $\mathfrak{C}$  large enough. Again, we apply Proposition 4.1 and compute the expectation of the previous term. We then observe that:

$$\mathbb{E}[A_2] = \frac{\sum_{k=0}^K \mathbb{E} \left[ \mathfrak{A}_k \sum_{\ell=0}^{k-1} \|\nu_{\ell+1}^\varepsilon - \nu_\ell^{\varepsilon}\|_{BL}^* \right]}{K+1} \leq \mathfrak{C} R_0^2 \eta \|\nu_0^\varepsilon\|_{TV} K. \quad (6.18)$$

Step 2c: Study of  $A_3$ . This last term deserves a specific study. We use a conditional expectation argument and observe that, for any  $k \in \{0, \dots, K\}$ ,  $j \in \{1, \dots, p\}$ ,

$$\mathbb{E} \left[ \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right] \middle| \mathfrak{F}_k \right] = \omega_j^k \mathbb{E} \left[ \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right] \middle| \mathfrak{F}_k \right].$$

We then apply Proposition C.2 and get:

$$\mathbb{E} \left[ \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right] \middle| \mathfrak{F}_k \right] \leq \mathfrak{C} \omega_j^k \alpha^2 (1 + \|\nu_k\|_{TV}^2)$$

We then sum the previous upper bounds from 0 to  $K$  with a global expectation and Proposition 4.1. We obtain that:

$$\mathbb{E}[A_3] = \frac{\sum_{k=0}^K \sum_{j=1}^p \mathbb{E} \left[ \omega_j^k \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \right] \right]}{\alpha(K+1)} \leq \mathfrak{C} \alpha R_0^3. \quad (6.19)$$

Step 3: End of the proof.

We gather Equations (6.16), (6.18), (6.19) and obtain that:

$$\mathbb{E} [J(\bar{\nu}_K) - J(\mu^*)] \leq \frac{\mathcal{H}(\nu_0, \nu_0^\varepsilon)}{\alpha K} + \mathfrak{C} R_0 \|\nu_0^\varepsilon - \mu^*\|_{BL}^* + \mathfrak{C} R_0^2 [\eta \|\nu_0^\varepsilon\|_{TV} K + \alpha R_0].$$

We then use the definition of  $\mathcal{Q}$  given in Equation (6.13) and observe that if  $\nu_0^\varepsilon$  is chosen in an optimal way, as given in Proposition C.3, then:

$$\begin{aligned} \mathbb{E} [J(\bar{\nu}_K) - J(\mu^*)] &\leq R_0 \mathcal{Q}_{\mu^*, \nu_0}(\alpha K R_0) + \mathfrak{C} R_0^2 [\eta \|\mu^*\|_{TV} K + \alpha R_0] \\ &\leq \mathfrak{C} \|\mu^*\|_{TV} \left[ \frac{d \left( 1 + \log \frac{\alpha K R_0}{2d} + \frac{\log |\mathcal{X}|}{d} \right)}{\alpha K} + \frac{\alpha R_0^3}{\|\mu^*\|_{TV}} + R_0^2 \eta K \right], \end{aligned}$$

where we have used the relationship  $\|\nu_0^\varepsilon\|_{TV} = \|\mu^*\|_{TV}$ . The choices

$$\alpha = \sqrt{\frac{d \|\mu^*\|_{TV}}{R_0^3 K}} \quad \text{and} \quad \eta = \sqrt{\frac{d R_0}{K^3 \|\mu^*\|_{TV}}}$$

then leads to

$$\mathbb{E} [J(\bar{\nu}_K) - J(\mu^*)] \leq \mathfrak{C} \sqrt{\frac{d \|\mu^*\|_{TV} R_0^3}{K}} \left[ \log(d \|\mu^*\|_{TV} R_0^3 K) + \frac{\log(|\mathcal{X}|)}{d} \right].$$

□

### 6.3 Proof of Theorem 4.2

*Proof.* The proof is splitted into three parts, and relies on a contraction argument with conditional expectation. In what follows, we will choose  $\alpha$  such that  $\alpha \mathcal{C}_1 (R_0 + 1) < 1$  where  $\mathcal{C}_1$  is involved in Lemma B.1, and we shall use  $|e^h - 1| \leq 2|h|$  which is valid when  $|h| \leq 1$ . We will frequently apply this inequality with  $h = -\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})$ .

Step 1: One-step evolution and second order term. Let  $k \in \mathbb{N}^*$  be fixed. According to Proposition B.1:

$$J(\nu_{k+1}) - J(\nu_k) = \int_{\mathcal{X}} J'_{\nu_k} d(\nu_{k+1} - \nu_k) + \frac{1}{2} \|\Phi(\nu_{k+1} - \nu_k)\|_{\mathbb{H}}^2.$$

Introducing the measure  $\tilde{\nu}_{k+1} = \sum_{j=1}^p \omega_j^{k+1} \delta_{\mathbf{t}_{j,k}}$ , we deduce that:

$$\begin{aligned} \|\Phi(\nu_{k+1} - \nu_k)\|_{\mathbb{H}}^2 &= \|\Phi(\nu_{k+1} - \tilde{\nu}_{k+1} + \tilde{\nu}_{k+1} - \nu_k)\|_{\mathbb{H}}^2, \\ &\leq 2\|\Phi(\nu_{k+1} - \tilde{\nu}_{k+1})\|_{\mathbb{H}}^2 + 2\|\Phi(\tilde{\nu}_{k+1} - \nu_k)\|_{\mathbb{H}}^2 \end{aligned}$$

First remark that,

$$\begin{aligned} \|\Phi(\nu_{k+1} - \tilde{\nu}_{k+1})\|_{\mathbb{H}}^2 &= \left\| \sum_{j=1}^p \omega_j^{k+1} (\varphi_{\mathbf{t}_j^{k+1}} - \varphi_{\mathbf{t}_j^k}) \right\|_{\mathbb{H}}^2, \\ &\leq \sum_{j=1}^p \omega_j^{k+1} \times \sum_{j=1}^p \omega_j^{k+1} \|\varphi_{\mathbf{t}_j^{k+1}} - \varphi_{\mathbf{t}_j^k}\|_{\mathbb{H}}^2, \\ &\leq Lip(\varphi) \|\nu_{k+1}\|_{\text{TV}} \sum_{j=1}^p \omega_j^{k+1} \|\mathbf{t}_j^{k+1} - \mathbf{t}_j^k\|_{\mathcal{X}}^2. \end{aligned}$$

We then use the update of  $\mathbb{T}^k$  to  $\mathbb{T}^{k+1}$  that is given in (2.15) and obtain that:

$$\|\Phi(\nu_{k+1} - \tilde{\nu}_{k+1})\|_{\mathbb{H}}^2 \leq Lip(\varphi) \|\nu_{k+1}\|_{\text{TV}} \eta^2 \sum_{j=1}^p \omega_j^k \|D_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})\|_{\mathcal{X}}^2. \quad (6.20)$$

In the same time, using (2.14), we obtain that:

$$\begin{aligned} \|\Phi(\tilde{\nu}_{k+1} - \nu_k)\|_{\mathbb{H}}^2 &= \left\| \sum_{j=1}^p (\omega_j^{k+1} - \omega_j^k) \varphi_{\mathbf{t}_{j,k}} \right\|_{\mathbb{H}}^2, \\ &= \left\| \sum_{j=1}^p \omega_j^k (e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1) \varphi_{\mathbf{t}_{j,k}} \right\|_{\mathbb{H}}^2, \\ &\leq \sum_{j=1}^p \omega_j^k \times \sum_{j=1}^p \omega_j^k (e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1)^2 \|\varphi_{\mathbf{t}_{j,k}}\|_{\mathbb{H}}^2, \end{aligned}$$

where the last line comes from the Jensen inequality. We then observe from Lemma B.1 that  $J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})$  is a bounded term so that if  $\alpha$  is chosen such that  $\alpha \mathcal{C}_1(R_0 + 1) < 1$ , then a constant  $C_\varphi$  large enough exists such that:

$$\|\Phi(\tilde{\nu}_{k+1} - \nu_k)\|_{\mathbb{H}}^2 \leq C_\varphi \|\nu_k\|_{\text{TV}} \alpha^2 \sum_{j=1}^p \omega_j^k |J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})|^2 \quad (6.21)$$

Gathering Equations (6.20) and (6.21), we then deduce that

$$\|\Phi(\nu_{k+1} - \nu_k)\|_{\mathbb{H}}^2 \leq C_\varphi \left( \eta^2 \|\nu_{k+1}\|_{\text{TV}} \sum_{j=1}^p \omega_j^k \|D_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})\|^2 + \alpha^2 \|\nu_k\|_{\text{TV}} \sum_{j=1}^p \omega_j^k |J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})|^2 \right) \quad (6.22)$$

