

# A stochastic Gauss-Newton algorithm for regularized semi-discrete optimal transport

Bernard Bercu<sup>1</sup>, Jérémie Bigot<sup>1</sup>, Sébastien Gadat<sup>2</sup>, Emilia Siviero<sup>3</sup>

<sup>1</sup>Institut de Mathématiques de Bordeaux et CNRS (UMR 5251), Université de Bordeaux

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

<sup>3</sup>LTCI, Télécom Paris, Institut Polytechnique de Paris

March 3, 2022

## Abstract

We introduce a new second order stochastic algorithm to estimate the entropically regularized optimal transport cost between two probability measures. The source measure can be arbitrary chosen, either absolutely continuous or discrete, while the target measure is assumed to be discrete. To solve the semi-dual formulation of such a regularized and semi-discrete optimal transportation problem, we propose to consider a stochastic Gauss-Newton algorithm that uses a sequence of data sampled from the source measure. This algorithm is shown to be adaptive to the geometry of the underlying convex optimization problem with no important hyperparameter to be accurately tuned. We establish the almost sure convergence and the asymptotic normality of various estimators of interest that are constructed from this stochastic Gauss-Newton algorithm. We also analyze their non-asymptotic rates of convergence for the expected quadratic risk in the absence of strong convexity of the underlying objective function. The results of numerical experiments from simulated data are also reported to illustrate the finite sample properties of this Gauss-Newton algorithm for stochastic regularized optimal transport, and to show its advantages over the use of the stochastic gradient descent, stochastic Newton and ADAM algorithms.

**Keywords:** Stochastic optimization; Stochastic Gauss-Newton algorithm; Optimal transport; Entropic regularization; Convergence of random variables.

**AMS classifications:** Primary 62G05; secondary 62G20.

## 1 Introduction

### 1.1 Computational optimal transport for data science

The use of optimal transport (OT) and Wasserstein distances for data science has recently gained an increasing interest in various research fields such as machine learning [1, 20, 21, 23, 25, 42, 44], statistics [5, 7, 8, 13, 30, 35, 40, 46, 49] and image processing or computer vision [4, 12, 19, 26, 39, 45]. Solving a problem of OT between two probability measures  $\mu$  and  $\nu$  is known to be computationally challenging, and entropic regularization [15, 16] has emerged as an efficient tool to approximate and smooth the variational Wasserstein problems arising in computational optimal transport for data science. A detailed presentation of the recent

research field of computational optimal transport is given in [38], while recent reviews on the application of optimal transport to statistics can be found in [9, 36].

Recently, approaches [5, 23] based on first order stochastic algorithms have gained popularity to solve (possibly regularized) OT problems using data sampled from  $\mu$ . These approaches are based on the semi-dual formulation [16] of regularized OT problems that can be rewritten as a *non-strongly convex* stochastic optimization problem. In this paper, for the purpose of obtaining stochastic algorithms for OT with faster convergence in practice, we introduce a second order stochastic algorithm to solve regularized semi-discrete OT between an arbitrary probability measure  $\mu$ , typically absolutely continuous, and a *known* discrete measure  $\nu$  with finite support of size  $J$ . More precisely, we focus on the estimation of an entropically regularized optimal transport cost  $W_\varepsilon(\mu, \nu)$  between such measures (where  $\varepsilon > 0$  is an entropic regularization parameter) using a class of stochastic quasi-Newton algorithms that we refer to as Gauss-Newton algorithms and which use the knowledge of a sequence  $(X_n)$  of independent random vectors sampled from  $\mu$ .

Applications of semi-discrete optimal transport can be found in computational geometry and computer graphics [33, 34], as well as in the problem of optimal allocation of resources from online observations [5]. For an introduction to semi-discrete optimal transport problems and related references, we also refer to [38, Chapter 5]. In a deterministic setting where the full knowledge of  $\mu$  is used and in the unregularized case, the convergence of a Newton algorithm for semi-discrete optimal transport has been studied in depth in [29]. An extension of the formulation of semi-discrete OT to include an entropic regularization is proposed in [16]. The main advantage of incorporating such a regularization term in classical OT is to obtain a dual formulation leading to a smooth convex minimization problem allowing the implementation of simple and numerically more stable algorithms as shown in [16]. The use of regularized semi-discrete OT has then found applications in image processing using generative adversarial networks [24, 43]. In these works, samples from the generative model are typically drawn from an absolutely continuous source measure in order to fit a discrete target distribution.

## 1.2 Main contributions and structure of the paper

As discussed above, we introduce a stochastic Gauss-Newton (SGN) algorithm for regularized semi-discrete OT for the purpose of estimating  $W_\varepsilon(\mu, \nu)$ , and the main goal of this paper is to study the statistical properties of such an approach. This algorithm is shown to be adaptive to the geometry of the underlying convex optimization problem with no important hyperparameter to be accurately tuned. Then, the main contributions of our work are to derive the almost sure rates of convergence, the asymptotic normality and the non-asymptotic rates of convergence (in expectation) of various estimators of interest that are constructed using the SGN algorithm to be described below. Although the underlying stochastic optimization problem is not strongly convex, fast rates of convergence can be obtained by combining the so-called notion of *generalized self-concordance* introduced in [2] that has been shown to hold for regularized OT in [5], and the Kurdyka-Łojasiewicz inequality as studied in [22]. We also report the results from various numerical experiments on simulated data to illustrate the finite sample properties of this algorithm, and to compare its performances with those of the stochastic gradient descent (SGD), stochastic Newton (SN) and ADAM [28] algorithms for stochastic regularized OT.

The paper is then organized as follows. The definitions of regularized semi-discrete OT

and the stochastic algorithms used for solving this problem are given in Section 2. The main results of the paper are stated in Section 3, while the important properties related to the regularized OT are given in Section 5. In Section 4, we describe a fast implementation of the SGN algorithm for regularized semi-discrete OT, and we assess the numerical ability of SGN to solve OT problems. In particular, we report numerical experiments on simulated data to compare the performances of various stochastic algorithms for regularized semi-discrete OT. The statistical properties of the SGN algorithm are established in an extended Section 6 that gathers the proof of our main results. Finally, two technical appendices A and B contain the proofs of auxiliary results.

## 2 A stochastic Gauss-Newton algorithm for regularized semi-discrete OT

In this section, we introduce the notion of regularized semi-discrete OT and the stochastic algorithm that we propose to solve this problem.

### 2.1 Notation, definitions and main assumptions on the OT problem

Let  $\mathcal{X}$  and  $\mathcal{Y}$  be two metric spaces. Denote by  $\mathcal{M}_+^1(\mathcal{X})$  and  $\mathcal{M}_+^1(\mathcal{Y})$  the sets of probability measures on  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. Let  $\mathbf{1}_J$  be the column vector of  $\mathbb{R}^J$  with all coordinates equal to one, and denote by  $\langle \cdot, \cdot \rangle$  and  $\| \cdot \|$  the standard inner product and norm in  $\mathbb{R}^J$ . We also use  $\lambda_{\max}(A)$  and  $\lambda_{\min}(A)$  to denote the largest and smallest eigenvalues of a symmetric matrix  $A$ , whose spectrum will be denoted by  $\text{Sp}(A)$  and Moore-Penrose inverse by  $A^-$ . By a slight abuse of notation, we sometimes denote by  $\lambda_{\min}(A)$  the smallest non-zero eigenvalue of a positive semi-definite matrix  $A$ . Finally,  $\|A\|_2$  and  $\|A\|_F$  stand for the operator and Frobenius norms of  $A$ , respectively. For  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$ , let  $\Pi(\mu, \nu)$  be the set of probability measures on  $\mathcal{X} \times \mathcal{Y}$  with marginals  $\mu$  and  $\nu$ . As formulated in [23], the problem of entropically regularized optimal transport between  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$  is defined as follows.

**Definition 2.1.** *For any  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y})$ , the Kantorovich formulation of the regularized optimal transport between  $\mu$  and  $\nu$  is the following convex minimization problem*

$$W_\varepsilon(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) + \varepsilon KL(\pi | \mu \otimes \nu), \quad (2.1)$$

where  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is a lower semi-continuous function referred to as the cost function of moving mass from location  $x$  to  $y$ ,  $\varepsilon \geq 0$  is a regularization parameter, and  $KL$  stands for the Kullback-Leibler divergence between  $\pi$  and a positive measure  $\xi$  on  $\mathcal{X} \times \mathcal{Y}$ , up to the additive term  $\int_{\mathcal{X} \times \mathcal{Y}} d\xi(x, y)$ , namely

$$KL(\pi | \xi) = \int_{\mathcal{X} \times \mathcal{Y}} \left( \log \left( \frac{d\pi}{d\xi}(x, y) \right) - 1 \right) d\pi(x, y).$$

For  $\varepsilon = 0$ , the quantity  $W_0(\mu, \nu)$  is the *standard OT cost*, while for  $\varepsilon > 0$ , we refer to  $W_\varepsilon(\mu, \nu)$  as the *regularized OT cost* between the two probability measures  $\mu$  and  $\nu$ . In

this framework, we shall consider cost functions that are lower semi-continuous and that satisfy the following standard assumption (see e.g. [48, Part I-4]), for all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ ,

$$0 \leq c(x, y) \leq c_{\mathcal{X}}(x) + c_{\mathcal{Y}}(y), \quad (2.2)$$

where  $c_{\mathcal{X}}$  and  $c_{\mathcal{Y}}$  are real-valued functions such that  $\int_{\mathcal{X}} c_{\mathcal{X}}(x) d\mu(x) < +\infty$  and  $\int_{\mathcal{Y}} c_{\mathcal{Y}}(y) d\nu(y) < +\infty$ . Under condition (2.2),  $W_{\varepsilon}(\mu, \nu)$  is finite regardless any value of the regularization parameter  $\varepsilon \geq 0$ . Moreover, note that  $W_{\varepsilon}(\mu, \nu)$  can be negative for  $\varepsilon > 0$ , and that we always have the lower bound  $W_{\varepsilon}(\mu, \nu) \geq -\varepsilon$ . In this paper, we concentrate on the *regularized case* where  $\varepsilon > 0$ , and on the *semi-discrete setting* where  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  is an arbitrary probability measure (e.g. either discrete or absolutely continuous with respect to the Lebesgue measure), and  $\nu$  is a discrete measure with finite support  $\mathcal{Y} = \{y_1, \dots, y_J\}$  that can be written as

$$\nu = \sum_{j=1}^J \nu_j \delta_{y_j}.$$

Here,  $\delta$  stands for the standard Dirac measure, the locations  $\{y_1, \dots, y_J\}$  as well as the positive weights  $\{\nu_1, \dots, \nu_J\}$  are assumed to be known and summing up to one. We shall also use the notation

$$\min(\nu) = \min_{1 \leq j \leq J} \{\nu_j\} \quad \text{and} \quad \max(\nu) = \max_{1 \leq j \leq J} \{\nu_j\}.$$

We shall also sometimes refer to the *discrete setting* when  $\mu$  is also a discrete measure. Let us now define the semi-dual formulation of the minimization problem (2.1) as introduced in [23]. In the semi-discrete setting and for  $\varepsilon > 0$ , using the semi-dual formulation of the minimization problem (2.1), it follows that  $W_{\varepsilon}(\mu, \nu)$  can be expressed as the following convex optimization problem

$$W_{\varepsilon}(\mu, \nu) = - \inf_{v \in \mathbb{R}^J} H_{\varepsilon}(v) \quad (2.3)$$

with

$$H_{\varepsilon}(v) = \mathbb{E}[h_{\varepsilon}(X, v)] = \int_{\mathcal{X}} h_{\varepsilon}(x, v) d\mu(x), \quad (2.4)$$

where  $X$  stands for a random variable drawn from the unknown distribution  $\mu$ , and for any  $(x, v) \in \mathcal{X} \times \mathbb{R}^J$ ,

$$h_{\varepsilon}(x, v) = \varepsilon + \varepsilon \log \left( \sum_{j=1}^J \exp \left( \frac{v_j - c(x, y_j)}{\varepsilon} \right) \nu_j \right) - \sum_{j=1}^J v_j \nu_j. \quad (2.5)$$

Throughout the paper, we shall assume that, for any  $\varepsilon > 0$ , there exists  $v^* \in \mathbb{R}^J$  that minimizes the function  $H_{\varepsilon}$ , leading to

$$W_{\varepsilon}(\mu, \nu) = -H_{\varepsilon}(v^*).$$

The above equality is the key result allowing to formulate (2.3) as a convex stochastic minimization problem, and to consider the issue of estimating  $W_{\varepsilon}(\mu, \nu)$  in the setting of

stochastic optimization. For a discussion on sufficient conditions implying the existence of such a minimizer  $v^*$ , we refer to [5, Section 2]. As discussed in Section 5, the function  $H_\varepsilon$  possesses a one-dimensional subspace of global minimizers, defined by  $\{v^* + t\mathbf{v}_J, t \in \mathbb{R}\}$  where

$$\mathbf{v}_J = \frac{1}{\sqrt{J}}\mathbf{1}_J.$$

Therefore, we will constrain our algorithm to live in  $\langle \mathbf{v}_J \rangle^\perp$ , which denotes the orthogonal complement of the one-dimensional subspace  $\langle \mathbf{v}_J \rangle$  of  $\mathbb{R}^J$  spanned by  $\mathbf{v}_J$ . In that setting, the OT problem (2.3) becomes identifiable.

## 2.2 Pre-conditioned stochastic algorithms

In the context of regularized OT, we first introduce a general class of stochastic pre-conditioned algorithms that are also referred to as quasi-Newton algorithms in the literature. Starting from Section 2.1, our approach is inspired by the recent works [5, 23], which use the property that

$$W_\varepsilon(\mu, \nu) = -H_\varepsilon(v^*) = -\min_{v \in \mathbb{R}^J} \mathbb{E}[h_\varepsilon(X, v)]$$

where  $h_\varepsilon(x, v)$  is the smooth function defined by (2.5) that is simple to compute. For a sequence  $(X_n)$  of independent and identically distributed random variables sampled from the distribution  $\mu$ , the class of pre-conditioned stochastic algorithms is defined as the following family of recursive stochastic algorithms to estimate the minimizer  $v^* \in \langle \mathbf{v}_J \rangle^\perp$  of  $H_\varepsilon$ . These algorithms can be written as

$$\widehat{V}_{n+1} = P_J \left( \widehat{V}_n - n^\alpha S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n) \right) \quad (2.6)$$

for some constant  $0 \leq \alpha < 1/2$ , where  $\nabla_v h_\varepsilon$  stands for the gradient of  $h_\varepsilon$  with respect to  $v$ ,  $\widehat{V}_0$  is a random vector belonging to  $\langle \mathbf{v}_J \rangle^\perp$ , and  $S_n$  is a symmetric and *positive definite*  $J \times J$  random matrix which is measurable with respect to the  $\sigma$ -algebra  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$ . In addition,  $P_J$  is the orthogonal projection matrix onto  $\langle \mathbf{v}_J \rangle^\perp$ ,

$$P_J = I_J - \mathbf{v}_J \mathbf{v}_J^T.$$

This stochastic algorithm allows us to estimate  $W_\varepsilon(\mu, \nu)$  by the recursive estimator

$$\widehat{W}_n = -\frac{1}{n} \sum_{k=1}^n h_\varepsilon(X_k, \widehat{V}_{k-1}). \quad (2.7)$$

The special case where  $S_n = s^{-1}nI_J$  with some constant  $s > 0$ , corresponds to the well-known stochastic gradient descent (SGD) algorithm that has been introduced in the context of stochastic OT in [23], and recently investigated in [5]. Following some recent contributions [6, 14] in stochastic optimization using Newton-type stochastic algorithms, another potential choice is  $S_n = \mathbb{S}_n$  where  $\mathbb{S}_n$  is the natural Newton recursion defined as

$$\mathbb{S}_n = I_J + \sum_{k=1}^n \nabla_v^2 h_\varepsilon(X_k, \widehat{V}_{k-1}) = \mathbb{S}_{n-1} + \nabla_v^2 h_\varepsilon(X_n, \widehat{V}_{n-1}) \quad (2.8)$$

where  $\nabla_v^2 h_\varepsilon$  stands for the Hessian matrix of  $h_\varepsilon$  with respect to  $v$ . We refer to the choice (2.8) for  $\mathbb{S}_n$  as the stochastic Newton (SN) algorithm. Unfortunately, from a computational point of view, a major limitation of this SN algorithm is the need to compute the inverse of  $\mathbb{S}_n$  at each iteration  $n$  in equation (2.6). As  $\mathbb{S}_n$  is given by the recursive equation (2.8), it is tempting to use the Sherman-Morrison-Woodbury (SMW) formula [27], that is recalled in Lemma A.1, in order to compute  $\mathbb{S}_n^{-1}$  from the knowledge of  $\mathbb{S}_{n-1}^{-1}$  in a recursive manner. However, as detailed in Section A.2 of Appendix A, the Hessian matrix  $\nabla_v^2 h_\varepsilon(X_n, \widehat{V}_{n-1})$  does not have a sufficiently low-rank structure that would lead to a fast recursive approach to compute  $\mathbb{S}_n^{-1}$ . Therefore, for the SN algorithm, the computational cost to evaluate  $\mathbb{S}_n^{-1}$  appears to be of order  $\mathcal{O}(J^3)$  which only leads to a feasible algorithm for very small values of  $J$ . This important computational limitation then drew our investigation towards the SGN algorithm instead of the SN approach.