Step 2: Study of the drift first order term. We expand the first order term and observe that:

$$\begin{aligned}
\int_{\mathcal{X}} J'_{\nu_k} d(\nu_{k+1} - \nu_k) &= \sum_{j=1}^p \left[ (\omega_j^{k+1} - \omega_j^k) J'_{\nu_k}(\mathbf{t}_j^k) + \omega_j^k (J'_{\nu_k}(\mathbf{t}_j^{k+1}) - J'_{\nu_k}(\mathbf{t}_j^k)) \right] \\
&+ \sum_{j=1}^p (\omega_j^{k+1} - \omega_j^k) (J'_{\nu_k}(\mathbf{t}_j^{k+1}) - J'_{\nu_k}(\mathbf{t}_j^k)) \\
&= \sum_{j=1}^p \left[ (\omega_j^{k+1} - \omega_j^k) J'_{\nu_k}(\mathbf{t}_j^k) + \omega_j^k \langle \mathbf{t}_j^{k+1} - \mathbf{t}_j^k, \nabla J'_{\nu_k}(\mathbf{t}_j^k) \rangle \right] \\
&+ \sum_{j=1}^p \left[ \omega_j^k \langle \mathbf{t}_j^{k+1} - \mathbf{t}_j^k, \nabla^2 J'_{\nu_k}(\mathbf{v}_j^k)(\mathbf{t}_j^{k+1} - \mathbf{t}_j^k) \rangle + (\omega_j^{k+1} - \omega_j^k) \langle \nabla J'_{\nu_k}(\tilde{\mathbf{v}}_j^k), \mathbf{t}_j^{k+1} - \mathbf{t}_j^k \rangle \right],
\end{aligned}$$

where  $\mathbf{v}_j^k$  and  $\tilde{\mathbf{v}}_j^k$  are some auxiliary points that belong to  $(\mathbf{t}_j^k, \mathbf{t}_j^{k+1})$  obtained with the help of first and second order Taylor expansions. Using Proposition C.1, we deduce that:

$$\begin{aligned}
\int_{\mathcal{X}} J'_{\nu_k} d(\nu_{k+1} - \nu_k) &\leq \sum_{j=1}^p \left[ (\omega_j^{k+1} - \omega_j^k) J'_{\nu_k}(\mathbf{t}_j^k) + \omega_j^k \langle \mathbf{t}_j^{k+1} - \mathbf{t}_j^k, \nabla J'_{\nu_k}(\mathbf{t}_j^k) \rangle \right] \\
&+ \sum_{j=1}^p \omega_j^k [Lip(\nabla\varphi)\|\nu_k\|_{TV} + Lip(\nabla\mathbf{y})] \|\mathbf{t}_j^{k+1} - \mathbf{t}_j^k\|_{\mathcal{X}}^2 \\
&+ \sum_{j=1}^p [Lip(\nabla\varphi)\|\nu_k\|_{TV} + Lip(\nabla\mathbf{y})] (\omega_j^{k+1} - \omega_j^k) \|\mathbf{t}_j^{k+1} - \mathbf{t}_j^k\|_{\mathcal{X}}.
\end{aligned}$$

Using the total variation upper bound stated in Proposition 4.1 by  $R_0$ , we then define for the sake of readability the constant  $A$  as:

$$A = Lip(\varphi)R_0 + Lip(\mathbf{y}) \vee Lip(\nabla\varphi)R_0 + Lip(\nabla\mathbf{y}).$$

We then derive:

$$\begin{aligned}
\int_{\mathcal{X}} J'_{\nu_k} d(\nu_{k+1} - \nu_k) &\leq \sum_{j=1}^p (\omega_j^{k+1} - \omega_j^k) J'_{\nu_k}(\mathbf{t}_j^k) + \omega_j^k \langle \mathbf{t}_j^{k+1} - \mathbf{t}_j^k, \nabla J'_{\nu_k}(\mathbf{t}_j^k) \rangle \\
&+ A \sum_{j=1}^p \omega_j^k \|\mathbf{t}_j^{k+1} - \mathbf{t}_j^k\|^2 + (\omega_j^{k+1} - \omega_j^k) \|\mathbf{t}_j^{k+1} - \mathbf{t}_j^k\|.
\end{aligned}$$

We now use the surrogate update stated in Algorithm 1 into the previous inequality and obtain that:

$$\begin{aligned}
\int_{\mathcal{X}} J'_{\nu_k} d(\nu_{k+1} - \nu_k) &\leq \sum_{j=1}^p \omega_j^k (e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1) J'_{\nu_k}(\mathbf{t}_j^k) - \eta \omega_j^k \langle D_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}), \nabla J'_{\nu_k}(\mathbf{t}_j^k) \rangle \\
&+ A \sum_{j=1}^p \eta^2 \omega_j^k \|D_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})\|_{\mathcal{X}}^2 + \eta \omega_j^k (e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1) \|D_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})\|_{\mathcal{X}}.
\end{aligned}$$

We then consider the conditional expectation at time  $k$  and apply Proposition C.1 (to upper bound some rest terms) and Proposition C.2 (to control the drift at iteration  $k$ ). We deduce that a large enough  $\mathfrak{C}$  such

that:

$$\begin{aligned}
\mathbb{E} \left[ \int_{\mathcal{X}} J'_{\nu_k} d(\nu_{k+1} - \nu_k) \middle| \mathfrak{F}_k \right] &\leq -\alpha \sum_{j=1}^p \omega_j^k J'_{\nu_k}(\mathbf{t}_j^k)^2 + \mathfrak{C} \alpha^2 \sum_{j=1}^p \omega_j^k (1 + \|\nu_k\|_{\text{TV}})^3 - \eta \sum_{j=1}^p \omega_j^k \|\nabla J'_{\nu_k}(\mathbf{t}_j^k)\|^2 \\
&\quad + \mathfrak{C} \eta^2 \sum_{j=1}^p \omega_j^k (1 + \|\nu_k\|_{\text{TV}})^2 + \mathfrak{C} \eta \alpha \sum_{j=1}^p \omega_j^k (1 + \|\nu_k\|_{\text{TV}})^3 \\
&\leq -\alpha \sum_{j=1}^p \omega_j^k J'_{\nu_k}(\mathbf{t}_j^k)^2 - \eta \sum_{j=1}^p \omega_j^k \|\nabla J'_{\nu_k}(\mathbf{t}_j^k)\|^2 \\
&\quad + \mathfrak{C} \|\nu_k\|_{\text{TV}} (1 + \|\nu_k\|_{\text{TV}})^3 (\alpha^2 + \eta^2),
\end{aligned}$$

where in the last line we used the Young inequality  $2\eta\alpha \leq \alpha^2 + \eta^2$  and some rough upper bounds on the rest terms. We now associate this last inequality with Equations (6.22) and obtain the descent property:

$$\mathbb{E} [J(\nu_{k+1}) \middle| \mathfrak{F}_k] \leq J(\nu_k) - \alpha \|J'_{\nu_k}\|_{\nu_k}^2 - \eta \|\nabla J'_{\nu_k}\|_{\nu_k}^2 + \mathfrak{C}(1 + R_0^4)(\alpha^2 + \eta^2) \quad (6.23)$$

Step 3: Conclusion of the proof. The rest of the proof proceeds with a standard argument. We use a telescopic sum + conditional expectation strategy and observe that:

$$\alpha \sum_{k=1}^K \mathbb{E}[\|J'_{\nu_k}\|_{\nu_k}^2] + \eta \sum_{k=1}^K \mathbb{E}[\|\nabla J'_{\nu_k}\|_{\nu_k}^2] \leq J(\nu_0) + \mathfrak{C}(1 + R_0^4)(K\alpha^2 + K\eta^2)$$

Choosing  $\alpha = \eta$ , we deduce that:

$$\frac{1}{K} \sum_{k=1}^K (\mathbb{E}[\|J'_{\nu_k}\|_{\nu_k}^2] + \mathbb{E}[\|\nabla J'_{\nu_k}\|_{\nu_k}^2]) \leq \frac{J(\nu_0)}{\alpha K} + \mathfrak{C}(1 + R_0^4)\alpha.$$

Finally, if  $\tau_K$  refers to a random variable uniformly distributed over  $\{1, \dots, K\}$ , independent from the sequence  $(\nu_k)_{k \geq 1}$ , the tuning  $\alpha = \eta = 1/\sqrt{K}$  yields:

$$\mathbb{E} \left[ \|J'_{\nu_{\tau_K}}\|_{\nu_{\tau_K}}^2 + \|\nabla J'_{\nu_{\tau_K}}\|_{\nu_{\tau_K}}^2 \right] \leq \frac{J(\nu_0) + \mathfrak{C}(1 + R_0^4)}{\sqrt{K}}.$$

□

**Acknowledgements** The authors would like to thank N. Jouvin for its remarks and time on the numerical aspects of this project. They are also in debt with L. Chizat for valuable discussions on CPGD during a seminar at Institut Henri Poincaré.