### 2.3 The stochastic Gauss-Newton algorithm

Historically, the Gauss-Newton adaptation of the Newton algorithm consists in replacing the Hessian matrix  $\nabla_v^2 h_\varepsilon(X_n, \widehat{V}_{n-1})$  by a tensor product of the gradient  $\nabla_v h_\varepsilon(X_n, \widehat{V}_{n-1})$ . In our framework, it leads to another pre-conditionned stochastic algorithm. We introduce  $S_n$  recursively as

$$\begin{aligned} S_n &= I_J + \sum_{k=1}^n \nabla_v h_\varepsilon(X_k, \widehat{V}_{k-1}) \nabla_v h_\varepsilon(X_k, \widehat{V}_{k-1})^T + \gamma \left(1 + \left\lfloor \frac{k}{J} \right\rfloor\right)^{-\beta} Z_k Z_k^T \\ &= S_{n-1} + \nabla_v h_\varepsilon(X_n, \widehat{V}_{n-1}) \nabla_v h_\varepsilon(X_n, \widehat{V}_{n-1})^T + \gamma \left(1 + \left\lfloor \frac{n}{J} \right\rfloor\right)^{-\beta} Z_n Z_n^T, \end{aligned} \quad (2.9)$$

for some constants  $\gamma > 0$  and  $0 < \beta < 1/2$ , and where  $(Z_1, \dots, Z_n)$  is a *deterministic* sequence of vectors defined, for all  $1 \leq k \leq n$ , by

$$Z_k = \sqrt{\nu_{\ell_k}} e_{\ell_k}$$

with  $\ell_k = 1 + (k-1) \pmod{J}$ , where  $(e_1, \dots, e_J)$  stands for the canonical basis of  $\mathbb{R}^J$ . We shall refer to the choice (2.9) for  $S_n$  as the regularized stochastic Gauss-Newton (SGN) algorithm and from now on, the notation  $S_n$  refers to this definition. We also use the convention that  $S_0 = I_J$ .

## 3 Main results on the SGN algorithm

Throughout this section, we investigate the statistical properties of the recursive sequence  $(\widehat{V}_n)$  defined by (2.6) with  $0 \leq \alpha < 1/2$ , where  $(S_n)$  is the sequence of random matrices defined by (2.9) with  $0 < \beta < 1/2$  that yields the SGN algorithm. The initial value  $\widehat{V}_0$  is assumed to be a square integrable random vector that belongs to  $\langle \mathbf{v}_J \rangle^\perp$ . Then, thanks to the projection step in equation (2.6), it follows that for all  $n \geq 1$ ,  $\widehat{V}_n$  also belongs to  $\langle \mathbf{v}_J \rangle^\perp$ . To derive the convergence properties of the SGN algorithm, we first need to introduce the matrix-valued function  $G_\varepsilon(v)$  defined as

$$G_\varepsilon(v) = \mathbb{E} [\nabla_v h_\varepsilon(X, v) \nabla_v h_\varepsilon(X, v)^T] \quad (3.1)$$

that will be shown to be a key quantity to analyze the SGN algorithm. In particular, we shall derive our results under the following assumption on the smallest eigenvalue of

$G_\varepsilon(v^*)$  associated to eigenvectors belonging to  $\langle \mathbf{v}_J \rangle^\perp$ .

**Invertibility assumption.** The matrix  $G_\varepsilon(v^*)$  satisfies

$$\min_{v \in \langle \mathbf{v}_J \rangle^\perp} \left\{ \frac{v^T G_\varepsilon(v^*) v}{\|v\|^2} \right\} > 0.$$

In all the sequel, we suppose that this invertibility assumption is satisfied. We denote by  $G_\varepsilon^-(v^*)$  the Moore-Penrose inverse of  $G_\varepsilon(v^*)$  and by  $G_\varepsilon^{-1/2}(v^*)$  its square-root. We now discuss, in what follows, the next keystone inequality.

**Proposition 3.1.** *Assume that the regularization parameter  $\varepsilon > 0$  satisfies*

$$\varepsilon \leq \frac{\min(\nu)}{\max(\nu) - \min(\nu)}. \quad (3.2)$$

*Then, in the sense of partial ordering between positive semi-definite matrices, we have*

$$G_\varepsilon(v^*) \leq \nabla^2 H_\varepsilon(v^*). \quad (3.3)$$

Inequality (3.3) is an important property of the SGN algorithm to prove its adaptivity to the geometry of the stochastic optimization problem (2.3). Of course, one can observe that no hyperparameter depending on the Hessian of  $H_\varepsilon$  needs to be tuned to run this algorithm provided that condition (3.2) holds. One can also remark that there is no restriction on the regularization parameter  $\varepsilon$  when  $\nu$  is the uniform distribution, that is when  $\nu_j = 1/J$ , for all  $1 \leq j \leq J$ , implying that  $\max(\nu) = \min(\nu)$ . Throughout the paper, we suppose that condition (3.2) holds true. Below, we denote by  $\lambda_{\min}^{\langle \mathbf{v}_J \rangle^\perp}(A)$  the smallest non-zero eigenvalue of a positive semi-definite matrix  $A$ , when the associated eigenvectors belong to  $\langle \mathbf{v}_J \rangle^\perp$ , the orthogonal complement of  $\mathbf{v}_J$ .

It immediately follows from inequality (3.3) that

$$1 \leq \lambda_{\min}^{\langle \mathbf{v}_J \rangle^\perp} (G_\varepsilon^{-1/2}(v^*) \nabla^2 H_\varepsilon(v^*) G_\varepsilon^{-1/2}(v^*)). \quad (3.4)$$

Inequality (3.4) will be a key property in this paper to derive the rates of convergence of the estimators obtained from the SGN algorithm. Note that the (pseudo) inverse of the Hessian matrix  $\nabla^2 H_\varepsilon(v^*)$  somehow represents an ideal deterministic pre-conditioning matrix, whose use would lead to the second order Newton algorithm: this ideal pre-conditioned algorithm is *non-adaptive* since it requires the use of  $\nabla^2 H_\varepsilon(v^*)$ , which is unknown in practice.

Indeed, in our SGN algorithm, adaptivity is tightly related to the limiting recursion induced by Equation (2.6). If we admit (temporarily) the almost sure convergence of the SGN algorithm towards  $v^*$  and of  $n^{-1}S_n$  towards  $G_\varepsilon(v^*)$ , the recursion induced by (2.6) looks very similar to a discretization of a dynamical system with a step size  $n^{-(1-\alpha)}$  and with a limiting linearized drift of the form  $-G_\varepsilon(v^*)^{-1} \nabla^2 H_\varepsilon(v^*) (v - v^*)$ . For further details on this point, we refer to the so-called ODE method (see *e.g.* [3]). The keystone property induced by Proposition 3.1 is that thanks to Equation (3.3) and Equation (3.4), the linearized drift of the limiting deterministic dynamical system has its eigenvalues that are lower bounded by 1, regardless of the value of the Hessian matrix  $\nabla^2 H_\varepsilon(v^*)$ . This translates an adaptation of the algorithm to the curvature of  $H_\varepsilon$  near the target point  $v^*$ . Therefore, the matrix  $G_\varepsilon(v^*)$  that is learnt on-line, and automatically adapts to the eigenspaces associated to the smallest eigenvalues of  $\nabla^2 H_\varepsilon(v^*)$ . Therefore, one may interpret (3.4) as the adaptivity of the SGN algorithm to the geometry of the semi-dual formulation of regularized OT.



### 3.1 Almost sure convergence

The almost sure convergence of the sequences  $(\widehat{V}_n)$ ,  $(\widehat{W}_n)$  and  $(\overline{S}_n)$  are as follows where  $\overline{S}_n = \frac{1}{n}S_n$ .

**Theorem 3.1.** *Assume that  $\alpha \in [0, 1/2[$  and  $\alpha + \beta < 1/2$ . Then, we have*

$$\lim_{n \rightarrow +\infty} \widehat{V}_n = v^* \quad a.s. \quad (3.5)$$

and

$$\lim_{n \rightarrow +\infty} \overline{S}_n = G_\varepsilon(v^*) \quad a.s. \quad (3.6)$$

The following result is an immediate corollary of Theorem 3.1, thanks to the continuity of the function  $h_\varepsilon$ .

**Corollary 3.1.** *Assume that  $\alpha \in [0, 1/2[$  and  $\alpha + \beta < 1/2$ . Suppose that the cost function  $c$  satisfies, for any  $1 \leq j \leq J$ ,*

$$\int_{\mathcal{X}} c^2(x, y_j) d\mu(x) < +\infty. \quad (3.7)$$

Then, we have

$$\lim_{n \rightarrow +\infty} \widehat{W}_n = W_\varepsilon(\mu, \nu) \quad a.s.$$

We now derive results on the almost sure rates of convergence of the sequences  $(\widehat{V}_n)$  and  $(\overline{S}_n)$  that are the keystone in the proof of the asymptotic normality of the estimator  $\widehat{V}_n$  studied in Section 3.2. We emphasize that we restrict our study to the case  $\alpha = 0$ , which yields the fastest rates of convergence and that corresponds to the meaningful situation from the numerical point of view.

**Theorem 3.2.** *Assume that  $\alpha = 0$ . Then, we have the almost sure rate of convergence*

$$\|\widehat{V}_n - v^*\|^2 = \mathcal{O}\left(\frac{\log n}{n}\right) \quad a.s. \quad (3.8)$$

In addition, we also have

$$\|\overline{S}_n - G_\varepsilon(v^*)\|_F = \mathcal{O}\left(\frac{1}{n^\beta}\right) \quad a.s. \quad (3.9)$$

and

$$\|\overline{S}_n^{-1} - G_\varepsilon^-(v^*)\|_F = \mathcal{O}\left(\frac{1}{n^\beta}\right) \quad a.s. \quad (3.10)$$

### 3.2 Asymptotic normality

The asymptotic normality of our estimates depends on the magnitude of the smallest eigenvalue (associated to eigenvectors belonging to  $\langle \mathbf{v}_J \rangle^\perp$ ) of the matrix

$$\Gamma_\varepsilon(v^*) = G_\varepsilon^{-1/2}(v^*) \nabla^2 H_\varepsilon(v^*) G_\varepsilon^{-1/2}(v^*). \quad (3.11)$$



Thanks to the key inequality (3.4), we have that the smallest eigenvalue of  $\Gamma_\varepsilon(v^*)$  is always greater than 1, in the sense that

$$\min_{v \in \langle v_J \rangle^\perp} \frac{v^T \Gamma_\varepsilon(v^*) v}{\|v\|^2} \geq 1.$$

One can observe that we also restrict our study to the case  $\alpha = 0$  which yields the usual  $\sqrt{n}$  rate of convergence for the central limit theorem that is stated below.

**Theorem 3.3.** *Assume that  $\alpha = 0$ . Then, we have the asymptotic normality*

$$\sqrt{n}(\widehat{V}_n - v^*) \xrightarrow{\mathcal{L}} \mathcal{N}\left(0, G_\varepsilon^{-1/2}(v^*)(2\Gamma_\varepsilon(v^*) - P_J)^- G_\varepsilon^{-1/2}(v^*)\right). \quad (3.12)$$

In addition, suppose that the cost function  $c$  satisfies, for any  $1 \leq j \leq J$ ,

$$\int_{\mathcal{X}} c^4(x, y_j) d\mu(x) < +\infty. \quad (3.13)$$

Then, we also have

$$\sqrt{n}(\widehat{W}_n - W_\varepsilon(\mu, \nu)) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \sigma_\varepsilon^2) \quad (3.14)$$

where the asymptotic variance  $\sigma_\varepsilon^2 = \mathbb{E}[h_\varepsilon^2(X, v^*)] - W_\varepsilon^2(\mu, \nu)$ .

In order to discuss the above result on the asymptotic normality of  $\widehat{V}_n$ , we denote by

$$\Sigma_\varepsilon(v^*) = G_\varepsilon^{-1/2}(v^*)(2\Gamma_\varepsilon(v^*) - P_J)^- G_\varepsilon^{-1/2}(v^*)$$

the asymptotic covariance matrix in (3.12). One can check that  $\Sigma_\varepsilon(v^*)$  satisfies the Lyapunov equation

$$\left(\frac{1}{2}P_J - A\nabla^2 H_\varepsilon(v^*)\right)\Sigma_\varepsilon(v^*) + \Sigma_\varepsilon(v^*)\left(\frac{1}{2}P_J - A\nabla^2 H_\varepsilon(v^*)\right)^T = -AG_\varepsilon(v^*)A \quad (3.15)$$

with  $A = G_\varepsilon^-(v^*)$ . Moreover, one can observe that

$$G_\varepsilon(v^*) = \lim_{n \rightarrow +\infty} \mathbb{E}[\varepsilon_{n+1}\varepsilon_{n+1}^T | \mathcal{F}_n] \quad \text{a.s.}$$

is the asymptotic covariance matrix of the martingale increment  $\varepsilon_{n+1} = \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n) - \nabla H_\varepsilon(\widehat{V}_n)$ . Hence, to better interpret the asymptotic covariance matrix  $\Sigma_\varepsilon(v^*)$ , let us consider the following sub-class of pre-conditionned stochastic algorithms

$$\widetilde{V}_{n+1} = P_J\left(\widetilde{V}_n - \frac{1}{n}A\nabla_v h_\varepsilon(X_{n+1}, \widetilde{V}_n)\right), \quad (3.16)$$

where  $A$  is a deterministic positive semi-definite matrix satisfying the stability condition

$$A\nabla^2 H_\varepsilon(v^*) \geq \frac{1}{2}P_J. \quad (3.17)$$

Then, adapting well-known results on stochastic optimisation (see e.g. [18, 37]), one may prove that

$$\sqrt{n}(\widetilde{V}_n - v^*) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \Sigma(A))$$

where  $\Sigma(A)$  is the solution of the Lyapunov equation (3.15) with  $\Sigma(A)$  instead of  $\Sigma_\varepsilon(v^*)$ . Hence, the asymptotic normality of the SGN algorithm coincides with the one of the pre-conditioned stochastic algorithm (3.16) for the choice  $A = G_\varepsilon^-(v^*)$ . Hence, the main advantage of the SGN algorithm is to be fully data-driven as  $G_\varepsilon^-(v^*)$  is obviously unknown. Among the deterministic pre-conditioning matrices satisfying condition (3.17), it is also known (see e.g. [18, 37]) that the best choice is to take  $A = \nabla^2 H_\varepsilon^-(v^*)$  that corresponds to an ideal Newton algorithm and which yields the optimal asymptotic covariance matrix

$$\Sigma^* = \Sigma(\nabla^2 H_\varepsilon(v^*)) = \nabla^2 H_\varepsilon(v^*)^- G_\varepsilon(v^*) \nabla^2 H_\varepsilon(v^*)^- \leq \Sigma_\varepsilon(v^*).$$

Therefore, the SGN algorithm does not yield an estimator  $\widehat{V}_n$  having an asymptotically optimal covariance matrix. Note that, as shown in [5, Theorem 3.4], using an average version of the standard SGD algorithm, that is with  $S_n = s^{-1}nI_J$  and  $0 < \alpha < 1/2$ , allows to obtain an estimator having an asymptotic distribution with optimal covariance matrix  $\Sigma^*$ . However, in numerical experiments, it appears that the choice of  $s$  for the averaged SGD algorithm is crucial but difficult to tune. The results from [5] suggests to take  $s = \varepsilon/(2\min(\nu))$  which follows from the property that

$$\lambda_{\min}(\nabla^2 H_\varepsilon(v^*)) \geq \frac{\min(\nu)}{\varepsilon},$$

that is discussed in Section 5. Hence, the choice  $s = \varepsilon/(2\min(\nu))$  ensures that the pre-conditioning matrix  $A = s^{-1}I_J$  satisfies the stability condition (3.17). However, as shown by the numerical experiments carried out in Section 4, it appears that the SGN algorithm automatically adapts to the geometry of the optimisation problem with better results than the SGD algorithm. Finally, we remark from the asymptotic normality (3.14) and [5, Theorem 3.5] that the asymptotic variance of the recursive estimator  $\widehat{W}_n$  is the same when  $\widehat{V}_n$  is either computed using the SGN or the SGD algorithm.

### 3.3 Non-asymptotic rates of convergence

The last contribution of our paper is to derive non-asymptotic upper bounds on the expected risk of various estimators arising from the use of the SGN algorithm when  $(S_n)$  is the sequence of positive definite matrices defined by (2.9). In particular, we derive the rate of convergence of the expected quadratic risks

$$\mathbb{E}[\|\widehat{V}_n - v^*\|^2] \quad \text{and} \quad \mathbb{E}[\|\overline{S}_n - G_\varepsilon(v^*)\|_F^2].$$

We also analyze the rate of convergence of the expected excess risk  $W_\varepsilon(\mu, \nu) - \mathbb{E}[\widehat{W}_n]$  of the recursive estimator  $\widehat{W}_n$  defined by (2.7) used to approximate the regularized OT cost  $W_\varepsilon(\mu, \nu)$ .