## References

- Jean-Marc Azais, Yohann De Castro, and Fabrice Gamboa. Spike detection from inaccurate samplings. *Applied and Computational Harmonic Analysis*, 38(2):177–195, 2015.
- Francis Bach and Lénaïc Chizat. Gradient descent on infinitely wide neural networks: Global convergence and generalization. *arXiv preprint arXiv:2110.08084*, 2021.
- Jérôme Bolte, Laurent Miclo, and Stéphane Villeneuve. Swarm gradient dynamics for global optimization: the mean-field limit case. *Mathematical Programming, to appear*, 2023.
- Claire Boyer, Antonin Chambolle, Yohann De Castro, Vincent Duval, Frédéric De Gournay, and Pierre Weiss. On representer theorems and convex regularization. *SIAM Journal on Optimization*, 29(2):1260–1281, 2019.
- Kristian Bredies and Hanna Katriina Pikkarainen. Inverse problems in spaces of measures. *ESAIM: Control, Optimisation and Calculus of Variations*, 19(01):190–218, 2013.
- Haim Brézis. *Functional analysis, Sobolev spaces and partial differential equations*, volume 2. Springer, 2011.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Emmanuel J. Candès and Carlos Fernandez-Granda. Super-resolution from noisy data. *Journal of Fourier Analysis and Applications*, 19(6):1229–1254, 2013.
- Émmanuel J Candès and Carlos Fernandez-Granda. Towards a mathematical theory of super-resolution. *Communications on pure and applied Mathematics*, 67(6):906–956, 2014.
- Lénaïc Chizat. Sparse optimization on measures with over-parameterized gradient descent. *Mathematical Programming*, 194(1-2):487–532, 2022.
- Lénaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- Yohann De Castro and Fabrice Gamboa. Exact reconstruction using beurling minimal extrapolation. *Journal of Mathematical Analysis and applications*, 395(1):336–354, 2012.
- Yohann De Castro, Sébastien Gadat, Clément Marteau, and C Maugis-Rabusseau. Supermix: sparse regularization for mixtures. *The Annals of Statistics*, 49(3):1779–1809, 2021.
- Bernard Delyon, Marc Lavielle, and Eric Moulines. Convergence of a stochastic approximation version of the em algorithm. *The Annals of Statistics*, 27(1):94–128, 1999.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- Quentin Denoyelle, Vincent Duval, Gabriel Peyré, and Emmanuel Soubies. The sliding frank-wolfe algorithm and its application to super-resolution microscopy. *Inverse Problems*, 2019.
- David L. Donoho. Superresolution via sparsity constraints. *SIAM Journal on Mathematical Analysis*, 23(5):1309–1331, 1992.
- Vincent Duval and Gabriel Peyré. Exact support recovery for sparse spikes deconvolution. *Foundations of Computational Mathematics*, 15(5):1315–1355, 2015.
- Carlos Fernandez-Granda. Support detection in super-resolution. In *The 10th International Conference on Sampling Theory and Applications (SampTA 2013)*, pages 145–148, 2013.

- Fabrice Gamboa and E Gassiat. Sets of superresolution and the maximum entropy method on the mean. *SIAM journal on mathematical analysis*, 27(4):1129–1152, 1996.
- Evarist Giné and Richard Nickl. *Mathematical foundations of infinite-dimensional statistical models*. Cambridge university press, 2021.
- Bernd Hofmann, Barbara Kaltenbacher, Christiane Poeschl, and Otmar Scherzer. A convergence rates result for tikhonov regularization in banach spaces with non-smooth operators. *Inverse Problems*, 23(3): 987, 2007.
- Nicolas Keriven, Anthony Bourrier, Rémi Gribonval, and Patrick Pérez. Sketching for large-scale learning of mixture models. *Information and Inference: A Journal of the IMA*, 7(3):447–508, 2018.
- Frederik Kunstner, Raunak Kumar, and Mark Schmidt. Homeomorphic-invariance of em: Non-asymptotic convergence in kl divergence for exponential families via mirror descent. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 3295–3303. PMLR, 13–15 Apr 2021.
- Guanghui Lan, Arkadij Semenovič Nemirovskij, and Alexander Shapiro. Validation analysis of mirror descent stochastic approximation method. *Mathematical programming*, 134(2):425–458, 2012.
- Laurent Miclo. On the convergence of global-optimization fraudulent stochastic algorithms. *Preprint*, 2023.
- Krikamol Muandet, Kenji Fukumizu, Bharath Sriperumbudur, Bernhard Schölkopf, et al. Kernel mean embedding of distributions: A review and beyond. *Foundations and Trends® in Machine Learning*, 10(1-2):1–141, 2017.
- Arkadij Semenovič Nemirovskij and David Borisovich Yudin. *Problem complexity and method efficiency in optimization*. Wiley-Interscience, 1983.
- Clarice Poon, Nicolas Keriven, and Gabriel Peyré. The geometry of off-the-grid compressed sensing. *Foundations of Computational Mathematics*, pages 1–87, 2021.
- Walter Rudin. *Real and complex analysis*. Mcgraw hill International Book Company, 1974.
- Ingo Steinwart and Andreas Christmann. *Support vector machines*. Springer Science & Business Media, 2008.
- Gongguo Tang, Badri N. Bhaskar, Parikshit Shah, and Benjamin Recht. Compressed sensing off the grid. *Information Theory, IEEE Transactions on*, 59(11):7465–7490, 2013.
- Gongguo Tang, Badri N. Bhaskar, and Benjamin Recht. Near minimax line spectral estimation. *Information Theory, IEEE Transactions on*, 61(1):499–512, 2015.

# Appendix of "FastPart: Over-Parameterized Stochastic Gradient Descent for Sparse optimisation on Measures"

## A Reformulations, proofs and technical lemmas

### A.1 The non-separable case and the interesting reformulation of the objective

There is a little subtlety on the properties that one should require on  $\mathbb{H}$ . At first glance, we need  $\mathbb{H}$  separable to prove Bochner integrability in Lemma A.2. But  $K$  is continuous on a compact space  $\mathcal{X}$ , hence its RKHS is separable [Steinwart and Christmann, 2008, Lemma 4.3.3]. This RKHS is isometric to a separable subspace of  $\mathbb{H}$  as proven by the next lemma.

**Lemma A.1.** *Let  $\mathbb{H}$  be Hilbert space and let  $\mathcal{X}$  be compact space. Under  $(\mathbf{A}_0)$ , there exists a separable closed vector subspace  $(\mathbb{H}_{\mathcal{F}}, \|\cdot\|_{\mathbb{H}})$  of  $\mathbb{H}$  which is isometric to  $(\mathcal{F}, \|\cdot\|_{\mathcal{F}})$ , the RKHS defined by  $K$ . Denote by  $\Pi$  the orthogonal projection onto  $\mathbb{H}_{\mathcal{F}}$  then for any  $\nu \in \mathcal{M}(\mathcal{X})$*

$$J(\nu) = \frac{1}{2} \|\mathbf{y} - \Pi(\mathbf{y})\|_{\mathbb{H}}^2 + \frac{1}{2} \|\Pi(\mathbf{y}) - \Phi(\nu)\|_{\mathbb{H}_{\mathcal{F}}}^2 + \lambda \|\nu\|_{\text{TV}}. \quad (\text{A.1})$$

*Proof.* We denote by  $\mathcal{F}$  the RKHS defined by  $K$ . By [Steinwart and Christmann, 2008, Theorem 4.21], one has that

$$\mathcal{F} := \left\{ f : \mathcal{X} \rightarrow \mathbb{H} : \exists h \in \mathbb{H}, \forall t \in \mathcal{X}, f(t) = \langle h, \Phi(t) \rangle_{\mathbb{H}} \right\}.$$

is the only RKHS defined by  $K$  and

$$\|f\|_{\mathcal{F}} = \inf \left\{ \|h\|_{\mathbb{H}} : h \in \mathbb{H} \text{ s.t. } f(\cdot) = \langle h, \Phi(\cdot) \rangle_{\mathbb{H}} \right\},$$

Consider the functions

$$f_h : x \mapsto \langle h, \Phi(x) \rangle_{\mathbb{H}}, \quad h = \sum_{j=1}^r \omega_j \Phi(t_j)$$

defining the pre-dual of  $\mathcal{F}$ . Observe that

$$\langle f_{h_1}, f_{h_2} \rangle_{\mathcal{F}} = \langle h_1, h_2 \rangle_{\mathbb{H}},$$

where we denote by  $\langle \cdot, \cdot \rangle_{\mathcal{F}}$ , the dot product of  $\mathcal{F}$ . Define  $\mathbb{H}_{\mathcal{F}}$  the vector subspace of  $\mathbb{H}$  defined as the closure (in  $\mathbb{H}$ ) of the span of  $\Phi(\mathcal{X})$ . The aforementioned equality shows that  $h \mapsto f_h$  is an isometry from  $(\mathbb{H}_{\mathcal{F}}, \|\cdot\|_{\mathbb{H}})$  onto  $(\mathcal{F}, \|\cdot\|_{\mathcal{F}})$ . Since  $\mathcal{F}$  is separable [Steinwart and Christmann, 2008, Lemma 4.3.3], we deduce that  $\mathbb{H}_{\mathcal{F}}$  is separable. The last statement is a consequence of the Pythagorean theorem.  $\square$

Note that in (A.1), the term  $\frac{1}{2} \|\mathbf{y} - \Pi(\mathbf{y})\|_{\mathbb{H}}^2$  is constant. Hence, up to a constant term, and without loss of generality, one can assume that  $\mathbb{H}$  is separable.

**Remark A.1.** *The proof of Lemma A.1 is a consequence of [Steinwart and Christmann, 2008, Theorem 4.21] and we choose to maintain it in this paper for sake of completeness. Moreover, it sheds light on an interesting reformulation of the quadratic term in (1.5), the objective  $J$ . Indeed, it holds*

$$J(\nu) = \frac{1}{2} \left\| (\sqsupset \circ \Pi)(\mathbf{y}) - (\sqsupset \circ \Phi)(\nu) \right\|_{\mathcal{F}}^2 + \lambda \|\nu\|_{\text{TV}}, \quad (\text{A.2})$$

up to a constant term and where  $\sqsupset$  denotes the isometry between  $(\mathbb{H}_{\mathcal{F}}, \|\cdot\|_{\mathbb{H}})$  and  $(\mathcal{F}, \|\cdot\|_{\mathcal{F}})$ .



## A.2 Existence of the kernel measure embedding

Kernel mean embedding is a standard notion in Machine Learning, see for instance Muandet et al. [2017]. Extending this notion of measure with finite total variation norm is straightforward. We referred to this notion as *Kernel Measure Embedding* as the two notions coincides on probability measures.

**Lemma A.2.** *Let  $\mathbb{H}$  be separable Hilbert space and let  $\mathcal{X}$  be compact metric space. Under  $(\mathbf{A}_0)$ , the operator  $\Phi$  defined by (1.3) is well defined and bounded linear as a function from  $\mathcal{M}(\mathcal{X})$  to  $\mathbb{H}$ . Furthermore, the dual of  $\Phi$  is given by*

$$\Phi^* : h \in \mathbb{H} \mapsto (t \mapsto \langle h, \varphi_t \rangle_{\mathbb{H}}) \in (\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty}), \quad (\text{A.3})$$

and for any  $(h, \nu) \in \mathbb{H} \times \mathcal{M}(\mathcal{X})$ ,

$$\langle h, \Phi(\nu) \rangle_{\mathbb{H}} = \int_{\mathcal{X}} \langle h, \varphi_t \rangle_{\mathbb{H}} d\nu(t) = \langle \Phi^*(h), \nu \rangle_{\mathcal{C}(\mathcal{X}), \mathcal{M}(\mathcal{X})} \leq \sup_t \sqrt{\mathbb{K}(t, t)} \|h\|_{\mathbb{H}} \|\nu\|_{\text{TV}}. \quad (\text{A.4})$$

**Remark A.2.** *A key result of Lemma A.2 is that  $\text{Im}(\Phi^*) \subseteq (\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty})$ , this latter being a subset of  $\mathcal{M}(\mathcal{X})^*$ , the topological dual of  $\mathcal{M}(\mathcal{X})$ .*

*Proof.* Let  $\nu \in \mathcal{M}(\mathcal{X})$ . We say that  $t \in \mathcal{X} \mapsto f(t) \in \mathbb{H}$  is *simple* if it is finitely valued, namely

$$f(t) = \sum_{i=1}^n h_i \mathbf{1}_{\{t \in B_i\}},$$

for some  $n \geq 1$ ,  $h_i \in \mathbb{H}$ , and  $B_i$  Borel set of  $\mathcal{X}$ . In this case, one has

$$\int_{\mathcal{X}} f d\nu = \sum_{i=1}^n h_i \nu(B_i).$$

Note that  $\|\varphi_t\|_{\mathbb{H}} = \sqrt{\mathbb{K}(t, t)}$  and, it holds that

$$\int_{\mathcal{X}} \|\varphi_t\|_{\mathbb{H}} d\nu(t) \leq \sup_t \sqrt{\mathbb{K}(t, t)} \|\nu\|_{\text{TV}} < \infty, \quad (\text{A.5})$$

using the fact that  $t \in \mathcal{X} \mapsto \sqrt{\mathbb{K}(t, t)}$  is a bounded continuous function by  $(\mathbf{A}_0)$ . We emphasize that this function not need to be vanishing at infinity.

From (A.5), we deduce that the map  $m : A \mapsto m(A) := \int_{\mathcal{X}} \|\varphi_t\|_{\mathbb{H}} \mathbf{1}_{\{\varphi_t \in A\}} d\nu(t)$  is a finite measure on the Borel sets of  $\mathbb{H}$  and hence, by Oxtoby-Ulam theorem (see [Giné and Nickl, 2021, Proposition 2.1.4] for instance), a tight Borel measure. Given  $0 < \varepsilon_n \rightarrow 0$ , let  $K_n$  be a compact set such that  $m(K_n^c) < \varepsilon_n/2$ , let  $A_{n,1}, \dots, A_{n,k_n}$  be a finite partition of  $K_n$  consisting of sets of diameter at most  $\varepsilon_n/2$ , pick up a point  $h_{n,k} \in A_{n,k}$  for each  $k$  and define the simple function

$$f_n(t) = \sum_{k=1}^{k_n} h_{n,k} \mathbf{1}_{\{\varphi_t \in A_{n,k}\}}.$$

Then

$$\int_{\mathcal{X}} \|\varphi_t - f_n(t)\|_{\mathbb{H}} d\nu(t) \leq \varepsilon_n/2 + m(K_n^c) < \varepsilon_n \rightarrow 0,$$

showing that  $t \in \mathcal{X} \mapsto \varphi_t \in \mathbb{H}$  is Bochner integrable, hence Petti's integrable, and both integrals coincide (see for instance [Giné and Nickl, 2021, Section 2.6.1]). We deduce that  $\Phi$  is well defined, using Bochner integration. Furthermore, one can deduce that

$$\left\| \int_{\mathcal{X}} \varphi_t d\nu(t) \right\|_{\mathbb{H}} \leq \int_{\mathcal{X}} \|\varphi_t\|_{\mathbb{H}} d\nu(t) \leq \sup_t \sqrt{\mathbb{K}(t, t)} \|\nu\|_{\text{TV}},$$

showing that  $\Phi$  is bounded linear.

Also, if  $h \in \mathbb{H}$  then

$$\left| \int_{\mathcal{X}} \langle h, f_n(t) \rangle_{\mathbb{H}} - \langle h, \varphi_t \rangle_{\mathbb{H}} d\nu(t) \right| \leq \|h\|_{\mathbb{H}} \int_{\mathcal{X}} \|\varphi_t - f_n(t)\|_{\mathbb{H}} d\nu(t) \rightarrow 0.$$

Hence,  $\int_{\mathcal{X}} \langle h, \varphi_t \rangle_{\mathbb{H}} d\nu(t) = \lim_n \int_{\mathcal{X}} \langle h, f_n(t) \rangle_{\mathbb{H}} d\nu(t)$  exists and is finite. We deduce that

$$\int_{\mathcal{X}} \langle h, \varphi_t \rangle_{\mathbb{H}} d\nu(t) = \langle h, \Phi(\nu) \rangle_{\mathbb{H}}, \quad (\text{A.6})$$

using that  $\int_{\mathcal{X}} \langle h, f_n(t) \rangle_{\mathbb{H}} d\nu(t) = \langle h, \int_{\mathcal{X}} f_n d\nu \rangle_{\mathbb{H}}$ .

Using (A.6) and Cauchy-Schwarz inequality, one gets that

$$\langle h, \Phi(\nu) \rangle_{\mathbb{H}} = \int_{\mathcal{X}} \langle h, \varphi_t \rangle_{\mathbb{H}} d\nu(t) \leq \sup_t \sqrt{\mathbb{K}(t, t)} \|h\|_{\mathbb{H}} \|\nu\|_{\text{TV}}, \quad (\text{A.7})$$

and hence, we can write

$$\langle h, \Phi(\nu) \rangle_{\mathbb{H}} = \langle \langle h, \varphi_t \rangle_{\mathbb{H}}, \nu \rangle_{\mathcal{M}(\mathcal{X})^*, \mathcal{M}(\mathcal{X})}$$

where  $\mathcal{M}(\mathcal{X})^*$  is the topological dual of  $\mathcal{M}(\mathcal{X})$ . It shows that the dual  $\Phi^*$  is given by  $\Phi^*(h)(t) = \langle h, \varphi_t \rangle_{\mathbb{H}}$ . As a function of  $t$ , it is clear that it is continuous by (A<sub>0</sub>) and that  $\|\Phi^*(h)\|_{\infty} \leq \sup_t \sqrt{\mathbb{K}(t, t)} \|h\|_{\mathbb{H}} < \infty$ , showing that it belongs to the space of bounded continuous functions.  $\square$

### A.3 Proof of Theorem 1.1

Let  $(\nu_n)$  be a minimizing sequence of measures of Program (1.6). Up to an extraction we can consider that  $L(\Phi(\nu_n)) + \lambda \|\nu_n\|_{\text{TV}} \leq 1 + \inf_{\nu} \{L(\Phi(\nu)) + \lambda \|\nu\|_{\text{TV}}\}$ . In particular, it holds that

$$\lambda \|\nu_n\|_{\text{TV}} \leq 1 + \inf_{\nu} \{L(\Phi(\nu)) + \lambda \|\nu\|_{\text{TV}}\}.$$

Up to an extraction, by Banach-Alaoglu theorem, we can consider that the sequence  $(\nu_n)$  converges for the weak- $\star$  topology and, up to another extraction, we can consider that  $\|\nu_n\|_{\text{TV}}$  converges.