**Theorem 3.4.** *Assume that  $\alpha \in ]0, 1/2[$  and that  $4\beta < 1 - 2\alpha$ . Then, there exists a positive constant  $c_\varepsilon$  such that for any  $n \geq 1$ ,*

$$\mathbb{E}[\|\widehat{V}_n - v^*\|^2] \leq \frac{c_\varepsilon}{n^{1-\alpha}} \quad \text{and} \quad \mathbb{E}[\|\overline{S}_n - G_\varepsilon(v^*)\|_F^2] \leq \frac{c_\varepsilon}{n^{2\beta}}. \quad (3.18)$$

Moreover, we also have

$$|\mathbb{E}[\widehat{W}_n] - W_\varepsilon(\mu, \nu)| \leq \frac{c_\varepsilon}{n^{1-\alpha}}, \quad (3.19)$$

and if the cost function  $c$  satisfies  $\int_{\mathcal{X}} c^2(x, y_j) d\mu(x) < +\infty$ , for any  $1 \leq j \leq J$ , then

$$\mathbb{E}[|\widehat{W}_n - W_\varepsilon(\mu, \nu)|] \leq \frac{c_\varepsilon}{\sqrt{n}}. \quad (3.20)$$

Note that the value of the constant  $c_\varepsilon$  appearing in Theorem 3.4 may also depend on  $\alpha$  and  $\beta$ , but we remove this dependency in the notation to simplify the presentation. One can observe that choosing  $\alpha > 0$  allows the algorithm to be fully adaptative in the sense that no important hyperparameter needs to be tune to obtain non-asymptotic rates of convergence. The case  $\alpha = 0$  could also be considered but this will require to introduce a multiplicative positive constant  $c$  in the definition of the SGN algorithm by replacing  $n^\alpha S_n^{-1}$  in equation (2.6) by  $c S_n^{-1}$ . Then, provided that  $c$  is sufficiently large, one may obtain faster rate of convergence for the expected quadratic risk of the order  $\mathcal{O}(\log(n)/n)$ . However, in our numerical experiments, we have found that introducing such a large multiplicative constant  $c$  makes the convergence of the SGN algorithm too slow. Therefore, results on non-asymptotic convergence rates in the case  $\alpha = 0$  are not reported here.

## 4 Implementation of the SGN algorithm and numerical experiments

In this section, we first discuss computational considerations on the implementation of the SGN algorithm, and we also make several remarks to justify its use. Then, we report the results of numerical experiments.

### 4.1 A fast recursive approach to compute $S_n^{-1}$ .

In this paragraph, we discuss on the computational benefits of using the Gauss-Newton method as an alternative to the Newton algorithm. A key point to define the SGN algorithm consists in replacing in equation (2.8) that defines the SN algorithm, the positive definite Hessian matrices  $\nabla_v^2 h_\varepsilon(X_k, \widehat{V}_{k-1})$  by the tensor product  $\nabla_v h_\varepsilon(X_k, \widehat{V}_{k-1}) \nabla_v h_\varepsilon(X_k, \widehat{V}_{k-1})^T$  of the gradient of  $h_\varepsilon$  at  $(X_k, \widehat{V}_{k-1})$ . A second important ingredient in the definition (2.9) of the SGN algorithm is the additive regularization terms  $\gamma (1 + \lfloor \frac{k}{J} \rfloor)^{-\beta} Z_k Z_k^T$  whose role is discussed in the sub-section below.

Note that  $Z_k Z_k^T = \nu_{\ell_k} e_{\ell_k} e_{\ell_k}^T$  (with  $\ell_k = (k-1) \pmod J + 1$ ) is a diagonal matrix, such that all its diagonal elements are equal to zero, except the  $\ell_k$ -th one which is equal to  $\nu_{\ell_k}$ . In this manner, the difference  $S_n - S_{n-1} = \phi_n \phi_n^T + \gamma (1 + \lfloor \frac{n}{J} \rfloor)^{-\beta} Z_n Z_n^T$  is thus the sum of two rank one matrices, where  $\phi_n = \nabla_v h_\varepsilon(X_n, \widehat{V}_{n-1})$ . Therefore, one may easily obtain  $S_n^{-1}$  from the knowledge of  $S_{n-1}^{-1}$  as follows. Introducing the intermediate matrix  $S_{n-\frac{1}{2}} = S_{n-1} + \gamma (1 + \lfloor \frac{n}{J} \rfloor)^{-\beta} Z_n Z_n^T$ , we observe that  $S_n = S_{n-\frac{1}{2}} + \phi_n \phi_n^T$ . Consequently, by applying the SMW formula (A.3), we first notice that

$$\begin{aligned} S_{n-\frac{1}{2}}^{-1} &= (S_{n-1} + \gamma (1 + \lfloor \frac{n}{J} \rfloor)^{-\beta} Z_n Z_n^T)^{-1} \\ &= S_{n-1}^{-1} - (Z_n^T S_{n-1}^{-1} Z_n + \gamma^{-1} (1 + \lfloor \frac{n}{J} \rfloor)^\beta)^{-1} S_{n-1}^{-1} Z_n Z_n^T S_{n-1}^{-1} \end{aligned}$$

Using that  $Z_n = \sqrt{\nu_{\ell_n}} e_{\ell_n}$ , we furthermore have that

$$S_{n-\frac{1}{2}}^{-1} = S_{n-1}^{-1} - \nu_{\ell_n} \frac{(S_{n-1}^{-1})_{\cdot, \ell_n} (S_{n-1}^{-1})_{\ell_n, \cdot}^T}{\nu_{\ell_n} (S_{n-1}^{-1})_{\ell_n, \ell_n} + \gamma^{-1} (1 + \lfloor \frac{n}{J} \rfloor)^\beta}. \quad (4.1)$$

Secondly, applying again the SMW formula (A.3), we obtain that

$$S_n^{-1} = S_{n-\frac{1}{2}}^{-1} - \frac{S_{n-\frac{1}{2}}^{-1} \phi_n \phi_n^T S_{n-\frac{1}{2}}^{-1}}{\phi_n^T S_{n-\frac{1}{2}}^{-1} \phi_n + 1}. \quad (4.2)$$

Hence, the recursive formulas (4.1) and (4.2) allow, at each iteration  $n$ , a much more faster computation of  $S_n^{-1}$  from the knowledge of  $S_{n-1}^{-1}$ , which is a key advantage of the SGN algorithm over the use of the SN algorithm. Indeed, the cost of computing  $S_n^{-1}$  using the above recursive formulas is that of matrix vector multiplication which is of order  $\mathcal{O}(J^2)$ .

## 4.2 The role of regularization.

Let us denote by  $R_n = \sum_{k=1}^n \gamma (1 + \lfloor \frac{k}{J} \rfloor)^{-\beta} Z_k Z_k^T$  the sum of the deterministic regularization terms in (2.9) implying that  $S_n$  can be decomposed as

$$S_n = I_J + \sum_{k=1}^n \nabla_v h_\varepsilon(X_k, \hat{V}_{k-1}) \nabla_v h_\varepsilon(X_k, \hat{V}_{k-1})^T + R_n. \quad (4.3)$$

If  $n = pJ$  for some integer  $p \geq 1$ , the regularization by the matrices  $\gamma (1 + \lfloor \frac{k}{J} \rfloor)^{-\beta} Z_k Z_k^T$  in (2.9) sum up to a simple expression given by

$$R_n = \left( \sum_{m=1}^p m^{-\beta} \right) \gamma \text{diag}(\nu).$$

The following two important comments can be made to clarify the role of this regularization effect:

- adding the supplementary matrix  $R_n$  in (4.3) implies that  $S_n$  is invertible as soon as  $n \geq J$  with a known lower bound on its smallest eigenvalue. Indeed, thanks to the condition  $0 < \beta < 1/2$  and to the property that

$$\left( \frac{1}{p} \sum_{m=1}^p m^{-\beta} \right) \sim \frac{1}{1-\beta} p^{-\beta},$$

the additive term  $R_n$  allows to regularize the smallest eigenvalue of  $S_n$ . This is important for the evolution of the stochastic algorithm: this regularization allows to show that  $\hat{V}_n$  converges almost surely to  $v^*$ . More precisely, while  $\hat{V}_{n+1} - \hat{V}_n$  is essentially modified in the direction  $-n^\alpha S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n)$ , it is well known that too large step sizes are prohibited to obtain a good behavior of stochastic algorithms. Therefore, taking a sufficiently small  $\beta$  guarantees a suitable upper bound of the increments of the SGN, that in turn limits the effect of the noise at each iteration of the algorithm.

- the growth of  $R_n$  is sublinear for large values of  $n$  whereas

$$\sum_{k=1}^n \nabla_v h_\varepsilon(X_k, \widehat{V}_{k-1}) \nabla_v h_\varepsilon(X_k, \widehat{V}_{k-1})^T$$

grows linearly with  $n$  so that this last term will become dominant in the decomposition of  $S_n$ , inducing a “learning” of the curvature of the landscape function  $H_\varepsilon$ . Recalling that  $\overline{S}_n = \frac{1}{n} S_n$ , it will be shown in Section 6 that

$$\lim_{n \rightarrow +\infty} \overline{S}_n = G_\varepsilon(v^*) \quad \text{a.s.}$$

with

$$G_\varepsilon(v^*) = \mathbb{E} [\nabla_v h_\varepsilon(X, v^*) \nabla_v h_\varepsilon(X, v^*)^T] = \text{diag}(\nu) - \nu \nu^T - \varepsilon \nabla^2 H_\varepsilon(v^*),$$

where the last equality above follows the proof of Proposition 3.1. Hence, when  $n \rightarrow +\infty$ , the regularization disappears as long as  $\beta > 0$ . Note that this would not be the case if  $\beta$  was chosen to be equal to 0.

To sum up, taking  $\beta \in (0, 1/2)$  will be a crucial assumption to derive the almost sure convergence rates that are stated in Theorem 3.2.

### 4.3 Numerical experiments

In this section, we report numerical results on the performances of stochastic algorithms for regularized optimal transport when the source measure  $\mu$  is either discrete or absolutely continuous. We shall compare the SGD, ADAM, SGN and SN algorithms. For the SGD algorithm, following the results in [5], we took  $\alpha = 1/2$  and  $S_n = s^{-1} n I_J$  with  $s = \varepsilon / (2 \min(\nu))$ . The ADAM algorithm has been implemented following the parametrization made in the seminal paper [28] except the value of the stepsize (as defined in [28, Algorithm 1]) that is set to 0.005 instead of 0.001, which improves the performances of ADAM in our numerical experiments. For the SGN algorithm, we set  $\alpha = 0$  and we have taken  $\gamma = 10^{-3}$  (a small value) and  $\beta = 0.49$ . For the results reported in this paper, we have found that the performances of the SGN algorithm are not very sensitive to the value of  $\beta \in (0, 1/2)$ . Finally, for the SN algorithm, we chose  $\alpha = 0$  and  $S_n = \mathbb{S}_n$  as defined by (2.8).

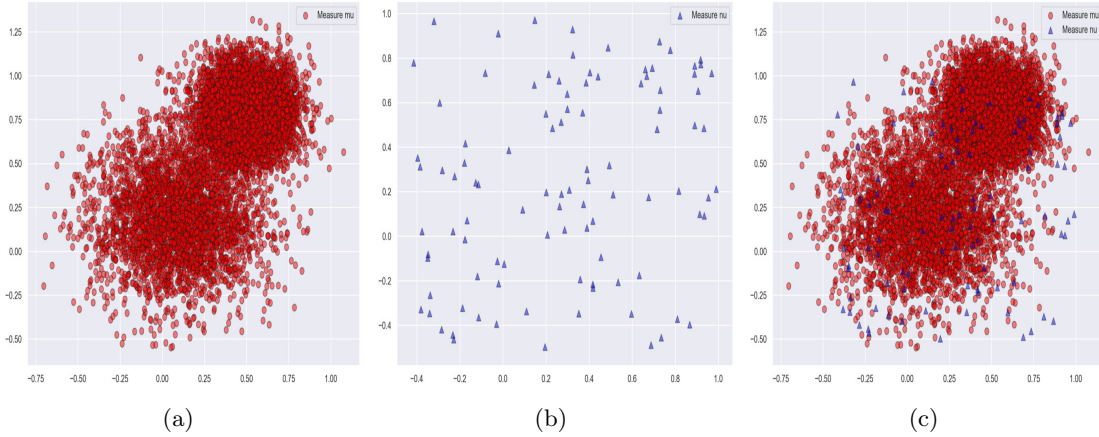
In the discrete setting, we shall also compare the performances of these stochastic algorithms to those of the Sinkhorn algorithm [15], which is a deterministic iterative procedure that uses the full knowledge of the measures  $\mu$  and  $\nu$  at each iteration. Let us recall that, for the SGD and ADAM algorithms, the computational cost of one iteration from  $n$  to  $n + 1$  is of order  $\mathcal{O}(J)$ , while it is of order  $\mathcal{O}(J^2)$  for the SGN algorithm and  $\mathcal{O}(J^3)$  for the SN algorithm. Each iteration of the Sinkhorn algorithm is of order  $\mathcal{O}(IJ)$ , where  $I$  denotes the size of the support of  $\mu$  in the discrete setting.

In these numerical experiments, we investigate the numerical behavior of the recursive estimators  $\widehat{W}_n$  and  $\widehat{V}_n$ . The performances of the various stochastic algorithms used to compute these estimators are compared in terms of the expected excess risks  $\mathbb{E}[|\widehat{W}_n - W_\varepsilon(\mu, \nu)|]$  and  $\mathbb{E}[\|\widehat{V}_n - v^*\|^2]$ . For the SGN algorithm, we also analyze the convergence of the estimator  $\overline{S}_n$  to the matrix  $G_\varepsilon(v^*)$ . The expected value involved in these expected risks is approximated using 100 Monte-Carlo replications. When the measure  $\mu$  is discrete,

we use the Sinkhorn algorithm [15] to preliminary compute  $W_\varepsilon(\mu, \nu)$  and  $v^*$ . When  $\mu$  is absolutely continuous, the regularized OT cost  $W_\varepsilon(\mu, \nu)$  is preliminary approximated by running the SN algorithm with a very large value of iterations (e.g.  $n = 10^6$ ). To the best of our knowledge, apart from stochastic approaches as in [23], there is no other method to evaluate  $W_\varepsilon(\mu, \nu)$  in the semi-discrete setting. Note that we shall compare the evolution of these excess risks as a function of the computational time (observed on the computer) of each algorithm. Moreover, the estimators  $\widehat{W}_n$  and  $\widehat{V}_n$  obviously depends on the regularization parameter  $\varepsilon$ . However, for the sake of simplicity, we have chosen to denote them as  $\widehat{W}_n$  and  $\widehat{V}_n$ , although we carry out numerical experiments for different values of  $\varepsilon$ . Finally, we also analyze the asymptotic distributions of  $\widehat{W}_n$  and  $\widehat{V}_n$  to illustrate the results on asymptotic normality given in Section 3.2.

#### 4.3.1 Discrete setting in dimension $d = 2$

In this section, the cost function is chosen as the squared Euclidean distance that is  $c(x, y) = \|x - y\|^2$ . We focus our attention when both  $\mu = \sum_{i=1}^I \mu_i \delta_{x_i}$  and  $\nu = \sum_{j=1}^J \nu_j \delta_{y_j}$  are uniform discrete measures supported on  $\mathbb{R}^2$ , that is  $\mu_i = \frac{1}{I}$  and  $\nu_j = \frac{1}{J}$ . The points  $(x_i)_{1 \leq i \leq I}$  (resp.  $(y_j)_{1 \leq j \leq J}$ ) are drawn randomly (once for all) from a Gaussian mixture with two components (resp. from the uniform distribution on  $[0, 1]^2$ ). An example of two such measures is displayed in Figure 1 for  $I = 10^4$  and  $J = 100$ . The number of iterations of the four stochastic algorithms is fixed to  $n = 10^5$  except for the experiments on the asymptotic distribution of  $\widehat{W}_n$  and  $\widehat{V}_n$ , where  $n$  is let being larger. Finally, the Sinkhorn algorithm is let running until convergence is reached to provide a reference value for  $W_\varepsilon(\mu, \nu)$  and  $v^*$  considered as the ground truth.



**Figure 1:** (a) Discrete measure  $\mu$  supported on  $I = 10^4$  points drawn from a mixture of two Gaussian distributions, and (b) discrete measure  $\nu$  supported on  $J = 100$  points randomly drawn from the uniform distribution on  $[0, 1]^2$ . (c) Superposition of  $\mu$  and  $\nu$ .