We denote by  $\mu^* \in \mathcal{M}(\mathcal{X})$  its limit. Using [Brézis, 2011, Proposition 3.13(iii)], the TV-norm is l.s.c. for the weak- $\star$  topology, and we get that

$$\lim_n \|\nu_n\|_{\text{TV}} \geq \|\mu^*\|_{\text{TV}}.$$

Using Lemma A.2, it holds that for any  $h \in \mathbb{H}$  and for any convergent sequence  $\nu_n \rightarrow \mu^*$  for the weak- $\star$  topology,

$$\langle h, \Phi(\nu_n) \rangle_{\mathbb{H}} = \langle \Phi^*(h), \nu_n \rangle_{\mathcal{C}(\mathcal{X}), \mathcal{M}(\mathcal{X})} \rightarrow \langle \Phi^*(h), \mu^* \rangle_{\mathcal{C}(\mathcal{X}), \mathcal{M}(\mathcal{X})},$$

proving that  $\Phi$  is continuous from  $\mathcal{M}(\mathcal{X})$  weak- $\star$  to  $\mathbb{H}$  weak (see Remark A.2). Since  $L$  is l.s.c. for the weak topology of  $\mathbb{H}$  [Brézis, 2011, Corollary 3.9], we get that

$$\lim_n L(\Phi(\nu_n)) \geq L(\Phi(\mu^*)).$$

Combining the aforementioned limits, we deduce that

$$(1.6) = \lim_n \{L(\Phi(\nu_n)) + \lambda \|\nu_n\|_{\text{TV}}\} \geq L(\Phi(\mu^*)) + \lambda \|\mu^*\|_{\text{TV}} \geq (1.6),$$

hence equality. The uniqueness of  $\Phi(\mu^*)$  follows by strict convexity.

**Remark A.3.** In this paper, we assume that  $\mathcal{X}$  is compact. Some of our bounds depend on the size of  $\mathcal{X}$  and do not hold for non-compact spaces. But, the existence of  $\mu^*$  can be proven in the non-compact case.

The subtlety is in (A.3). To get the proof of Theorem 1.1 work when  $\mathcal{X}$  is a Polish space (not necessarily compact), one needs that  $\text{Im}(\Phi^*) \in (\mathcal{C}_0(\mathcal{X}), \|\cdot\|_{\infty})$ , the space of continuous functions vanishing at infinity. We already know that  $\text{Im}(\Phi^*) \in (\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty})$  by Condition (A<sub>0</sub>). We have the following result.

**Theorem A.3.** Let  $\mathbb{H}$  be Hilbert space and let  $\mathcal{X}$  be Polish space. Assume that

- Assumption (**A**<sub>0</sub>) holds;
- the RKHS  $\mathcal{F}$  (defined by  $\mathbb{K}$ ) is contained in  $\mathcal{C}_0(\mathcal{X})$ ;
- and  $\sup_t \sqrt{\mathbb{K}(t, t)} < \infty$ ;

then there exists a measure  $\mu^* \in \mathcal{M}(\mathcal{X})$  such that

$$J(\mu^*) = \min_{\mu \in \mathcal{M}(\mathcal{X})} J(\mu).$$

Furthermore, the vector  $\Phi(\mu^*) \in \mathbb{H}$  is unique.

**Remark A.4.** The same argument can be used to prove that Program (1.6) restricted to  $\mathcal{M}(\mathcal{X})_+$  admits solutions. Indeed, take  $(\nu_n)$  a sequence of nonnegative measures such that the objective converges towards the infimum. We can use the above proof to show the existence of  $\mu^*$ . The only point left to prove is that the measure  $\mu^*$  is nonnegative, which is straight forward using weak-star convergence and Riesz representation theorem [Rudin, 1974, Chapter 2] of nonnegative linear functional defined by nonnegative continuous functions with compact support (which are included in  $\mathcal{C}_0(\mathcal{X})$ ).

**Remark A.5.** A similar result can be found in [Chizat, 2022, Proposition 3.1] using Prokhorov's theorem.

## B Gradients of the objective

### B.1 In the space of (nonnegative) measures

We first consider the variation of  $J$  in  $\mathcal{M}(\mathcal{X})_+$  in terms of its Fréchet differential.

**Proposition B.1.** If  $\nu + \sigma \in \mathcal{M}(\mathcal{X})_+$  and  $\nu \in \mathcal{M}(\mathcal{X})_+$  then

$$J(\nu + \sigma) - J(\nu) = \int_{\mathcal{X}} J'_\nu d\sigma + \frac{1}{2} \|\Phi(\sigma)\|_{\mathbb{H}}^2, \quad (\text{B.1})$$

where  $J'_\nu := \Phi^*(\Phi(\nu) - \mathbf{y}) + \lambda$  and  $\Phi^* : (\mathbb{H}, \|\cdot\|_{\mathbb{H}}) \rightarrow (\mathcal{C}(\mathcal{X}), \|\cdot\|_{\infty})$  is the dual of  $\Phi$ .

*Proof of B.1.* The proof follows from the expansion of  $J(\nu + \sigma)$ :

$$\begin{aligned} J(\nu + \sigma) &= \frac{1}{2} \|\mathbf{y} - \Phi(\nu + \sigma)\|_{\mathbb{H}}^2 + \lambda \|\nu + \sigma\|_{\text{TV}} \\ &= \frac{1}{2} \|\mathbf{y} - \Phi(\nu) - \Phi(\sigma)\|_{\mathbb{H}}^2 + \lambda \|\nu + \sigma\|_{\text{TV}} \\ &= \frac{1}{2} \|\mathbf{y} - \Phi(\nu)\|_{\mathbb{H}}^2 - \langle \mathbf{y} - \Phi(\nu), \Phi(\sigma) \rangle_{\mathbb{H}} + \frac{1}{2} \|\Phi(\sigma)\|_{\mathbb{H}}^2 + \lambda \|\nu + \sigma\|_{\text{TV}} \\ &= J(\nu) - \langle \Phi^*(\mathbf{y} - \Phi(\nu)), \sigma \rangle_{\mathbb{H}} + \frac{1}{2} \|\Phi(\sigma)\|_{\mathbb{H}}^2 + \lambda [\|\nu + \sigma\|_{\text{TV}} - \|\nu\|_{\text{TV}}]. \end{aligned}$$

Using  $\text{Sign}(\nu)$  as a subgradient of the TV-norm at point  $\nu$  ( $\nu$ -almost everywhere equal to the sign of  $\nu$  and with infinity norm less than one), we then observe that:

$$\|\nu + \sigma\|_{\text{TV}} - \|\nu\|_{\text{TV}} = \langle \text{Sign}(\nu), \sigma \rangle_{\mathcal{M}(\mathcal{X})^*, \mathcal{M}(\mathcal{X})} + \mathcal{D}_\nu(\sigma),$$

where  $\mathcal{D}_\nu(\sigma)$  is the second order Bregman divergence of the TV-norm between  $\nu$  and  $\nu + \sigma$  using the subgradient  $\text{Sign}(\nu)$ , given by:

$$\mathcal{D}_\nu(\sigma) := \|\nu + \sigma\|_{\text{TV}} - \|\nu\|_{\text{TV}} - \langle \text{Sign}(\nu), \sigma \rangle_{\mathcal{M}(\mathcal{X})^*, \mathcal{M}(\mathcal{X})},$$

with  $\mathcal{M}(\mathcal{X})^* \subseteq (L^\infty(\mathcal{X}), \|\cdot\|_{\infty})$  the topological dual of  $\mathcal{M}(\mathcal{X})$ . Gathering all the pieces, we obtain that:

$$J(\nu + \sigma) - J(\nu) = \langle J'_\nu, \sigma \rangle_{\mathcal{M}(\mathcal{X})^*, \mathcal{M}(\mathcal{X})} + q(\sigma), \quad (\text{B.2})$$

where  $J'_\nu$  is given in the statement of Proposition B.1 and  $q$  is a second order term given by:

$$q(\sigma) := \frac{1}{2} \|\Phi(\sigma)\|_{\mathbb{H}}^2 + \lambda \mathcal{D}_\nu(\sigma).$$

Finally, we remark that when  $\nu$  is nonnegative, one possible choice for the TV subgradient is  $\text{Sign}(\nu) = 1$ . In this case, the previous decomposition may be simplified as:

$$J'_\nu = \Phi^*(\Phi(\nu) - \mathbf{y}) + \lambda$$

and  $\mathcal{D}_\nu(\sigma) = 0$  when  $\nu + \sigma \in \mathcal{M}(\mathcal{X})_+$  and  $\nu \in \mathcal{M}(\mathcal{X})_+$ .  $\square$

**Remark B.1.** *The above proof shows that for any  $\mu, \nu \in \mathcal{M}(\mathcal{X})$ ,*

$$J(\nu + \sigma) - J(\nu) = \int_{\mathcal{X}} J'_\nu d\sigma + \frac{1}{2} \|\Phi(\sigma)\|_{\mathbb{H}}^2 + \lambda \mathcal{D}_\nu(\sigma),$$

where  $\mathcal{D}_\nu(\sigma) := \|\nu + \sigma\|_{\text{TV}} - \|\nu\|_{\text{TV}} - \langle \text{Sign}(\nu), \sigma \rangle_{\mathcal{M}(\mathcal{X}), \mathcal{M}(\mathcal{X})}$  and  $\mathcal{M}(\mathcal{X})^* \subseteq (L^\infty(\mathcal{X}), \|\cdot\|_\infty)$  the topological dual of  $\mathcal{M}(\mathcal{X})$ .