**Convergence of the excess risks.** We first report results for  $I = 10^4$  (size of the support of  $\mu$ ) and  $J \in \{100, 400\}$  (size of the support of  $\nu$ ), and two small values of the regularization

parameter  $\varepsilon \in \{0.01, 0.005\}$ . For different combinations of these hyperparameters, we display from Figure 2 to Figure 5 the value of the expected excess risks (in logarithmic scale)  $\mathbb{E}[\|\widehat{W}_n - W_\varepsilon(\mu, \nu)\|]$  and  $\mathbb{E}[\|\widehat{V}_n - v^*\|^2]$  as functions of the averaged (along the 100 Monte Carlo replications) computational time of each iteration of the stochastic algorithms. We also draw the evolution of the metrics (in logarithmic scale)  $|W_k - W_\varepsilon(\mu, \nu)|$  and  $\|V_k - v^*\|^2$  as functions of the computational time of the iterations of the Sinkhorn algorithm, where  $W_k \in \mathbb{R}$  and  $V_k \in \mathbb{R}^J$  are the output of the Sinkhorn algorithm at its  $k$ -th iteration.

In Figure 2 to Figure 5, the various curves are displayed as functions of the computational time until the convergence of the Sinkhorn algorithm is reached, that is until  $k = k_{\max}$  (the maximum number of Sinkhorn iterations). Note that for  $k \approx k_{\max}$  then  $|W_k - W_\varepsilon(\mu, \nu)| \approx 0$  and  $\|V_k - v^*\|^2 \approx 0$ . Hence, for such large values of  $k$ , these metrics have necessarily smaller values than those that are used to evaluate the stochastic algorithms. In the discussion that follows, we thus consider that the stochastic algorithms have reached convergence when the values of either  $\mathbb{E}[\|\widehat{W}_n - W_\varepsilon(\mu, \nu)\|]$  or  $\mathbb{E}[\|\widehat{V}_n - v^*\|^2]$  stabilize, although these metrics may be larger than the metrics used to evaluate the Sinkhorn algorithm for large values of the computational time. This is due to the randomness of the stochastic algorithms and their resulting positive variance (even for large values of  $n$ ). Then, the following comments can be made from the output of these numerical experiments.

- For  $I = 10^4$ ,  $J = 100$  and  $\varepsilon = 0.01$ , the four stochastic algorithms reach convergence faster than the Sinkhorn algorithm. The convergence is much faster for the metric  $\mathbb{E}[\|\widehat{V}_n - v^*\|^2]$  than for the metric  $\mathbb{E}[\|\widehat{W}_n - W_\varepsilon(\mu, \nu)\|]$ .
- For  $I = 10^4$ ,  $J = 100$  and  $\varepsilon = 0.005$ , SGD fails to converge either for the estimator  $\widehat{W}_n$  or the estimator  $\widehat{V}_n$ . For this smallest value of  $\varepsilon$ , SGN and SN have similar performances for the metric  $\mathbb{E}[\|\widehat{V}_n - v^*\|^2]$ . The SN algorithm is slightly better than SGN for the metric  $\mathbb{E}[\|\widehat{W}_n - W_\varepsilon(\mu, \nu)\|]$ . We also observe that SN and SGN converge much faster than Sinkhorn, and that they have better performances than ADAM for the two metrics.
- For  $I = 10^4$ ,  $J = 400$  and  $\varepsilon \in \{0.01, 0.005\}$ , it can be seen that the SGD algorithm does not converge. For the metric  $\mathbb{E}[\|\widehat{V}_n - v^*\|^2]$ , the convergence of the SN and SGN algorithms is much faster than Sinkhorn, and these two algorithms outperform ADAM.

Therefore, these numerical experiments suggest that the SGN algorithm has interesting benefits over the SGD, ADAM and Sinkhorn algorithms for moderate values of  $J$  and for small values of the regularization parameter  $\varepsilon$ . In these settings, SGN seems to be particularly relevant for the estimation of  $v^*$ , and it reaches performances similar to those of SN for the metric  $\mathbb{E}[\|\widehat{V}_n - v^*\|^2]$ . The SGN algorithm may also converge much faster than the Sinkhorn algorithm for either the estimation of  $W_\varepsilon(\mu, \nu)$  or  $v^*$  as the size  $I$  of the support of  $\mu$  is large.

**Asymptotic distribution of the stochastic algorithms.** Now, we illustrate the results from Section 3.2 on the asymptotic distributions of  $\widehat{W}_n$  and  $\widehat{V}_n$ . To this end, we consider the setting  $I = 10^3$  and  $J = 50$ . Then, we display in Figure 6 for  $\varepsilon = 0.1$  and  $n = 2 \times 10^5$



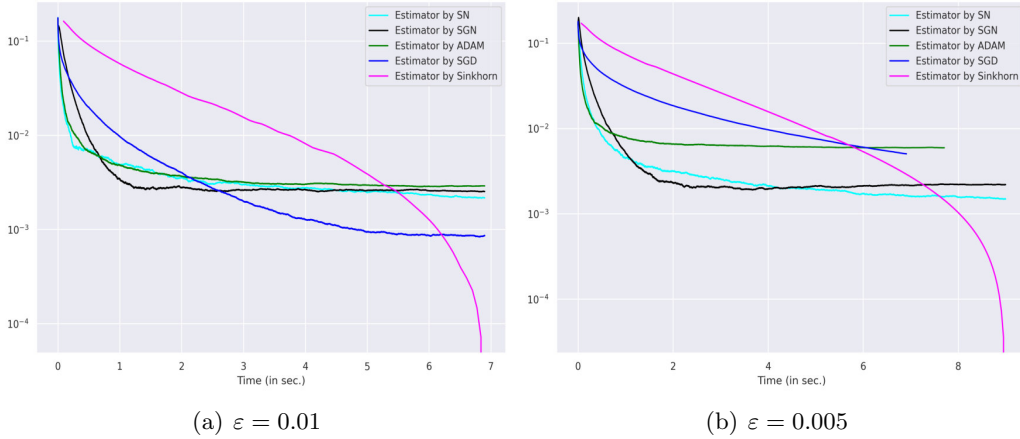
iterations (resp. Figure 7 for  $\varepsilon = 0.01$  and  $n = 4 \times 10^5$  iterations) the histograms of 200 independent realizations of

$$\widetilde{W}_n = \frac{\sqrt{n} \left( \widehat{W}_n - W_\varepsilon(\mu, \nu) \right)}{\widehat{\sigma}_n}$$

using each of the four stochastic algorithms, where

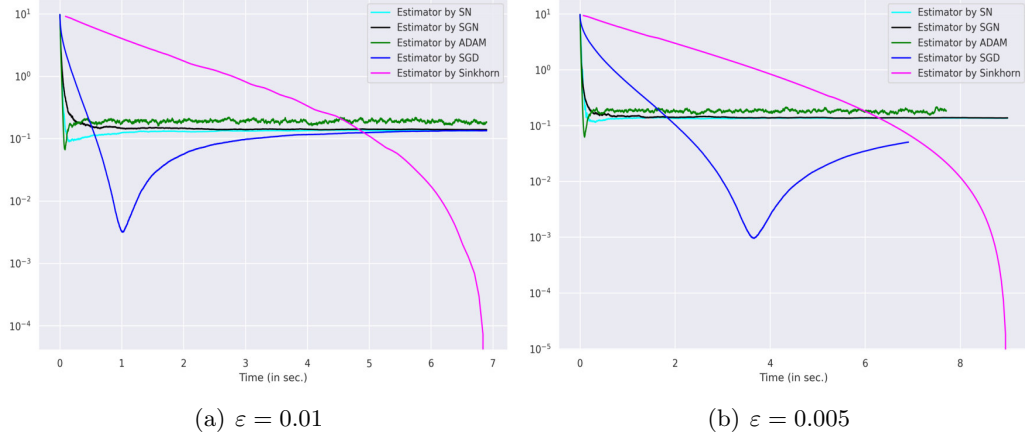
$$\widehat{\sigma}_n^2 = \frac{1}{n} \sum_{k=1}^n h_\varepsilon^2(X_k, \widehat{V}_{k-1}) - \widehat{W}_n^2,$$

is a recursive estimator of the asymptotic variance of  $\widehat{W}_n$  that has been introduced in [5]. For all the algorithms, it can be seen in Figure 6 and Figure 7 that  $\widehat{W}_n$  is normally distributed. For the SGD and the SN algorithms, the histograms of  $\widetilde{W}_n$  are very close to the standard Gaussian distribution, while the SGN is seen to be slightly biased. The bias is much more important for the ADAM algorithm.

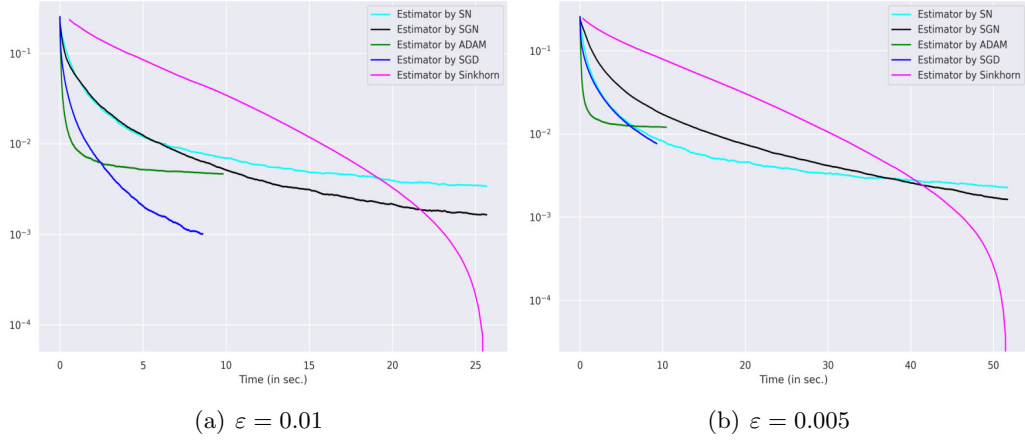


**Figure 2:** Discrete setting with  $I = 10^4$  and  $J = 100$  with  $n = 10^5$  iterations. Expected excess risk (in logarithmic scale)  $\log(\mathbb{E}[|\widehat{W}_n - W_\varepsilon(\mu, \nu)|])$  (resp. metric  $\log(|W_k - W_\varepsilon(\mu, \nu)|)$ ) as a function of the averaged computational cost of the iterations of the four stochastic algorithms (resp. the Sinkhorn algorithm) for different values of the regularization parameter  $\varepsilon$ .

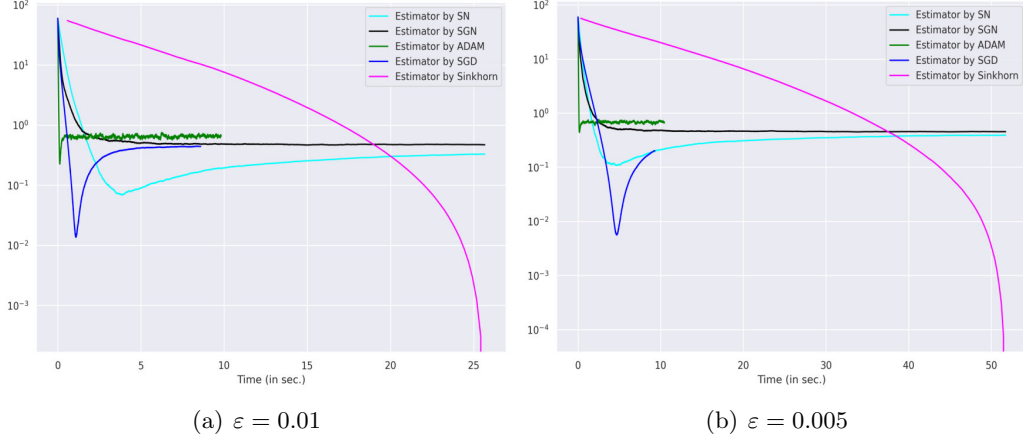
In Figure 8 and Figure 9, we also display the histograms of 200 independent realizations of  $\widetilde{V}_n = n \|\widehat{V}_n - v^*\|^2$  for the four stochastic algorithms. The distribution of  $\widetilde{V}_n$  has the shape of a  $\chi^2$ -distribution but the “number of degrees of freedom” is highly varying from one algorithm to the other. It can be seen from Figure 8 and Figure 9, that  $\widetilde{V}_n$  reaches its smallest variance for the SN algorithm, and that the second smallest variance is obtained with the SGN algorithm. The SGD and the ADAM algorithms finally have a much larger variance.



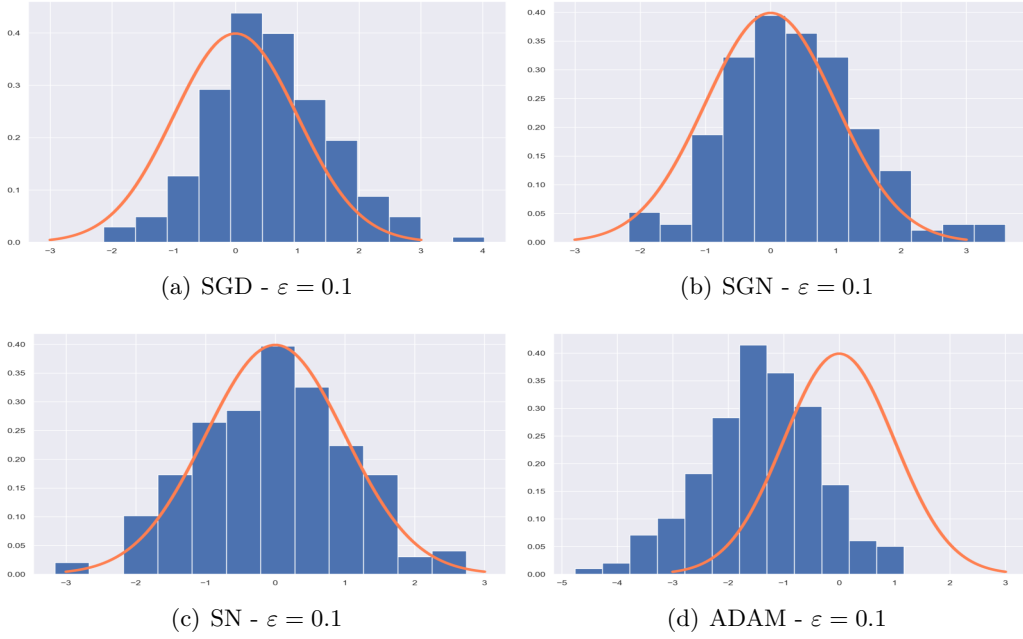
**Figure 3:** Discrete setting with  $I = 10^4$  and  $J = 100$  with  $n = 10^5$  iterations. Expected excess risk (in logarithmic scale)  $\log(\mathbb{E}[\|\widehat{V}_n - v^*\|^2])$  (resp. metric  $\log(\|V_k - v^*\|^2)$ ) as a function of the averaged computational cost of the iterations of the four stochastic algorithms (resp. the Sinkhorn algorithm) for different values of the regularization parameter  $\varepsilon$ .



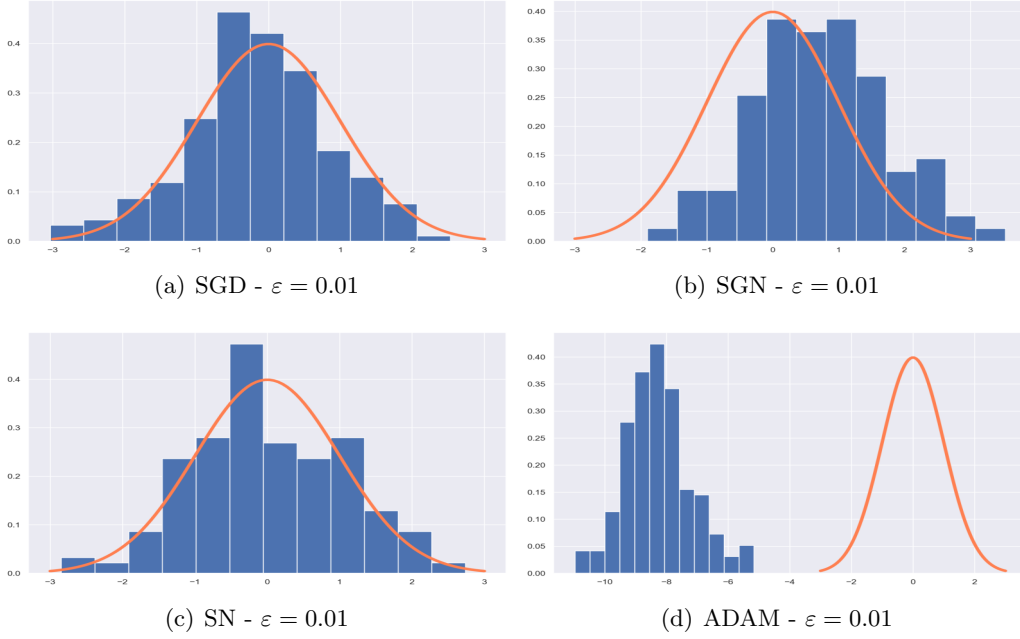
**Figure 4:** Discrete setting with  $I = 10^4$  and  $J = 400$  with  $n = 10^5$  iterations. Expected excess risk (in logarithmic scale)  $\log(\mathbb{E}[\|\widehat{W}_n - W_\varepsilon(\mu, \nu)\|])$  (resp. metric  $\log(\|W_k - W_\varepsilon(\mu, \nu)\|)$ ) as a function of the averaged computational cost of the iterations of the four stochastic algorithms (resp. the Sinkhorn algorithm) for different values of the regularization parameter  $\varepsilon$ .



**Figure 5:** Discrete setting with  $I = 10^4$  and  $J = 400$  with  $n = 10^5$  iterations. Expected excess risk (in logarithmic scale)  $\log(\mathbb{E}[\|\hat{V}_n - v^*\|^2])$  (resp. metric  $\log(\|\hat{V}_k - v^*\|^2)$ ) as a function of the averaged computational cost of the iterations of the four stochastic algorithms (resp. the Sinkhorn algorithm) for different values of the regularization parameter  $\varepsilon$ .



**Figure 6:** Discrete setting with  $I = 10^3$ ,  $J = 50$  and  $\varepsilon = 0.1$ . Histogram of 200 independent realizations of  $\frac{\sqrt{n}(\hat{W}_n - W_\varepsilon(\mu, \nu))}{\hat{\sigma}_n}$  with  $n = 2 \times 10^5$  using each of the four stochastic algorithms. The orange curve is the density of the standard Gaussian distribution.



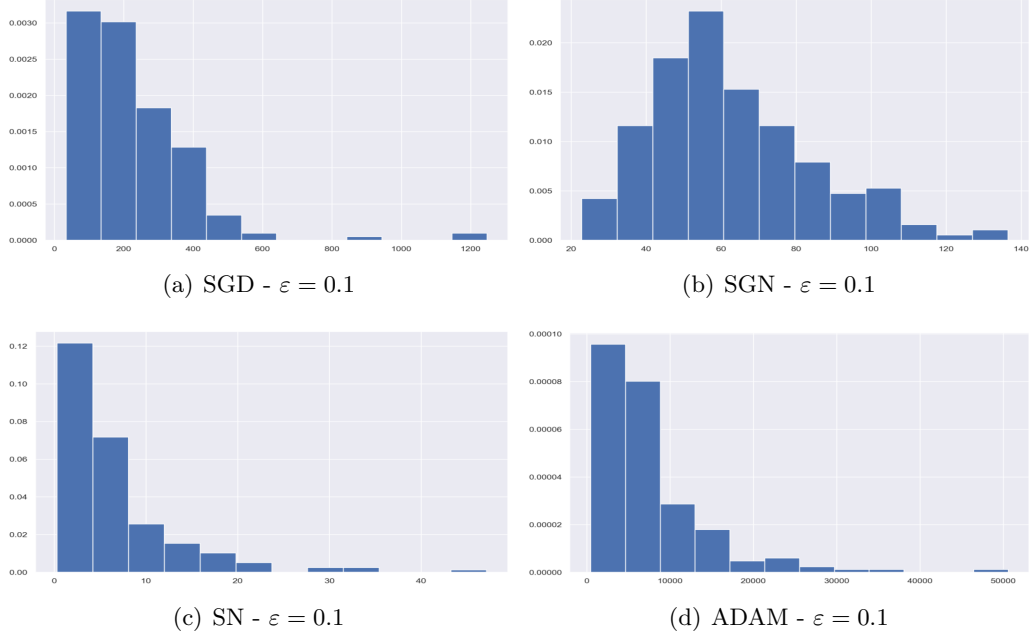
**Figure 7:** Discrete setting with  $I = 10^3$ ,  $J = 50$  and  $\varepsilon = 0.01$ . Histogram of 200 independent realizations of  $\frac{\sqrt{n}(\widehat{W}_n - W_\varepsilon(\mu, \nu))}{\widehat{\sigma}_n}$  with  $n = 4 \times 10^5$  iterations using each of the four stochastic algorithms. The orange curve is the density of the standard Gaussian distribution.

Therefore, these numerical experiments clearly show that using the SGN algorithm has interesting benefits as it outperforms SGD and ADAM for the estimation of  $v^*$  since it yields an estimator  $\widetilde{V}_n$  with a smaller variance.

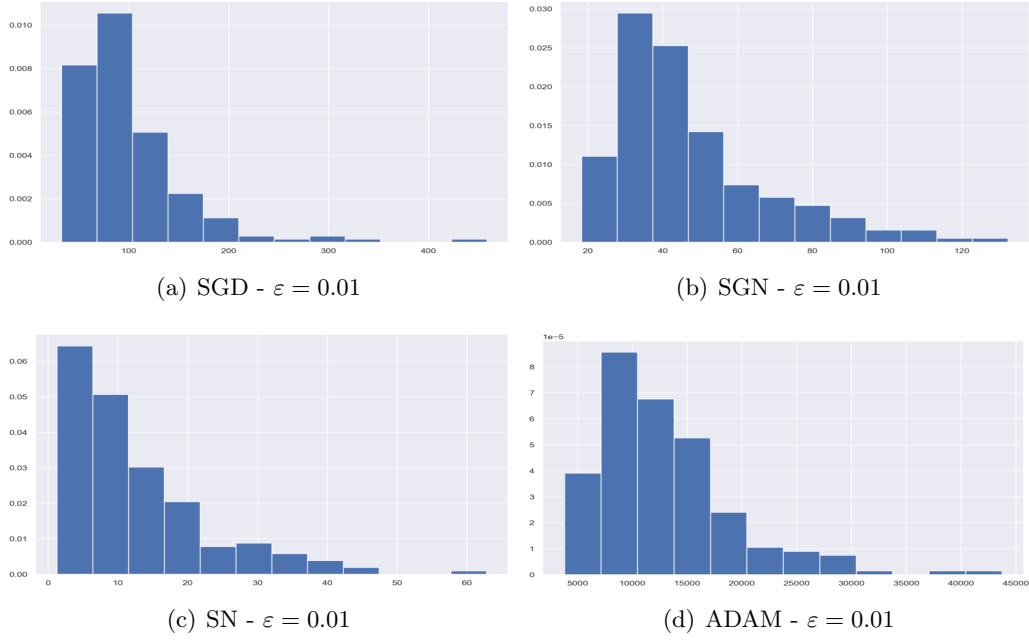
**Convergence of the pre-conditioning matrix for the SN algorithm.** Finally, we report in Figure 10 numerical results (with  $I = 10^4$  and  $J \in \{100, 200\}$ ) on the convergence of  $\overline{S}_n$  to  $G_\varepsilon(v^*)$  as a function the computational time of the SGN algorithm (using  $n = 10^5$  iterations) for different values of the regularization parameter  $\varepsilon$ . We observe that the convergence becomes slower as  $\varepsilon$  decreases.

#### 4.3.2 Semi-discrete setting

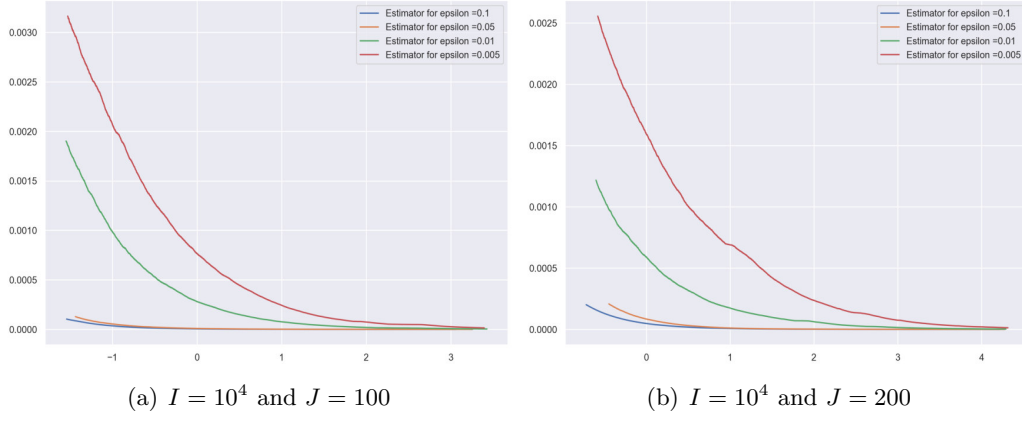
In this section, the cost function is chosen as the following normalized squared Euclidean distance  $c(x, y) = \frac{1}{d}\|x - y\|^2$ . We now consider the framework where  $\mu$  is a mixture of three Gaussian densities in dimension  $d$ . In these numerical experiments, the size  $J = 100$  of the support of  $\nu$  is held fixed, and it is chosen as the uniform discrete probability measure supported on  $J$  points drawn uniformly on the hypercube  $[0, 1]^d$ . The value of the dimension  $d$  is let growing, and we analyze its influence on the performances of the stochastic algorithms with either  $n = 5 \times 10^5$  or  $n = 10^6$  iterations. We also study the performances of the Sinkhorn algorithm from a full-batch sample, that is using the



**Figure 8:** Discrete setting with  $I = 10^3$ ,  $J = 50$  and  $\varepsilon = 0.1$ . Histogram of 200 independent realizations of  $n\|\widehat{V}_n - v^*\|^2$  with  $n = 2 \times 10^5$  iterations using each of the four stochastic algorithms.



**Figure 9:** Discrete setting with  $I = 10^3$ ,  $J = 50$  and  $\varepsilon = 0.01$ . Histogram of 200 independent realizations of  $n\|\widehat{V}_n - v^*\|^2$  with  $n = 4 \times 10^5$  iterations using each of the four stochastic algorithms.



**Figure 10:** Discrete setting and convergence of  $\log \left( \|\bar{S}_n - G_\varepsilon(v^*)\|_F^2 \right)$  as a function the computational time of the SGN algorithm (using  $n = 10^5$  iterations) for different values of the regularization parameter  $\varepsilon$ .

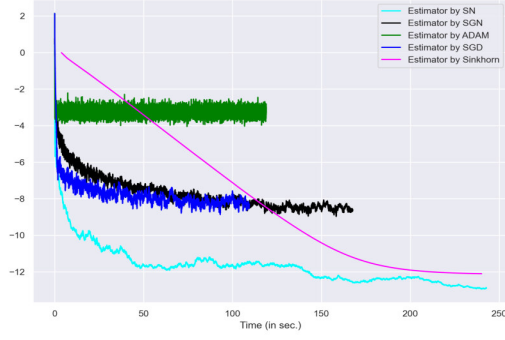
empirical measure

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}, \quad \text{where } X_1, \dots, X_n \sim_{iid} \mu,$$

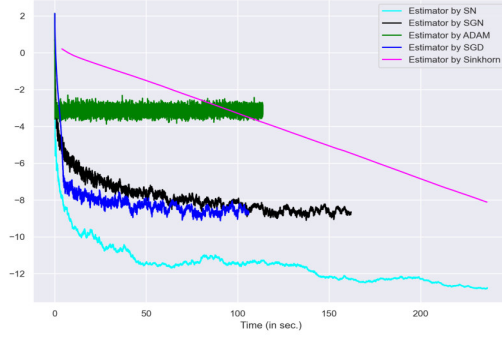
to compute a solution of the regularized OT problem  $W_\varepsilon(\hat{\mu}_n, \nu)$  as an approximation of  $W_\varepsilon(\mu, \nu)$ . At each iteration, the cost of this Sinkhorn algorithm is thus  $\mathcal{O}(nJ)$  as  $n$  is the size of the support of  $\hat{\mu}_n$ . We have chosen to display results for only one simulation as averaging over Monte-Carlo replications does not change our main conclusions from these numerical experiments. All stochastic algorithms (resp. Sinkhorn algorithm) are compared for the metric  $\|\hat{V}_n - v^*\|^2$  (resp.  $\|V_k - v^*\|^2$ ), where the vector  $v^*$  is preliminary approximated by running the SN algorithm with a large value of iterations  $n_{\max} = 10^6$ .

In Figure 11 (for  $n = 5 \times 10^5$  iterations) and Figure 12 (for  $n = n_{\max} = 10^6$  iterations), we display these metrics (in logarithmic scale) as functions of the computational time of the stochastic and the Sinkhorn algorithms for different values of dimension  $d \in \{5, 10, 50\}$  and  $\varepsilon \in \{0.01, 0.005\}$ . First, for either  $n = 5 \times 10^5$  or  $n = 10^6$ , it can be observed that the SN algorithm has always the best performances, while ADAM has the worst ones. Moreover, the SGN algorithm has slightly better performances than SGD. In Figure 12, the very fast decay of the error after 450 seconds for SN is due to the fact that the ground truth value  $v^*$  has been preliminary computed with the SN algorithm with  $n = n_{\max} = 10^6$  iterations. For either  $n = 5 \times 10^5$  or  $n = 10^6$ , we also remark that the SGN outperforms Sinkhorn for  $\varepsilon = 0.005$  and  $d \in \{5, 10\}$  in the sense that when the SGN stops the value reached by  $\|\hat{V}_n - v^*\|^2$  is smaller than  $\|V_k - v^*\|^2$  obtained with Sinkhorn. In larger dimension  $d = 50$ , the Sinkhorn algorithm appears to have better performances than SGN and SGD, but we recall that Sinkhorn uses the full sample at each iteration.

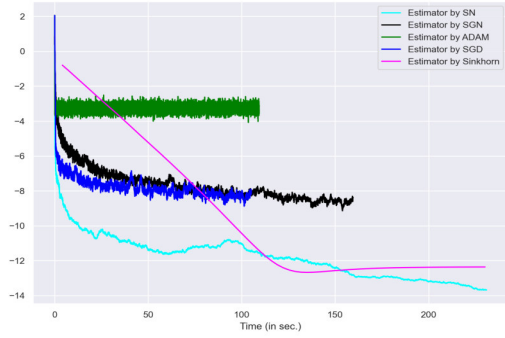
Therefore, these numerical experiments show that SGN has interesting benefits over Sinkhorn in small dimension  $d$  when combined with small values of the regularization parameter  $\varepsilon$  and large values of  $n$ . Moreover, we observe that the best results are always



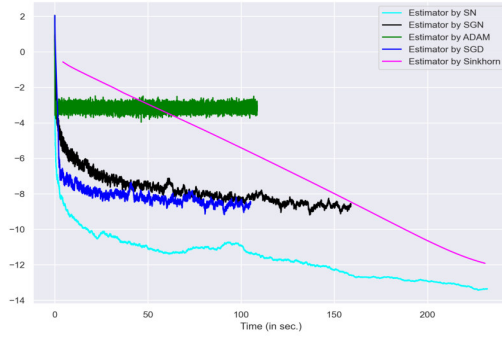
(a)  $d = 5$  and  $\varepsilon = 0.01$



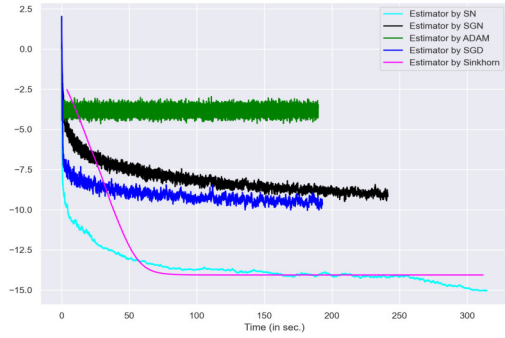
(b)  $d = 5$  and  $\varepsilon = 0.005$



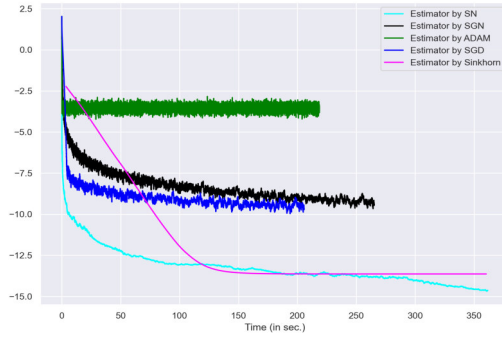
(c)  $d = 10$  and  $\varepsilon = 0.01$



(d)  $d = 10$  and  $\varepsilon = 0.005$



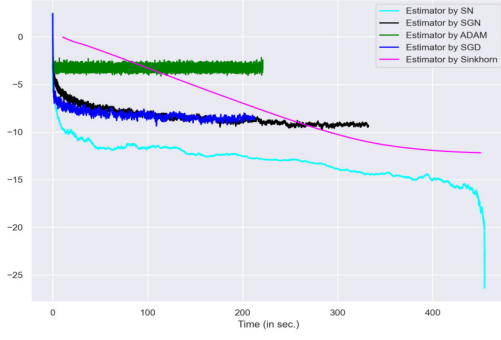
(e)  $d = 50$  and  $\varepsilon = 0.01$



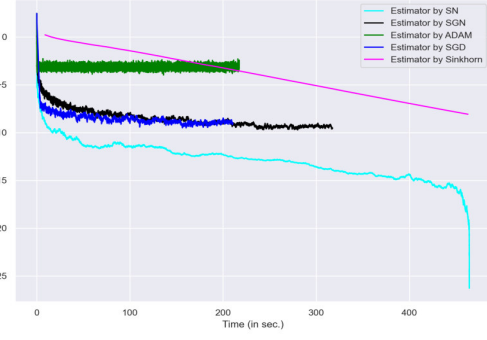
(f)  $d = 50$  and  $\varepsilon = 0.005$

**Figure 11:** Semi-discrete setting where  $\mu$  is a mixture of three Gaussian densities,  $J = 100$  and  $n = 5 \times 10^5$  iterations. Excess risk (in logarithmic scale)  $\log(\|\hat{V}_n - v^*\|^2)$  (resp. metric  $\log(\|V_k - v^*\|^2)$ ) as a function of the computational cost of the iterations of the four stochastic algorithms (resp. the Sinkhorn algorithm) for different values of the dimension  $d$  and the regularization parameter  $\varepsilon$ .