Besides the expression of the Fréchet differential of  $J$  on the space  $\mathcal{M}(\mathcal{X})_+$ , it is possible to explicit the value of  $J'_\nu$  at any point  $t \in \mathcal{X}$ . Using Lemma A.2, it holds,

$$J'_\nu(t) = \langle \varphi_t, \Phi(\nu) - \mathbf{y} \rangle_{\mathbb{H}} + \lambda, \quad \forall t \in \mathcal{X}. \quad (\text{B.3})$$

Note that  $J'_\nu$  depends on  $\nu$  through  $\Phi(\nu)$ . Recall that  $\Phi(\mu^*)$  is constant across all solutions  $\mu^*$  of Program  $(\mathcal{B}_+)$ , see Theorem 1.1. Hence, the function

$$J'_*(x) := \langle \varphi_x, \Phi(\mu^*) - \mathbf{y} \rangle_{\mathbb{H}} + \lambda, \quad x \in \mathcal{X},$$

is well defined and does not depend on the choice of the solution  $\mu^*$  (it is the same function across all possible choice of  $\mu^*$  solution to  $(\mathcal{B}_+)$ ). The next proposition gives the first order condition of Program  $(\mathcal{B}_+)$ .

**Proposition B.2.** *It holds that  $J'_* \geq 0$  and, for any solution  $\mu^*$  to Program  $(\mathcal{B}_+)$ ,*

$$\text{Supp}(\mu^*) \subseteq \{x \in \mathcal{X} : J'_*(x) = 0\}.$$

*Conversely, if a measure  $\nu \in \mathcal{M}(\mathcal{X})_+$  is such that  $J'_\nu \geq 0$  then  $\nu$  is a solution to Program  $(\mathcal{B}_+)$  and  $J'_* = J'_\nu$ .*

*Proof.* Let  $x \in \mathcal{X}$  and  $\varepsilon > 0$ . By Proposition B.1, one has

$$J'_*(x) = \frac{J(\mu^* + \varepsilon \delta_x) - J(\mu^*)}{\varepsilon} - \frac{\varepsilon}{2} \|\varphi_x\|_{\mathbb{H}}^2,$$

hence

$$J'_*(x) \geq \liminf_{\varepsilon \downarrow 0} \left\{ \frac{J(\mu^* + \varepsilon \delta_x) - J(\mu^*)}{\varepsilon} \right\} \geq 0,$$

since  $J(\mu^* + \varepsilon \delta_x) - J(\mu^*) \geq 0$ .

Assume now that there exists a point  $x \in \mathcal{X}$  such that  $x \in \text{Supp}(\mu^*)$  and  $J'_*(x) > 0$ . Since  $J'_*$  is continuous and  $J'_*(x) > 0$  there exists  $\varepsilon > 0$  and a open neighborhood  $U_x$  of  $x$  such that

$$\forall t \in U_x, \quad J'_*(t) > \sqrt{\varepsilon}$$

By Jordan decomposition theorem, there exists two nonnegative measures  $\mu^*_+$  and  $\mu^*_-$  with disjoint supports such that  $\mu^* = \mu^*_+ - \mu^*_-$  and  $\text{Supp}(\mu^*) = \text{Supp}(\mu^*_+) \sqcup \text{Supp}(\mu^*_-)$ . Without loss of generality, we assume that  $x \in \text{Supp}(\mu^*_+)$ . Taking  $\varepsilon > 0$  and  $U_x$  sufficiently smalls, one has  $U_x \cap \text{Supp}(\mu^*_-) = \emptyset$  and

$$\mu^*(U_x) > \sqrt{\varepsilon}.$$

Let  $B$  be a Borelian of  $\mathcal{X}$  and define  $\sigma \in \mathcal{M}(\mathcal{X})$  by  $\sigma(B) := -\mu_+^*(B \cap U_x)$ . Remark that  $\mathcal{D}_{\mu^*}(\sigma) = 0$ , this latter being straightforward when  $\mu^* \in \mathcal{M}(\mathcal{X})_+$  (in this case  $\mu^* + \sigma \in \mathcal{M}(\mathcal{X})_+$ ). By Proposition B.1, one has

$$0 \leq J(\mu^* + \sigma) - J(\mu^*) = \int_{\mathcal{X}} J'_* d\sigma = - \int_{\mathcal{X} \cap U_x} J'_* d\mu_+^*,$$

and also

$$\int_{\mathcal{X} \cap U_x} J'_* d\mu_+^* \geq \varepsilon,$$

which is a contradiction. □

The lemma displayed below provides some bounds on the Frechet differential of the objective function and on its stochastic estimate.

**Lemma B.1.** *There exists a positive constant  $\mathcal{C}_0 = \mathcal{C}_0(y, \varphi, \lambda)$  such that*

$$\|J'_\nu\|_\infty := \sup_{t \in \mathcal{X}} |J'_\nu(t)| \leq \mathcal{C}_0(\|\nu\|_{TV} + 1) \quad \forall \nu \in \mathcal{M}(\mathcal{X})_+.$$

Moreover, provided Assumption (A<sub>1</sub>) is satisfied, we have almost surely for any  $\nu \in \mathcal{M}(\mathcal{X})_+$

$$\sup_{t \in \mathcal{X}} |J'_\nu(t, Z)| \leq \mathcal{C}_1(\|\nu\|_{TV} + 1) \quad \text{and} \quad \sup_{t \in \mathcal{X}} |\xi_\nu(t, Z)| \leq \mathcal{C}_2(\|\nu\|_{TV} + 1),$$

for some constants  $\mathcal{C}_1$  and  $\mathcal{C}_2$  depending only on  $y, \varphi$  and  $\lambda$ .

*Proof.* Let  $\nu \in \mathcal{M}(\mathcal{X})_+$  be fixed. We denote by  $\tilde{\nu}$  the normalized measure  $\tilde{\nu} = \nu/\|\nu\|_{TV} \in \mathcal{M}(\mathcal{X})_+$ . According to (B.3), we have

$$\begin{aligned} \sup_{t \in \mathcal{X}} |J'_\nu(t)| &\leq \lambda + \sup_{t \in \mathcal{X}} |\langle \varphi_t, y \rangle_{\mathbb{H}}| + \|\nu\|_{TV} \sup_{t \in \mathcal{X}} |\langle \varphi_t, \Phi(\tilde{\nu}) \rangle_{\mathbb{H}}|, \\ &\leq \lambda + \sup_{t \in \mathcal{X}} |\langle \varphi_t, y \rangle_{\mathbb{H}}| + \|\nu\|_{TV} \sup_{t, s \in \mathcal{X}} |\langle \varphi_t, \varphi_s \rangle_{\mathbb{H}}|, \\ &\leq \lambda + \|\varphi\|_{\infty, \mathbb{H}} \|y\|_{\mathbb{H}} + \|\nu\|_{TV} \|\varphi\|_{\infty, \mathbb{H}}^2, \\ &\leq \mathcal{C}_0(\|\nu\|_{TV} + 1), \end{aligned}$$

with

$$\mathcal{C}_0 = \max(\lambda + \|\varphi\|_{\infty, \mathbb{H}} \|y\|_{\mathbb{H}}; \|\varphi\|_{\infty, \mathbb{H}}^2). \quad (\text{B.4})$$

Concerning the second part of the lemma, we first remark that, provided Assumption (A<sub>1</sub>) is satisfied, we have for any  $\nu \in \mathcal{M}(\mathcal{X})_+$

$$\sup_{t \in \mathcal{X}} |J'_\nu(t, Z)| \leq \|\nu\|_{TV} \sup_{t \in \mathcal{X}} |g_{t, T}(U)| + \sup_{t \in \mathcal{X}} |h_t(V)| + \lambda \leq \mathcal{C}_1(\|\nu\|_{TV} + 1),$$

with

$$\mathcal{C}_1 := \max(\|g\|_{\infty}; \|h\|_{\infty} + \lambda). \quad (\text{B.5})$$

The last results is obtained thanks to a basic triangle inequality

$$\sup_{t \in \mathcal{X}} |\xi_\nu(t, Z)| \leq \sup_{t \in \mathcal{X}} |J'_\nu(t)| + \sup_{t \in \mathcal{X}} |J'_\nu(t, Z)| \leq \mathcal{C}_2(\|\nu\|_{TV} + 1),$$

with

$$\mathcal{C}_2 = \mathcal{C}_0 + \mathcal{C}_1 = \max(\lambda + \|\varphi\|_{\infty, \mathbb{H}} \|y\|_{\mathbb{H}}; \|\varphi\|_{\infty, \mathbb{H}}^2) + \max(\|g\|_{\infty}; \|h\|_{\infty} + \lambda). \quad (\text{B.6})$$

□

## B.2 In the space of particles

We consider any set of positions  $\mathbb{T}$  and their associate weights  $\mathbb{W}$ . In order to compute the derivatives of  $F$  w.r.t.  $\mathbb{W}$  and  $\mathbb{T}$ , our starting point is Equation (1.9) and we observe that the gradient with respect to  $\mathbb{W}$  is easily computed:

$$\nabla_{\omega} F(\omega, \mathbb{T}) = \lambda - \mathbf{k}_{\mathbb{T}} + \mathbb{K}_{\mathbb{T}} \mathbb{W}.$$

Nevertheless, the interpretation in terms of Fréchet derivative and of  $J$  allows to obtain the next result.