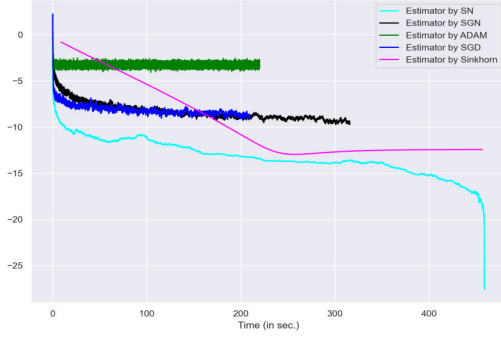




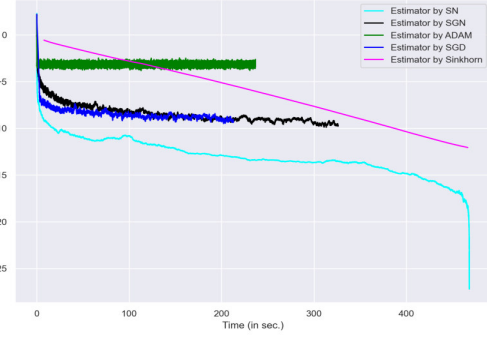
(a)  $d = 5$  and  $\varepsilon = 0.01$



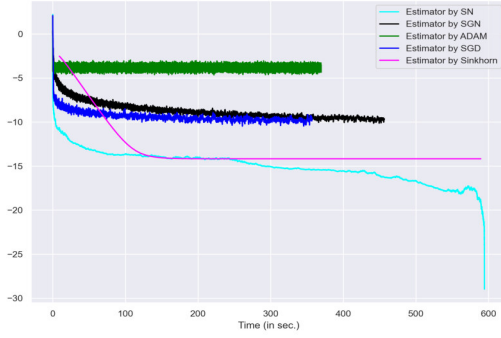
(b)  $d = 5$  and  $\varepsilon = 0.005$



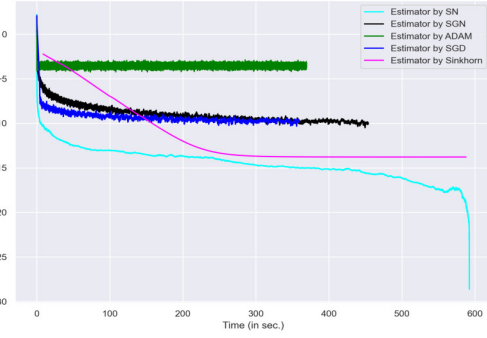
(c)  $d = 10$  and  $\varepsilon = 0.01$



(d)  $d = 10$  and  $\varepsilon = 0.005$



(e)  $d = 50$  and  $\varepsilon = 0.01$



(f)  $d = 50$  and  $\varepsilon = 0.005$

**Figure 12:** Semi-discrete setting where  $\mu$  is a mixture of three Gaussian densities,  $J = 100$  and  $n = n_{\max} = 10^6$  iterations. Excess risk (in logarithmic scale)  $\log(\|\hat{V}_n - v^*\|^2)$  (resp. metric  $\log(\|V_k - v^*\|^2)$ ) as a function of the computational cost of the iterations of the four stochastic algorithms (resp. the Sinkhorn algorithm) for different values of the dimension  $d$  and the regularization parameter  $\varepsilon$ .

obtained with the SN algorithm.

## 5 Properties of the objective function $H_\varepsilon$

The purpose of this section is to discuss various keystone properties of the objective function  $H_\varepsilon$  that are needed to establish our main results.

### 5.1 Gradient properties

Let us first remark that, for any  $x \in \mathcal{X}$ , the function  $v \mapsto h_\varepsilon(x, v)$ , defined by (2.5), is twice differentiable. For a fixed  $x \in \mathcal{X}$ , the gradient vector and Hessian matrix of the function  $h_\varepsilon$ , with respect to its second argument, are given by

$$\nabla_v h_\varepsilon(x, v) = \pi(x, v) - \nu, \quad (5.1)$$

and

$$\nabla_v^2 h_\varepsilon(x, v) = \frac{1}{\varepsilon} \left( \text{diag}(\pi(x, v)) - \pi(x, v)\pi(x, v)^T \right), \quad (5.2)$$

where the  $j^{\text{th}}$  component of the vector  $\pi(x, v) \in \mathbb{R}^J$  is such that

$$\pi_j(x, v) = \left( \sum_{k=1}^J \nu_k \exp\left(\frac{v_k - c(x, y_k)}{\varepsilon}\right) \right)^{-1} \nu_j \exp\left(\frac{v_j - c(x, y_j)}{\varepsilon}\right).$$

Consequently, the gradient vector and the Hessian matrix of the function  $H_\varepsilon$ , defined by (2.3), are as follows

$$\nabla H_\varepsilon(v) = \mathbb{E}[\nabla_v h_\varepsilon(X, v)] = \mathbb{E}[\pi(X, v)] - \nu, \quad (5.3)$$

and

$$\nabla^2 H_\varepsilon(v) = \mathbb{E}[\nabla_v^2 h_\varepsilon(X, v)] = \frac{1}{\varepsilon} \mathbb{E}[\text{diag}(\pi(X, v)) - \pi(X, v)\pi(X, v)^T]. \quad (5.4)$$

Note that the minimizer  $v^*$  satisfies  $\nabla H_\varepsilon(v^*) = 0$ , leading to

$$\mathbb{E}[\pi(X, v^*)] = \nu,$$

which allows us to simplify the expression for the Hessian of  $H_\varepsilon$  at  $v^*$ ,

$$\nabla^2 H_\varepsilon(v^*) = \frac{1}{\varepsilon} \left( \text{diag}(\nu) - \mathbb{E}[\pi(X, v^*)\pi(X, v^*)^T] \right). \quad (5.5)$$

We now discuss some properties of the above gradient vectors and Hessian matrices that will be of interest to study the SGN algorithm.

### 5.2 Convexity of $H_\varepsilon$ and related properties

First of all, the baseline remark is that  $\nabla_v^2 H_\varepsilon(v)$  is a positive semi-definite matrix for any  $v \in \mathbb{R}^J$ , which entails the convexity of  $H_\varepsilon$ .

**Minimizers and rank of the Hessian.** It is clear from (5.4) that for any  $v \in \mathbb{R}^J$ , the smallest eigenvalue of the Hessian matrix  $\nabla_v^2 H_\varepsilon(v)$  associated to the eigenvector  $\mathbf{v}_J$  is equal to zero. Therefore, as indicated in the end of Section 2.1, for any  $t \in \mathbb{R}$ , the vector  $v^* + t\mathbf{v}_J$  is also a minimizer of (2.4). Nevertheless, it is well-known [15] that the minimizer  $v^*$  of (2.4) is unique up to a scalar translation of its coordinates. We shall thus denote by  $v^*$  the minimizer of (2.3) satisfying  $\langle v^*, \mathbf{v}_J \rangle = 0$ . It means that  $v^*$  belongs to  $\langle \mathbf{v}_J \rangle^\perp$ , and that the function  $H_\varepsilon$  admits a unique minimizer over the  $J-1$  dimensional subspace  $\langle \mathbf{v}_J \rangle^\perp$ . However, as already shown in [23] and further discussed in [5, Section 3.3], the objective function  $H_\varepsilon$  is *not strongly convex*, even by restricting the maximization problem (2.3) to the subspace  $\langle \mathbf{v}_J \rangle^\perp$  since it may be shown that  $v \mapsto H_\varepsilon(v)$  may have some vanishing curvature, leading to a flat landscape, *i.e.* to eigenvalues of the Hessian matrix that are arbitrarily close to 0 for large values of  $\|v\|$  in  $\langle \mathbf{v}_J \rangle^\perp$ .

Moreover, for any  $(x, v) \in \mathcal{X} \times \mathbb{R}^J$ , it follows from [5, Lemma A.1] that the matrices  $\nabla_v^2 h_\varepsilon(x, v)$  and  $\nabla_v^2 H_\varepsilon(v)$  are of rank  $J-1$ , and therefore, all their eigenvectors associated to non-zero eigenvalues belong to  $\langle \mathbf{v}_J \rangle^\perp$ . Finally, one also has that  $\nabla_v h_\varepsilon(x, v) \in \langle \mathbf{v}_J \rangle^\perp$  for any  $(x, v) \in \mathcal{X} \times \mathbb{R}^J$ .

**Useful upper and lower bounds.** We conclude this section by stating a few inequalities that we repeatedly use in the proofs of our main results. Since  $\nu$  and  $\pi(x, v)$  are vectors with positive entries that sum up to one, it follows from (5.1) that for any  $(x, v) \in \mathcal{X} \times \mathbb{R}^J$ ,

$$\|\nabla_v h_\varepsilon(x, v)\| \leq \|\nu\| + \|\pi(x, v)\| \leq 2, \quad (5.6)$$

and that the gradient of  $H_\varepsilon$  is always bounded for any  $v \in \mathbb{R}^J$ ,

$$\|\nabla H_\varepsilon(v)\| \leq 2. \quad (5.7)$$

Moreover, thanks to the property that

$$\lambda_{\max} \left( \frac{1}{\varepsilon} (\text{diag}(\pi(x, v)) - \pi(x, v)\pi(x, v)^T) \right) \leq \frac{1}{\varepsilon} \text{Tr}(\text{diag}(\pi(x, v))) = \frac{1}{\varepsilon},$$

we obtain that for any  $(x, v) \in \mathcal{X} \times \mathbb{R}^J$ ,

$$\lambda_{\max}(\nabla_v^2 h_\varepsilon(x, v)) \leq \frac{1}{\varepsilon} \quad \text{and} \quad \lambda_{\max}(\nabla_v^2 H_\varepsilon(v)) \leq \frac{1}{\varepsilon}. \quad (5.8)$$

Finally, by [5, Lemma A.1], the second smallest eigenvalue of  $\nabla^2 H_\varepsilon(v^*)$  is positive, and one has that

$$\lambda_{\min}(\nabla^2 H_\varepsilon(v^*)) = \min_{v \in \langle \mathbf{v}_J \rangle^\perp} \left\{ \frac{v^T \nabla^2 H_\varepsilon(v^*) v}{\|v\|^2} \right\} \geq \frac{1}{\varepsilon} \min(\nu). \quad (5.9)$$

### 5.3 Generalized self-concordance for regularized semi-discrete OT

Let us now introduce the so-called notion of generalized self-concordance proposed in Bach [2] for the purpose of obtaining fast rates of convergence for stochastic algorithms with non-strongly convex objective functions. Generalized self-concordance has been shown to hold for regularized semi-discrete OT in [5], and we discuss below its implications of some

key properties for the analysis of the SGN algorithm studied in this paper. To this end, for any  $v \in \langle \mathbf{v}_J \rangle^\perp$  and for all  $t$  in the interval  $[0, 1]$ , we denote  $v_t = v^* + t(v - v^*)$ , and we define the function  $\varphi$ , for all  $t \in [0, 1]$ , as

$$\varphi(t) = H_\varepsilon(v_t).$$

The second-order Taylor expansion of  $\varphi$  with integral remainder is given by

$$\varphi(1) = \varphi(0) + \varphi'(0) - \int_0^1 (t-1)\varphi''(t) dt. \quad (5.10)$$

Using that  $\varphi(1) = H_\varepsilon(v)$ ,  $\varphi(0) = H_\varepsilon(v^*)$  and  $\varphi'(0) = \langle v - v^*, \nabla H_\varepsilon(v^*) \rangle = 0$ , it has been first remarked in [5] that inequality (5.8) implies that

$$H_\varepsilon(v) - H_\varepsilon(v^*) \leq \frac{1}{2\varepsilon} \|v - v^*\|^2 \quad (5.11)$$

Moreover, it is shown in the proof of [5, Lemma A.2] that the following inequality holds

$$|\varphi'''(t)| \leq \frac{\sqrt{2}}{\varepsilon} \varphi''(t) \|v - v^*\|. \quad (5.12)$$

It means that the function  $\varphi$  satisfies the so-called generalized self-concordance property with constant  $s_\varepsilon = \sqrt{2}/\varepsilon$  as defined in Appendix B of [2]. As a consequence of inequality (5.12) and thanks to the arguments in the proof of [5, Lemma A.2], the error of linearizing the gradient  $\nabla H_\varepsilon(v) \approx \nabla^2 H_\varepsilon(v^*)(v - v^*)$  is controlled as follows,

$$\|\nabla H_\varepsilon(v) - \nabla^2 H_\varepsilon(v^*)(v - v^*)\| \leq 2s_\varepsilon \|v - v^*\|^2. \quad (5.13)$$

Moreover, generalized self-concordance also implies the following result (which is a consequence of the arguments in the proof of Lemma A.2 in [5]) that may be interpreted as a local strong convexity property of the function  $H_\varepsilon$  in the neighborhood of  $v^*$ .

**Lemma 5.1.** *For any  $v \in \langle \mathbf{v}_J \rangle^\perp$ , we have*

$$\langle \nabla H_\varepsilon(v), v - v^* \rangle \geq \frac{1 - \exp(-\delta(v))}{\delta(v)} (v - v^*)^T \nabla^2 H_\varepsilon(v^*)(v - v^*), \quad (5.14)$$

where  $\delta(v) = s_\varepsilon \|v - v^*\|$ .

Finally, if we now consider the matrix-valued function  $G_\varepsilon(v)$  introduced in equation (3.1), we have the following result which can be interpreted as a local Lipschitz property of  $G_\varepsilon(v)$  around  $v = v^*$ . The proof of this lemma is postponed to Appendix A.

**Lemma 5.2.** *For any  $v \in \langle \mathbf{v}_J \rangle^\perp$ , we have that*

$$-\frac{4}{\varepsilon} \|v - v^*\| I_J \leq G_\varepsilon(v) - G_\varepsilon(v^*) \leq \frac{4}{\varepsilon} \|v - v^*\| I_J, \quad (5.15)$$

in the sense of partial ordering between positive semi-definite matrices.

## 6 Proofs of the main results

This section contains the proofs of our main results that are stated in Section 3. Our results are based on previous important contributions on self-concordance functions (see *e.g.* [2, 5]), regularization of second order algorithms [6] and on the Kurdyka-Łojasiewicz inequality adapted to stochastic algorithms [22]. More specifically, almost sure convergence and almost sure convergence rates crucially depend on the adaptive property (3.4), which induces a contraction rate of the sequence  $G_\varepsilon^{1/2}(v^*)(\widehat{V}_n - v^*)$  (see Equation (6.11) below). The non-asymptotic study is then based on both the self-concordance property, stated in Lemma 5.1 and on the KL inequality stated in Proposition 6.1. The combination of these two properties is an essential novelty brought by our work in order to build a key Lyapunov function in Equation (6.45). We emphasize that to obtain the results stated below, we have derived quantitative computations that are specific to the regularized OT problem. In particular, if the use of the KL inequality is borrowed from [22], the exact values of the constant  $m_\varepsilon$  and  $M$  used in Proposition 6.1 crucially depend on the self-concordance property of the regularized OT problem.

### 6.1 Keystone property

We start with the proof of inequality (3.4) that states the adaptivity of the SGN algorithm to the local geometry of  $H_\varepsilon$ .

*Proof of Proposition 3.1.* First of all, one can remark that  $G_\varepsilon(v^*)$  is a positive semi-definite matrix whose smallest eigenvalue is equal to zero and associated to the eigenvector  $\mathbf{v}_J$ . Thus, all the eigenvectors of  $G_\varepsilon(v^*)$  associated to non-zero eigenvalues belong to  $\langle \mathbf{v}_J \rangle^\perp$ . We already saw from (5.5) that

$$\nabla^2 H_\varepsilon(v^*) = \frac{1}{\varepsilon} (\text{diag}(\nu) - \mathbb{E}[\pi(X, v^*)\pi(X, v^*)^T]).$$

Consequently, it follows from (3.1) and (5.1) that

$$G_\varepsilon(v^*) = \mathbb{E}[\pi(X, v^*)\pi(X, v^*)^T] - \nu\nu^T = \text{diag}(\nu) - \nu\nu^T - \varepsilon\nabla^2 H_\varepsilon(v^*),$$

which implies that

$$G_\varepsilon(v^*) = \nabla^2 H_\varepsilon(v^*) + \Sigma_\varepsilon^*$$

where

$$\Sigma_\varepsilon^* = \text{diag}(\nu) - \nu\nu^T - (1 + \varepsilon)\nabla^2 H_\varepsilon(v^*).$$

On the one hand, it is easy to see that  $\mathbf{v}_J^T \Sigma_\varepsilon^* \mathbf{v}_J = 0$ . On the other hand, we deduce from inequality (5.9) that for all  $v \in \langle \mathbf{v}_J \rangle^\perp$ ,

$$v^T \Sigma_\varepsilon^* v \leq \max(\nu) - \left( \frac{1 + \varepsilon}{\varepsilon} \right) \min(\nu).$$

Finally, condition (3.2) on the regularization parameter  $\varepsilon$  leads to  $G_\varepsilon(v^*) \leq \nabla^2 H_\varepsilon(v^*)$ , which completes the proof of Proposition 3.1.  $\square$

## 6.2 Proofs of the most sure convergence results

*Proof of Theorem 3.1.* In what follows, we borrow some arguments from the proof of [6, Theorem 4.1] to establish the almost sure convergence of the regularized versions of the SGN algorithm as an application of the Robbins-Siegmund Theorem [41].

• We already saw that for all  $n \geq 0$ ,  $\widehat{V}_n$  belongs to  $\langle v_J \rangle^\perp$ . We clearly have from (5.1) that for all  $n \geq 0$ ,  $\nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n)$  also belong to  $\langle v_J \rangle^\perp$ . Hence, we have from (2.6) that for all  $n \geq 0$ ,

$$\widehat{V}_{n+1} = \widehat{V}_n - n^\alpha P_J S_n^{-1} P_J (\nabla H_\varepsilon(\widehat{V}_n) + \varepsilon_{n+1}), \quad (6.1)$$

where the martingale increment  $\varepsilon_{n+1}$  is given by

$$\varepsilon_{n+1} = \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n) - \mathbb{E}[\nabla(h_\varepsilon(X_{n+1}, \widehat{V}_n)) | \mathcal{F}_n] = \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n) - \nabla H_\varepsilon(\widehat{V}_n)$$

with  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$ . Moreover, it follows from the Taylor-Lagrange formula that

$$H_\varepsilon(\widehat{V}_{n+1}) = H_\varepsilon(\widehat{V}_n) + \nabla H_\varepsilon(\widehat{V}_n)^T (\widehat{V}_{n+1} - \widehat{V}_n) + \frac{1}{2} (\widehat{V}_{n+1} - \widehat{V}_n)^T \nabla^2 H_\varepsilon(\xi_{n+1}) (\widehat{V}_{n+1} - \widehat{V}_n), \quad (6.2)$$

where  $\xi_{n+1} = \widehat{V}_n + t(\widehat{V}_{n+1} - \widehat{V}_n)$  with  $t \in ]0, 1[$ . Consequently, we deduce from (6.1) and (6.2) that for all  $n \geq 0$ ,

$$\begin{aligned} H_\varepsilon(\widehat{V}_{n+1}) &= H_\varepsilon(\widehat{V}_n) - n^\alpha \langle \nabla H_\varepsilon(\widehat{V}_n), P_J S_n^{-1} P_J (\nabla H_\varepsilon(\widehat{V}_n) + \varepsilon_{n+1}) \rangle \\ &\quad + \frac{n^{2\alpha}}{2} (P_J S_n^{-1} P_J (\nabla H_\varepsilon(\widehat{V}_n) + \varepsilon_{n+1}))^T \nabla^2 H_\varepsilon(\xi_{n+1}) (P_J S_n^{-1} P_J (\nabla H_\varepsilon(\widehat{V}_n) + \varepsilon_{n+1})). \end{aligned}$$

Taking the conditional expectation with respect to  $\mathcal{F}_n$  on both sides of the previous equality, we obtain that for all  $n \geq 0$ ,

$$\begin{aligned} \mathbb{E}[H_\varepsilon(\widehat{V}_{n+1}) | \mathcal{F}_n] &= H_\varepsilon(\widehat{V}_n) - n^\alpha \nabla H_\varepsilon(\widehat{V}_n)^T P_J S_n^{-1} P_J \nabla H_\varepsilon(\widehat{V}_n) \\ &\quad + \frac{n^{2\alpha}}{2} \mathbb{E} \left[ \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n)^T P_J S_n^{-1} P_J \nabla^2 H_\varepsilon(\xi_{n+1}) P_J S_n^{-1} P_J \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n) \middle| \mathcal{F}_n \right] \\ &= H_\varepsilon(\widehat{V}_n) - n^\alpha \nabla H_\varepsilon(\widehat{V}_n)^T P_J S_n^{-1} P_J \nabla H_\varepsilon(\widehat{V}_n) \\ &\quad + \frac{n^{2\alpha}}{2} \mathbb{E} \left[ \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n)^T S_n^{-1} \nabla^2 H_\varepsilon(\xi_{n+1}) S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n) \middle| \mathcal{F}_n \right] \end{aligned} \quad (6.3)$$

using the elementary fact that  $P_J \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n) = \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n)$  as well as  $\nabla^2 H_\varepsilon(\xi_{n+1}) v_J = 0$  which implies that  $P_J \nabla^2 H_\varepsilon(\xi_{n+1}) P_J = \nabla^2 H_\varepsilon(\xi_{n+1})$ . On the one hand, we have from inequality (5.6) that  $\|\nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n)\| \leq 2$ . On the other hand, it follows from inequality (5.8) that  $\lambda_{\max}(\nabla^2 H_\varepsilon(\xi_{n+1})) \leq 1/\varepsilon$ . Therefore, we deduce from (6.3) that for all  $n \geq 0$ ,

$$\mathbb{E}[H_\varepsilon(\widehat{V}_{n+1}) | \mathcal{F}_n] \leq H_\varepsilon(\widehat{V}_n) + A_n - B_n \quad \text{a.s.}$$

where the two positive random variables  $A_n$  and  $B_n$  are given by

$$A_n = \frac{2n^{2\alpha}}{\varepsilon(\lambda_{\min}(S_n))^2}$$

and  $B_n = n^\alpha \nabla H_\varepsilon(\widehat{V}_n)^T P_J S_n^{-1} P_J \nabla H_\varepsilon(\widehat{V}_n)$ . Our purpose is now to show that

$$\sum_{n=1}^{\infty} A_n < +\infty \quad \text{a.s.}$$

We already saw from (4.3) that for all  $n \geq 1$ ,

$$S_n = I_J + \sum_{k=1}^n \nabla_v h_\varepsilon(X_k, \widehat{V}_{k-1}) \nabla_v h_\varepsilon(X_k, \widehat{V}_{k-1})^T + R_n$$

where

$$R_n = \sum_{k=1}^n \gamma \left(1 + \left\lfloor \frac{k}{J} \right\rfloor\right)^{-\beta} Z_k Z_k^T.$$

We clearly have  $\lambda_{\min}(S_n) \geq \lambda_{\min}(R_n)$ . Let  $p_n$  be the largest integer such that  $p_n J \leq n$ . One can remark that

$$R_n = \left( \sum_{m=1}^{p_n} m^{-\beta} \right) \gamma \text{diag}(\nu) + \sum_{k=p_n J + 1}^n \gamma \left(1 + \left\lfloor \frac{k}{J} \right\rfloor\right)^{-\beta} Z_k Z_k^T, \quad (6.4)$$

which implies that

$$\lambda_{\min}(S_n) \geq \gamma \min(\nu) \left( \sum_{m=1}^{p_n} m^{-\beta} \right).$$

However, for any  $0 < \beta < 1/2$  and for all  $p_n \geq 4$ ,

$$\sum_{m=1}^{p_n} \frac{1}{m^\beta} \geq \frac{p_n^{1-\beta}}{2(1-\beta)}.$$

Consequently, using that  $p_n \geq nJ^{-1} - 1$ , we obtain that

$$\sum_{n=1}^{\infty} A_n \leq \frac{8(1-\beta)^2}{\varepsilon(\gamma \min(\nu))^2} \sum_{n=1}^{\infty} \frac{n^{2\alpha}}{p_n^{2(1-\beta)}} \leq \frac{16(1-\beta)^2 J^{2(1-\beta)}}{\varepsilon(\gamma \min(\nu))^2} \sum_{n=1}^{\infty} \frac{1}{n^{2(1-\alpha-\beta)}} < +\infty \quad \text{a.s.}$$

since the assumption  $0 < \alpha + \beta < 1/2$  implies that  $2(1 - \alpha - \beta) > 1$ . Therefore, we can apply the Robbins-Siegmund Theorem [41] to conclude that the sequence  $(H_\varepsilon(\widehat{V}_n))$  converges almost surely to a finite random variable and that the series

$$\sum_{n=1}^{\infty} B_n < +\infty \quad \text{a.s.}$$

leading to

$$\sum_{n=1}^{\infty} n^\alpha \frac{\|\nabla H_\varepsilon(\widehat{V}_n)\|^2}{\lambda_{\max}(S_n)} < +\infty \quad \text{a.s.} \quad (6.5)$$

One can very from inequality (5.6) and (4.3) that for all  $n \geq 1$ ,

$$\lambda_{\max}(S_n) \leq 1 + 4n + \gamma \max(\nu) \sum_{k=1}^n \left(1 + \left\lfloor \frac{k}{J} \right\rfloor\right)^{-\beta} \leq 1 + (4 + \gamma \max(\nu))n.$$



Since  $\alpha \geq 0$ , it implies that

$$\sum_{n=1}^{\infty} \frac{n^\alpha}{\lambda_{\max}(S_n)} = +\infty \quad \text{a.s.} \quad (6.6)$$

The rest of the proof proceeds from standard arguments combining (6.5) and (6.6). Let

$$H_{\varepsilon, \infty} = \lim_{n \rightarrow +\infty} H_\varepsilon(\widehat{V}_n) \quad \text{a.s.}$$

and assume by contradiction that  $H_{\varepsilon, \infty} > H_\varepsilon(v^*)$  where

$$H_\varepsilon(v^*) = \min_{v \in \langle \mathbf{v}_J \rangle^\perp} H_\varepsilon(v).$$

Since  $H_\varepsilon$  is a convex function with a unique minimizer  $v^*$  on  $\langle \mathbf{v}_J \rangle^\perp$ , we necessarily have

$$\lim_{\|v\| \rightarrow +\infty} H_\varepsilon(v) = +\infty.$$

It means that  $(\widehat{V}_n)$  is almost surely bounded since  $H_{\varepsilon, \infty}$  is finite. Therefore, we can find a compact set  $K$  such that  $v^* = \arg \min_{v \in \langle \mathbf{v}_J \rangle^\perp} H_\varepsilon(v) \notin K$  and  $\widehat{V}_n \in K$  for all  $n$  large enough. Using the continuity of  $\|\nabla H_\varepsilon\|$  and the compactness of  $K$ , we conclude that  $\|\nabla H_\varepsilon\|$  attains its lower bound, which is strictly positive on  $K$ . It ensures the existence of a constant  $c > 0$ , such that, for all  $n$  large enough,

$$\|\nabla H_\varepsilon(\widehat{V}_n)\| \geq c > 0.$$

The above lower bound associated with (6.5) and (6.6) yields a contradiction. Hence, we can conclude that

$$\lim_{n \rightarrow +\infty} \|\nabla H_\varepsilon(\widehat{V}_n)\| = 0 \quad \text{a.s.}$$

It clearly implies that equation (3.5) holds true since  $(\widehat{V}_n)$  is a bounded sequence with a unique adherence point  $v^*$ .

- It now remains to investigate the almost sure convergence of the matrix  $\overline{S}_n$ . We observe from equation (4.3) that  $S_n$  can be splitted into two terms,

$$S_n = M_n + \Sigma_n \quad (6.7)$$

with

$$M_n = \sum_{k=1}^n \Phi_k \Phi_k^T - G_\varepsilon(\widehat{V}_{k-1}) \quad \text{and} \quad \Sigma_n = I_J + \sum_{k=1}^n G_\varepsilon(\widehat{V}_{k-1}) + R_n,$$

where the vector  $\Phi_k$  stands for  $\Phi_k = \nabla_v h_\varepsilon(X_k, \widehat{V}_{k-1})$ . Using the assumption  $0 < \beta < 1/2$ , we have from (6.4) that

$$\lim_{n \rightarrow +\infty} \frac{1}{n} (I_J + R_n) = 0.$$

Moreover, it follows from (5.3) that  $G_\varepsilon$  is a continuous function from  $\mathbb{R}^J$  to  $\mathbb{R}^{J \times J}$ . Consequently, we deduce from convergence (3.5) together with the Cesaro mean convergence theorem that

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{k=1}^n G_\varepsilon(\widehat{V}_{k-1}) = G_\varepsilon(v^*) \quad \text{a.s.}$$

which implies that

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \Sigma_n = G_\varepsilon(v^*) \quad \text{a.s.} \quad (6.8)$$

Hereafter, we focus our attention on the first term  $M_n$  in the right-hand side of (6.7). For any  $u \in \mathbb{R}^J$ , let

$$M_n(u) = u^T M_n u = \sum_{k=1}^n \xi_k(u)$$

where, for all  $n \geq 1$ ,  $\xi_n(u) = u^T (\Phi_n \Phi_n^T - G_\varepsilon(\widehat{V}_{n-1}))u$ . It follows from (3.1) that for all  $n \geq 1$ ,  $\mathbb{E}[\Phi_{n+1} \Phi_{n+1}^T | \mathcal{F}_n] = G_\varepsilon(\widehat{V}_n)$ . Hence, for all  $n \geq 1$ ,  $\mathbb{E}[\xi_{n+1}(u) | \mathcal{F}_n] = 0$ . Furthermore, we obtain from (5.6) and (3.1) that for all  $n \geq 1$ ,  $\mathbb{E}[\xi_{n+1}^2(u) | \mathcal{F}_n] \leq 16\|u\|^2$ . Consequently,  $(M_n(u))$  is a locally square-integrable martingale with predictable quadratic variation satisfying

$$\langle M(u) \rangle_n = \sum_{k=1}^n \mathbb{E}[\xi_k^2(u) | \mathcal{F}_{k-1}] \leq 16n\|u\|^4.$$

We deduce from the strong law of large numbers for martingales given (e.g. by Theorem 1.3.24 in [18]) that

$$\lim_{n \rightarrow +\infty} \frac{1}{n} M_n(u) = 0 \quad \text{a.s.}$$

which may be translated immediately into the matricial form

$$\lim_{n \rightarrow +\infty} \frac{1}{n} M_n = 0 \quad \text{a.s.} \quad (6.9)$$

Finally, the convergence (3.6) follows from the decomposition (6.7) together with (6.8) and (6.9), which completes the proof of Theorem 3.1.  $\square$

It is straightforward to obtain the almost sure convergence of  $\widehat{W}_n$  as follows.

*Proof of Corollary 3.1.* By Theorem 3.1 one has that  $\widehat{V}_n$  converges a.s. to  $v^*$  under the assumption that  $\alpha + \beta < 1/2$ . Then, the almost sure convergence of  $\widehat{W}_n$  to  $W_\varepsilon(\mu, \nu)$  follows from assumption (3.7) and the arguments in the proof of [5, Theorem 3.5].  $\square$

### 6.3 Proofs of the almost sure rates of convergence

We now establish the almost sure rates of convergence rates for the SGN algorithm. In contrast with the previous results, we emphasize that the regularization parameter  $\varepsilon$  must now be small enough, in the sense of condition (3.2). This entails the key inequality (3.4) deduced from Proposition 3.1.

*Proof of Theorem 3.2.* To alleviate the notation, we denote by  $\|A\|$  either the operator norm  $\|A\|_2$  or the Frobenius norm  $\|A\|_F$  all along the proof. Since these two norms are equivalent and verify that  $\|\cdot\|_2^2 \leq \|\cdot\|_F^2 \leq J\|\cdot\|_2^2$ , the upper bounds derived below might hold up to multiplicative constant depending on  $J$ , which will not affect the results that are purely asymptotic.