**Proposition B.3.** *For any  $\mathbb{W}$  and  $\mathbb{T}$ , denote  $\nu = \nu(\mathbb{W}, \mathbb{T})$ , one has:*

(i) *Gradient w.r.t. weights: for any  $j \in \{1, \dots, p\}$ , one has  $\nabla_{\omega_j} F(\omega, \mathbb{T}) = J'_{\nu}(\mathbf{t}_j)$*

(ii) *Gradient w.r.t. positions: for any  $j \in \{1, \dots, p\}$ , one has  $\nabla_{\mathbf{t}_j} F(\omega, \mathbb{T}) = \omega_j \nabla_{\mathbf{t}_j} J'_{\nu}(\mathbf{t}_j)$*

*Proof.* The starting point is the Fréchet derivative that, if we consider  $\sigma \in \mathcal{M}(\mathcal{X})_+$  and  $\varepsilon > 0$  small enough:

$$J(\nu + \varepsilon\sigma) = J(\nu) + \varepsilon \langle J'_{\nu}, \sigma \rangle_{\mathbb{H}} + o(\varepsilon).$$

Proof of (i): Considering any particle  $j \in \{1, \dots, p\}$  and  $\sigma = \delta_{\mathbf{t}_j}$ , we then obtain that:

$$\lim_{\varepsilon \rightarrow 0} \frac{J(\nu + \varepsilon\sigma) - J(\nu)}{\varepsilon} = \langle J'_{\nu}, \delta_{\mathbf{t}_j} \rangle_{\mathbb{H}} = J'_{\nu}(\mathbf{t}_j),$$

where the last equality comes from the reproducing kernel property. In the meantime, we observe that

$$\lim_{\varepsilon \rightarrow 0} \frac{J(\nu + \varepsilon\sigma) - J(\nu)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{F(\mathbb{W} + \varepsilon \mathbf{1}_j, \mathbb{T}) - F(\mathbb{W}, \mathbb{T})}{\varepsilon} = \frac{\partial F(\mathbb{W}, \mathbb{T})}{\partial \omega_j}.$$

We then conclude using the Fréchet derivative of  $J$  that:

$$J'_{\nu}(\mathbf{t}_j) = \nabla_{\omega_j} F(\omega, \mathbb{T}).$$

Proof of (ii): Using the same consideration on the positions of the particles, we then consider any perturbed set of positions  $\tilde{\mathbb{T}}_{\varepsilon,j} = \mathbb{T} + \varepsilon \mathbf{1}_j$  where only the coordinate  $j$  of  $\mathbb{T}$  is modified. We then write the partial derivative of  $F$ :

$$\lim_{\varepsilon \rightarrow 0} \frac{F(\mathbb{W}, \tilde{\mathbb{T}}_{\varepsilon,j}) - F(\mathbb{W}, \mathbb{T})}{\varepsilon} = \frac{\partial F(\mathbb{W}, \mathbb{T})}{\partial \mathbf{t}_j}.$$

In the meantime, we observe that with the Fréchet derivative of  $J$  that:

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \frac{F(\mathbb{W}, \tilde{\mathbb{T}}_{\varepsilon,j}) - F(\mathbb{W}, \mathbb{T})}{\varepsilon} &= \lim_{\varepsilon \rightarrow 0} \frac{J(\nu(\mathbb{W}, \tilde{\mathbb{T}}_{\varepsilon,j})) - J(\nu(\mathbb{W}, \mathbb{T}))}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{J(\nu(\mathbb{W}, \mathbb{T}) + \omega_j [\delta_{\mathbf{t}_j + \varepsilon} - \delta_{\mathbf{t}_j}]) - J(\nu(\mathbb{W}, \mathbb{T}))}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{\langle J'_{\nu(\mathbb{W}, \mathbb{T})}, \omega_j \delta_{\mathbf{t}_j + \varepsilon} \rangle_{\mathbb{H}} - \langle J'_{\nu(\mathbb{W}, \mathbb{T})}, \omega_j \delta_{\mathbf{t}_j} \rangle_{\mathbb{H}}}{\varepsilon} \\ &= \omega_j \nabla_{\mathbf{t}_j} J'_{\nu(\mathbb{W}, \mathbb{T})}(\mathbf{t}_j). \end{aligned}$$

□

Finally, it is possible to quantify the way  $F$  is modified when we change  $(\mathbb{W}, \mathbb{T})$  to  $(\mathbb{W}', \mathbb{T}')$  thanks to the next proposition. where  $\nu = \sum_{k=1}^p \omega_k \delta_{\mathbf{t}_k}$ .

**Proposition B.4.** *Consider two pairs  $(\mathbb{W}, \mathbb{T})$  and  $(\mathbb{W}', \mathbb{T}')$  of weights/positions and denote  $\nu = \nu(\mathbb{W}, \mathbb{T})$  defined in Equation (1.7), then:*

$$F(\omega', \mathbb{T}') - F(\omega, \mathbb{T}) = \sum_{j=1}^p (\omega'_j J'_{\nu}(\mathbf{t}'_j) - \omega_j J'_{\nu}(\mathbf{t}_j)) + \frac{1}{2} (\omega', -\omega)^{\top} \mathbb{K}_{(\mathbb{T}', \mathbb{T})} (\omega', -\omega), \quad (\text{B.7})$$

where  $\mathbb{K}_{(\mathbb{T}', \mathbb{T})}$  is a  $(2p \times 2p)$  symmetric matrix with  $(p \times p)$  diagonal blocks  $\langle \varphi_{\mathbf{t}'_j}, \varphi_{\mathbf{t}'_j} \rangle_{\mathbb{H}}$  and  $\langle \varphi_{\mathbf{t}_j}, \varphi_{\mathbf{t}_j} \rangle_{\mathbb{H}}$ , and  $(p \times p)$  off-diagonal block  $\langle \varphi_{\mathbf{t}'_j}, \varphi_{\mathbf{t}_j} \rangle_{\mathbb{H}}$ .

*Proof.* We denote  $\nu' = \nu(\mathbb{W}', \mathbb{T}')$  and apply Equation (B.1) with  $\nu' = \nu + \sigma$  with  $\sigma = \nu' - \nu$ , we obtain that:

$$F(\mathbb{W}', \mathbb{T}') - F(\mathbb{W}, \mathbb{T}) = J(\nu + \sigma) - J(\nu) = \int_{\mathcal{X}} J'_\nu d\sigma + \frac{1}{2} \|\Phi(\sigma)\|_{\mathbb{H}}^2 = \int_{\mathcal{X}} J'_\nu d\sigma + \frac{1}{2} \|\Phi(\nu') - \Phi(\nu)\|_{\mathbb{H}}^2$$

Note that the last term rewrites

$$(\boldsymbol{\omega}', -\boldsymbol{\omega}')^\top \mathbb{K}_{(\mathbb{T}', \mathbb{T}')}(\boldsymbol{\omega}', -\boldsymbol{\omega}') = \|\Phi(\nu') - \Phi(\nu)\|_{\mathbb{H}}^2,$$

which is the squared Maximum Mean Discrepancy (MMD) between  $\nu$  and  $\nu'$  for the kernel  $K$ .  $\square$

## C Technical results for Theorem 4.1

**Proposition C.1.** Consider  $(\mathbf{t}, \tilde{\mathbf{t}}) \in \mathcal{X}^2$ , the following technical inequalities hold:

- i)  $|\nabla J'_\nu(\mathbf{t}) - \nabla J'_\nu(\tilde{\mathbf{t}})| \leq [\text{Lip}(\nabla\varphi)\|\nu\|_{\mathcal{T}\mathcal{V}} + \text{Lip}(\nabla y)] \|\mathbf{t} - \tilde{\mathbf{t}}\|$ ,
- ii)  $\|\nabla_{\mathbf{t}} J'_\nu(\mathbf{t})\|_{\mathcal{X}} \leq (\|\nu\|_{\mathcal{T}\mathcal{V}}\|\varphi\|_{\infty, \mathbb{H}} + \|y\|_{\mathbb{H}})\|\varphi'\|_{\mathbb{H}}$  with  $\|\varphi'\|_{\infty, \mathbb{H}} := \sup_{\mathbf{t}} \sup_{\psi: \|\psi\|_{\mathbb{H}} \leq 1} \|\nabla_{\mathbf{t}} \langle \varphi_{\mathbf{t}}, \psi \rangle\|_{\mathcal{X}}$ ,
- iii)  $\|D_\nu(\mathbf{t}, Z)\|_{\mathcal{X}} \leq \|\nu\|_{\mathcal{T}\mathcal{V}}\|\mathbf{g}'\|_{\infty, \mathbb{H}} + \|\mathbf{h}'\|_{\mathbb{H}}$ ,

where Assumption (A<sub>2</sub>) is required for inequality iii).

*Proof of Proposition C.1.* We consider any finite measure  $\nu$ .

Proof of i) We provide an upper bound on the Lipschitz constant of  $\nabla J'_\nu$ : consider  $(\mathbf{t}, \tilde{\mathbf{t}}) \in \mathcal{X}^2$ , we repeat the same arguments as above and observe that:

$$|\nabla J'_\nu(\mathbf{t}) - \nabla J'_\nu(\tilde{\mathbf{t}})| \leq [\text{Lip}(\nabla\varphi)\|\nu\|_{\mathcal{T}\mathcal{V}} + \text{Lip}(\nabla y)] \|\mathbf{t} - \tilde{\mathbf{t}}\|.$$

Proof of ii) Using a rough bound, we get

$$\begin{aligned} \|\nabla_{\mathbf{t}} J'_\nu(\mathbf{t})\|_{\mathcal{X}} &= \left\| \sum_{j=1}^p \omega_j \nabla_{\mathbf{t}} \langle \varphi_{\mathbf{t}}, \varphi_{t_j} \rangle_{\mathbb{H}} - \nabla_{\mathbf{t}} \langle \varphi_{\mathbf{t}}, y \rangle_{\mathbb{H}} \right\|_{\mathcal{X}}, \\ &\leq \sum_{j=1}^p \omega_j \|\nabla_{\mathbf{t}} \langle \varphi_{\mathbf{t}}, \varphi_{t_j} \rangle_{\mathbb{H}}\|_{\mathcal{X}} + \|\nabla_{\mathbf{t}} \langle \varphi_{\mathbf{t}}, y \rangle_{\mathbb{H}}\|_{\mathcal{X}}, \\ &\leq (\|\nu\|_{\mathcal{T}\mathcal{V}}\|\varphi\|_{\infty, \mathbb{H}} + \|y\|_{\mathbb{H}}) \times \sup_{\psi: \|\psi\|_{\mathbb{H}} \leq 1} \|\nabla_{\mathbf{t}} \langle \varphi_{\mathbf{t}}, \psi \rangle\|_{\mathcal{X}}, \end{aligned}$$

which provides the desired result.

Proof of iii) Using Assumption (A<sub>2</sub>), and in particular the boundedness of the derivative of  $\mathbf{g}$  and  $\mathbf{h}$ , we get

$$\begin{aligned} \|D_\nu(\mathbf{t}, Z)\|_{\mathcal{X}} &= \|\|\nu\|_{\mathcal{T}\mathcal{V}} \nabla_{\mathbf{t}} \mathbf{g}_{\mathbf{t}, \mathcal{T}}(U) - \nabla_{\mathbf{t}} \mathbf{h}_{\mathbf{t}}(V)\|_{\mathcal{X}}, \\ &\leq \|\nu\|_{\mathcal{T}\mathcal{V}} \sup_{t, s, u} \|\nabla_{\mathbf{t}} \mathbf{g}_{t, s}(u)\|_{\mathcal{X}} + \sup_{t, v} \|\nabla_{\mathbf{t}} \mathbf{h}_{\mathbf{t}}(v)\|_{\mathcal{X}}. \end{aligned}$$

$\square$

**Proposition C.2.** A large enough constant  $\mathfrak{C}$  exists such that for any iteration  $k \in \mathbb{N}$  and any particle  $j \in \{1, \dots, p\}$ , :

$$\left| \mathbb{E} \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \mid \mathfrak{F}_k \right] \right| \leq \mathfrak{C} \alpha^2 (1 + \|\nu_k\|_{\mathcal{T}\mathcal{V}})^2.$$

*Proof.* This key technical argument relies on the Hoeffding inequality. We shall write:

$$\begin{aligned} \mathbb{E} \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \mid \mathfrak{F}_k \right] &= \alpha J'_{\nu_k}(\mathbf{t}_j^k) - 1 + \mathbb{E} \left[ e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} \mid \mathfrak{F}_k \right] \\ &= \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) - 1 + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k)} \mathbb{E} \left[ e^{-\alpha [J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) - J'_{\nu_k}(\mathbf{t}_j^k)]} \mid \mathfrak{F}_k \right] \right] \\ &= \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) - 1 + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k)} \right] \\ &\quad + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k)} \mathbb{E} \left[ e^{-\alpha [J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) - J'_{\nu_k}(\mathbf{t}_j^k)]} - 1 \mid \mathfrak{F}_k \right] \end{aligned}$$

To derive an upper bound, we apply the Hoeffding Lemma to the random variable  $J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) - J'_{\nu_k}(\mathbf{t}_j^k)$  that is centered and bounded by  $\mathbf{T} = \mathfrak{C}(1 + \|\nu_k\|_{\text{TV}})$  from according to Lemma B.1. We obtain that:

$$\left| \mathbb{E} \left[ e^{-\alpha[J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1}) - J'_{\nu_k}(\mathbf{t}_j^k)]} - 1 \middle| \mathfrak{F}_k \right] \right| \leq e^{\frac{\tau^2 \alpha^2}{8}} - 1 \leq \mathfrak{C} \alpha^2 (1 + \|\nu_k\|_{\text{TV}}^2).$$

Using that  $|e^x - 1 - x| \leq c|x|^2$  for bounded  $x = \alpha J'_{\nu_k}(\mathbf{t}_j^k)$  and  $c$  large enough, we finally obtain that:

$$\left| \mathbb{E} \left[ \alpha J'_{\nu_k}(\mathbf{t}_j^k) + e^{-\alpha J'_{\nu_k}(\mathbf{t}_j^k, Z^{k+1})} - 1 \middle| \mathfrak{F}_k \right] \right| \leq \mathfrak{C} \alpha^2 (1 + \|\nu_k\|_{\text{TV}}^2).$$

□

We recall here the result essentially due to Chizat [2022], which is stated in a simplest way for our purpose.

**Proposition C.3.** *Assume that  $\mu^*$  is discrete and that  $\nu_0$  is a uniform distribution over a grid of size  $\delta = \frac{2d}{\tau}$  where  $d$  is the dimension of  $\mathcal{X}$ , then:*

$$Q_\tau(\mu^*, \nu_0) \leq \frac{\|\mu^*\|_{\text{TV}} d}{\tau} \left( 1 + \log \frac{\tau}{2d} + \frac{\log |\mathcal{X}|}{d} \right)$$

Moreover, the measure  $\nu_0^\delta$  that meets this upper bound satisfies  $\|\nu_0^\delta\|_{\text{TV}} = \|\mu^*\|_{\text{TV}}$ .

*Proof.* We define  $m$  as the size of the support of  $\nu_0$  and

$$\nu_0 = m^{-1} \sum_{i=1}^m \delta_{x_i},$$

where  $(x_i)_{1 \leq i \leq m}$  refers to the uniform grid of size  $\delta$  on  $\mathcal{X}$ . Since  $\mu^*$  is discrete, it may be written as:

$$\mu^* = \sum_{j=1}^{m^*} \mu_j^* \delta_{z_j^*}.$$

For any support point  $z_j^*$  of  $\mu^*$ , we then consider  $i_j \in \{1, \dots, m\}$  such that  $\|x_{i_j} - z_j^*\| \leq \delta/2$  and we define  $\nu_0^\delta$  as:

$$\nu_0^\delta := \sum_{j=1}^{m^*} \mu_j^* \delta_{x_{i_j}}$$

We observe that by construction,  $\|\nu_0^\delta\|_{\text{TV}} = \|\mu^*\|_{\text{TV}}$  and

$$\begin{aligned} \mathcal{H}(\nu_0, \nu_0^\delta) &= \sum_{j=1}^{m^*} \nu_0^\delta(x_{i_j}) \log \left( \frac{\nu_0^\delta(x_{i_j})}{\nu_0(x_{i_j})} \right) \\ &= \sum_{j=1}^{m^*} \mu_j^* \log \left( \frac{\mu_j^*}{\nu_0(x_{i_j})} \right) \\ &\leq -\text{Ent}(\mu^*) + \|\mu^*\|_{\text{TV}} \left[ d \log \left( \frac{1}{\delta} \right) + \log |\mathcal{X}| \right], \end{aligned} \tag{C.1}$$

where we used the entropy of a discrete measure defined as

$$\text{Ent}(\mu) = - \sum_{x \in \text{Supp}(\mu)} \mu(x) \log(\mu(x))$$

and a lower bound of  $\nu_0(x_{i_j})$ , which is of the order  $\delta^d |\mathcal{X}|^{-1}$  where  $|\mathcal{X}|$  refers to the Lebesgue measure of  $\mathcal{X}$ .



In the meantime, we also observe that the BL dual norm between  $\nu_0^\delta$  and  $\mu^*$  can be easily upper bounded. Indeed

$$\begin{aligned}
\|\nu_0^\delta - \mu^*\|_{BL}^* &= \sup_{\|f\|_{BL} \leq 1} \int_{\mathcal{X}} f d[\nu_0^\delta - \mu^*] \\
&= \sup_{\|f\|_{BL} \leq 1} \sum_{j=1}^{m^*} \mu_j^* [f(x_{i_j}) - f(z_j)] \\
&\leq \frac{\|\mu^*\|_{TV} \delta}{2}
\end{aligned} \tag{C.2}$$

We then add the two upper bounds (C.1) and (C.2) and minimize

$$\delta \mapsto \frac{\|\mu^*\|_{TV} \delta}{2} + \frac{1}{\tau} \left( -\text{Ent}(\mu^*) + \|\mu^*\|_{TV} \left[ d \log \left( \frac{1}{\delta} \right) + \log |\mathcal{X}| \right] \right).$$

We are led to choose  $\delta = 2d\tau^{-1}$  and we obtain the following upper bound:

$$\frac{1}{\tau} \mathcal{H}(\nu_0, \nu_0^\delta) + \|\nu_0^\delta - \mu^*\|_{BL}^* \leq \frac{\|\mu^*\|_{TV} d}{\tau} \left( 1 + \log \frac{\tau}{2d} + \frac{\log |\mathcal{X}|}{d} \right),$$

which ends the proof of the proposition. □