• Our starting point when  $\alpha = 0$  is equation (2.6) written with  $S_n = n\bar{S}_n$ . We recall that the martingale increment is  $\varepsilon_{n+1} = \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) - \nabla H_\varepsilon(\hat{V}_n)$ , so that for all  $n \geq 0$ ,

$$\begin{aligned}\hat{V}_{n+1} - v^* &= \hat{V}_n - v^* - P_J S_n^{-1} P_J (\nabla H_\varepsilon(\hat{V}_n) + \varepsilon_{n+1}) \\ &= \hat{V}_n - v^* - \frac{1}{n} (P_J (\bar{S}_n^{-1} - G_\varepsilon^-(v^*)) P_J) (\nabla H_\varepsilon(\hat{V}_n) + \varepsilon_{n+1}) \\ &\quad - \frac{1}{n} G_\varepsilon^-(v^*) (\nabla H_\varepsilon(\hat{V}_n) + \varepsilon_{n+1})\end{aligned}\tag{6.10}$$

where we decomposed  $S_n^{-1} = n^{-1} \bar{S}_n^{-1} = n^{-1} G_\varepsilon^-(v^*) + n^{-1} (\bar{S}_n^{-1} - G_\varepsilon^-(v^*))$  and  $P_J G_\varepsilon(v^*) = G_\varepsilon(v^*)$  which implies that  $P_J G_\varepsilon^-(v^*) P_J = G_\varepsilon^-(v^*)$ . The rest of the proof then consists in a linearization of  $\nabla H_\varepsilon(\hat{V}_n)$  around  $v^*$ . For that purpose, denote

$$D_n = (P_J (\bar{S}_n^{-1} - G_\varepsilon^-(v^*)) P_J) \quad \text{and} \quad \delta_n = \nabla H_\varepsilon(\hat{V}_n) - \nabla^2 H_\varepsilon(v^*) (\hat{V}_n - v^*).$$

We obtain from (6.10) that for all  $n \geq 0$ ,

$$\begin{aligned}\hat{V}_{n+1} - v^* &= \left( P_J - \frac{1}{n} G_\varepsilon^-(v^*) \nabla^2 H_\varepsilon(v^*) \right) (\hat{V}_n - v^*) - \frac{1}{n} P_J \bar{S}_n^{-1} P_J \varepsilon_{n+1} \\ &\quad - \frac{1}{n} P_J \bar{S}_n^{-1} P_J \delta_n - \frac{1}{n} D_n \nabla^2 H_\varepsilon(v^*) (\hat{V}_n - v^*).\end{aligned}$$

Hence, by setting  $\hat{U}_n = G_\varepsilon^{1/2}(v^*) (\hat{V}_n - v^*)$ , we obtain that for all  $n \geq 0$ ,

$$\hat{U}_{n+1} = \left( P_J - \frac{1}{n} \Gamma_\varepsilon(v^*) \right) \hat{U}_n - \frac{1}{n} A_n \varepsilon_{n+1} - \frac{1}{n} T_n\tag{6.11}$$

where  $\Gamma_\varepsilon(v^*) = G_\varepsilon^{-1/2}(v^*) \nabla^2 H_\varepsilon(v^*) G_\varepsilon^{-1/2}(v^*)$  and  $T_n = A_n \delta_n + B_n (\hat{V}_n - v^*)$  with

$$A_n = G_\varepsilon^{1/2}(v^*) P_J \bar{S}_n^{-1} P_J,\tag{6.12}$$

$$B_n = G_\varepsilon^{1/2}(v^*) D_n \nabla^2 H_\varepsilon(v^*).\tag{6.13}$$

Thanks to inequality (3.4), we have that  $\lambda_{\min}^{\langle v_J \rangle^\perp}(\Gamma_\varepsilon(v^*)) \geq 1$ . For all  $0 \leq k \leq n$ , let

$$P_k^n = \prod_{i=k+1}^n \left( P_J - \frac{1}{i} \Gamma_\varepsilon(v^*) \right)\tag{6.14}$$

with the usual convention that  $P_n^n = P_J$ . We deduce from (6.11) that for all  $n \geq 0$ ,

$$\hat{U}_{n+1} = P_0^n \hat{U}_1 - \sum_{k=1}^n \frac{1}{k} P_k^n A_k \varepsilon_{k+1} - \sum_{k=1}^n \frac{1}{k} P_k^n T_k.\tag{6.15}$$

The first term of (6.15) is easy to handle. If  $\rho$  stands for the minimal eigenvalue of  $\Gamma_\varepsilon(v^*)$  when restricted to act on the subspace  $\langle v_J \rangle^\perp$ , a simple diagonalization of the matrix  $\Gamma_\varepsilon(v^*)$  leads, for all  $0 \leq k \leq n$ , to

$$\|P_k^n\| \leq \kappa \left( \frac{k}{n} \right)^\rho\tag{6.16}$$

where  $\kappa > 0$ . Concerning the middle term of (6.15), let  $(M_n)$  be the multidimensional martingale defined by  $M_1 = 0$  and, for all  $n \geq 1$ ,

$$M_{n+1} = \sum_{k=1}^n A_k \varepsilon_{k+1}.$$

We infer from (5.3) and (3.1) that  $\|\varepsilon_{n+1}\| \leq 4$ ,  $\mathbb{E}[\varepsilon_{n+1}|\mathcal{F}_n] = 0$  and

$$\mathbb{E}[\varepsilon_{n+1} \varepsilon_{n+1}^T | \mathcal{F}_n] = G_\varepsilon(\widehat{V}_n) - \nabla H_\varepsilon(\widehat{V}_n) \nabla H_\varepsilon(\widehat{V}_n)^T.$$

Moreover, it follows from (3.5)

$$\lim_{n \rightarrow +\infty} \nabla H_\varepsilon(\widehat{V}_n) = \nabla H_\varepsilon(v^*) = 0 \quad \text{and} \quad \lim_{n \rightarrow +\infty} G_\varepsilon(\widehat{V}_n) = G_\varepsilon(v^*) \quad \text{a.s.}$$

which ensures via (3.6) and (6.12) that

$$\lim_{n \rightarrow +\infty} A_n \mathbb{E}[\varepsilon_{n+1} \varepsilon_{n+1}^T | \mathcal{F}_n] A_n^T = P_J \quad \text{a.s.}$$

Consequently, we have from the Cesaro mean convergence theorem that the predictable quadratic variation of the multidimensional martingale  $(M_n)$  satisfies

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \langle M \rangle_n = \lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{k=2}^n A_{k-1} \mathbb{E}[\varepsilon_k \varepsilon_k^T | \mathcal{F}_{k-1}] A_{k-1}^T = P_J \quad \text{a.s.}$$

Hence, we deduce from the strong law of large numbers for multidimensional martingales given by Theorem 4.3.16 in [18] that

$$\|M_n\|^2 = \mathcal{O}(n \log n) \quad \text{a.s.}$$

Therefore, there exists a finite positive random variable  $C$  such that for all  $n \geq 1$ ,

$$\|M_{n+1}\| \leq C \sqrt{n \log n} \quad \text{a.s.} \quad (6.17)$$

Hereafter, denote by  $Q_{n+1}$  the middle term of (6.15). We obtain from a simple Abel transform that

$$\begin{aligned} Q_{n+1} &= \sum_{k=1}^n \frac{1}{k} P_k^n (M_{k+1} - M_k) = \frac{1}{n} M_{n+1} + \sum_{k=1}^{n-1} \frac{1}{k} P_k^n M_{k+1} - \sum_{k=2}^n \frac{1}{k} P_k^n M_k \\ &= \frac{1}{n} M_{n+1} + \sum_{k=1}^{n-1} \left( \frac{1}{k} P_k^n - \frac{1}{k+1} P_{k+1}^n \right) M_{k+1} \\ &= \frac{1}{n} M_{n+1} + \sum_{k=1}^{n-1} \frac{1}{k(k+1)} (P_J - \Gamma_\varepsilon(v^*)) P_{k+1}^n M_{k+1} \end{aligned} \quad (6.18)$$

It follows from (6.16), (6.17), (6.18) that for all  $n \geq 1$ ,

$$\begin{aligned} \|Q_{n+1}\| &\leq C \left( \frac{\sqrt{n \log n}}{n} + \frac{\kappa}{n^\rho} \sum_{k=1}^{n-1} \frac{(k+1)^\rho}{k(k+1)} \sqrt{k \log k} \right) \quad \text{a.s.} \\ &\leq C \left( \left( \frac{\log n}{n} \right)^{1/2} + \frac{\kappa \sqrt{\log n}}{n^\rho} \sum_{k=1}^n \frac{1}{k^a} \right) \quad \text{a.s.} \end{aligned}$$

where  $a = 3/2 - \rho < 1$ . Consequently, we deduce that for all  $n \geq 1$ ,

$$\|Q_{n+1}\| \leq C \left( \left( \frac{\log n}{n} \right)^{1/2} + \frac{\kappa n^{1-a} \sqrt{\log n}}{(1-a)n^\rho} \right) \leq D \left( \frac{\log n}{n} \right)^{1/2} \quad \text{a.s.} \quad (6.19)$$

where

$$D = \frac{C(1-a+\kappa)}{1-a}.$$

The last term of (6.15) is much more difficult to handle. Denote for all  $n \geq 1$ ,

$$\Delta_n = \sum_{k=1}^n \frac{1}{k} P_k^n T_k \quad (6.20)$$

We recall that  $T_n = A_n \delta_n + B_n (\widehat{V}_n - v^*)$  where  $A_n$  and  $B_n$  are given by (6.12) and (6.13). We already saw from (5.13) that

$$\|\delta_n\| \leq \frac{2\sqrt{2}}{\varepsilon} \|\widehat{V}_n - v^*\|^2$$

which implies that

$$\|T_n\| \leq \frac{2\sqrt{2}}{\varepsilon} \|A_n\| \|\widehat{V}_n - v^*\|^2 + \|B_n\| \|\widehat{V}_n - v^*\|. \quad (6.21)$$

Moreover, it follows from (3.5) and (3.6) that

$$\lim_{n \rightarrow +\infty} \|A_n\| = \|G_\varepsilon^{-1/2}(v^*)\| \quad \text{and} \quad \lim_{n \rightarrow +\infty} \|B_n\| = 0 \quad \text{a.s.}$$

Consequently, we obtain from (3.5) and (6.21) that it exists a positive constant  $b = (4\kappa)^{-1}$  where  $\kappa$  is introduced in (6.16), such that for  $n$  large enough,

$$\|T_n\| \leq b \|\widehat{V}_n - v^*\| \quad \text{a.s.} \quad (6.22)$$

Define for all  $n \geq 1$ ,

$$L_n = \frac{1}{n} \sum_{k=1}^n \|\widehat{V}_k - v^*\|. \quad (6.23)$$

We deduce from (6.20) together with (6.16) and (6.22) that for all  $n \geq 1$ ,

$$\|\Delta_n\| \leq \frac{\kappa}{n^\rho} \sum_{k=1}^n \frac{k^\rho}{k} \|T_k\| \leq \frac{E}{n^\rho} + \frac{\kappa b}{n^\rho} \sum_{k=1}^n \frac{k^\rho}{k} \|\widehat{V}_k - v^*\| \leq \frac{E}{n^\rho} + \kappa b L_n \quad \text{a.s.} \quad (6.24)$$

where  $E$  is a finite positive random variable. Putting together the three contributions (6.16), (6.19) and (6.24), we obtain from (6.15) that for all  $n \geq 1$ ,

$$\|\widehat{U}_{n+1}\| \leq \frac{\tau \|\widehat{U}_1\| + E}{n^\rho} + D \left( \frac{\log n}{n} \right)^{1/2} + \kappa b L_n \quad \text{a.s.}$$

which implies that a finite positive random variable  $F$  exists and a constant  $0 < c < 1/2$  such that for all  $n \geq 1$ ,

$$\|\widehat{V}_{n+1} - v^*\| \leq F \left( \frac{\log n}{n} \right)^{1/2} + c L_n \quad \text{a.s.} \quad (6.25)$$

Herafter, we have from (6.23) and (6.25) that or all  $n \geq 1$ ,

$$\begin{aligned} L_{n+1} &= \left(1 - \frac{1}{n+1}\right)L_n + \frac{1}{n+1}\|\widehat{V}_{n+1} - v^*\|, \\ &\leq \left(1 - \frac{1}{n+1}\right)L_n + \frac{1}{n+1}\left(F\left(\frac{\log n}{n}\right)^{1/2} + cL_n\right) \quad \text{a.s} \\ &\leq \left(1 - \frac{d}{n+1}\right)L_n + \frac{F}{(n+1)}\left(\frac{\log n}{n}\right)^{1/2} \quad \text{a.s} \end{aligned}$$

where  $d = 1 - c$ . A straightforward induction yields that for all  $n \geq 1$ ,

$$L_n \leq \prod_{k=2}^n \left(1 - \frac{d}{k}\right)L_1 + \sum_{k=2}^n \prod_{i=k+1}^n \left(1 - \frac{d}{i}\right) \frac{F}{(k+1)} \left(\frac{\log k}{k}\right)^{1/2} \quad \text{a.s} \quad (6.26)$$

However, it is well-known that

$$\prod_{k=2}^n \left(1 - \frac{d}{k}\right) \leq \left(\frac{2}{n+1}\right)^d \quad \text{and} \quad \prod_{i=k+1}^n \left(1 - \frac{d}{i}\right) \leq \left(\frac{k+1}{n+1}\right)^d.$$

Hence, we obtain from (6.26) that for all  $n \geq 1$ ,

$$L_n \leq \left(\frac{2}{n+1}\right)^d L_1 + F \left(\frac{1}{n+1}\right)^d \sum_{k=2}^n \frac{(k+1)^d}{(k+1)} \left(\frac{\log k}{k}\right)^{1/2} \quad \text{a.s}$$

Since  $1/2 < d < 1$ , it implies that

$$L_n \leq \frac{2^d L_1}{n^d} + \frac{F(\log n)^{1/2}}{n^d} \sum_{k=1}^n \frac{1}{k^{3/2-d}} \leq \frac{2^d L_1}{n^d} + \frac{2F}{2d-1} \left(\frac{\log n}{n}\right)^{1/2} \quad \text{a.s}$$

leading to

$$L_n = \mathcal{O}\left(\left(\frac{\log n}{n}\right)^{1/2}\right) \quad \text{a.s} \quad (6.27)$$

Finally, it follows from (6.25) and (6.27) that

$$\|\widehat{V}_n - v^*\|^2 = \mathcal{O}\left(\frac{\log n}{n}\right) \quad \text{a.s.}$$

which completes the proof of (3.8).

• We now focus our attention on (3.9). We have from (6.7) that

$$\bar{S}_n - G_\varepsilon(v^*) = \frac{1}{n}M_n + \frac{1}{n}(I_J + R_n) + \frac{1}{n} \sum_{k=1}^n (G_\varepsilon(\widehat{V}_{k-1}) - G_\varepsilon(v^*)). \quad (6.28)$$

On the one hand, let  $M_n(u) = u^T M_n u$  where  $u \in \mathbb{R}^J$ . We already saw that  $(M_n(u))$  is a locally square-integrable martingale with increments bounded by  $8\|u\|^2$ . Moreover, its predictable quadratic variation satisfies  $\langle M(u) \rangle_n \leq 16n\|u\|^4$ . Therefore, we obtain from the third part of Theorem 1.3.24 in [18] that

$$|M_n(u)|^2 = \mathcal{O}(n \log n) \quad \text{a.s.}$$

which implies that

$$\frac{1}{n} \|M_n\| = \mathcal{O}\left(\left(\frac{\log n}{n}\right)^{1/2}\right) \quad \text{a.s.} \quad (6.29)$$

On the other hand, we already saw from Lemma 5.2 that

$$\frac{1}{n} \sum_{k=1}^n \|G_\varepsilon(\widehat{V}_{k-1}) - G_\varepsilon(v^*)\| \leq \frac{4L_n}{\varepsilon}$$

which ensures via (6.27) that

$$\frac{1}{n} \sum_{k=1}^n \|G_\varepsilon(\widehat{V}_{k-1}) - G_\varepsilon(v^*)\| = \mathcal{O}\left(\left(\frac{\log n}{n}\right)^{1/2}\right) \quad \text{a.s.} \quad (6.30)$$

Furthermore, we also have

$$\frac{1}{n} \|R_n\| \leq \left(\frac{\gamma \max(\nu)}{1 - \beta}\right) \frac{1}{n^\beta} \quad (6.31)$$

where  $\beta < 1/2$ . Consequently, we deduce the almost sure rate of convergence (3.9) for  $\overline{S}_n$  from the conjunction of (6.28), (6.29), (6.30) and (6.31). Finally, we obtain the almost sure rate of convergence (3.9) for  $\overline{S}_n^{-1}$  from the identity

$$\overline{S}_n^{-1} - G_\varepsilon^-(v^*) = \overline{S}_n^{-1} (G_\varepsilon(v^*) - \overline{S}_n) G_\varepsilon^-(v^*),$$

which completes the proof of Theorem 3.2.  $\square$

## 6.4 Proofs of the asymptotic normality results

*Proof of Theorem 3.3.* We now prove the asymptotic normality for the SGN algorithm.

- We recall from (6.15) that for all  $n \geq 1$ ,

$$\sqrt{n}(\widehat{V}_{n+1} - v^*) = -\sqrt{n}G_\varepsilon^{-1/2}(v^*)Q_{n+1} + R_n \quad (6.32)$$

where  $R_n = \sqrt{n}G_\varepsilon^{-1/2}(v^*)(P_0^n \widehat{U}_1 - \Delta_n)$  with  $\widehat{U}_1 = G_\varepsilon^{1/2}(v^*)(\widehat{V}_1 - v^*)$ ,

$$Q_{n+1} = \sum_{k=1}^n \frac{1}{k} P_k^n A_k \varepsilon_{k+1},$$

$$\Delta_n = \sum_{k=1}^n \frac{1}{k} P_k^n T_k,$$

On the one hand, we claim that the remainder  $R_n$  vanishes almost surely,

$$\lim_{n \rightarrow \infty} R_n = 0 \quad \text{a.s.} \quad (6.33)$$

As a matter of fact, we obviously have from (6.16) that

$$\lim_{n \rightarrow \infty} \sqrt{n} P_0^n \widehat{U}_1 = 0 \quad \text{a.s.}$$



Moreover, we deduce from the proof of Theorem 3.2 together with (6.21) that

$$\|T_n\| = \mathcal{O}\left(\frac{\log n}{n}\right) + \mathcal{O}\left(\frac{1}{n^\beta}\left(\frac{\log n}{n}\right)^{1/2}\right) = \mathcal{O}\left(\frac{1}{n^\beta}\left(\frac{\log n}{n}\right)^{1/2}\right) \quad \text{a.s.} \quad (6.34)$$

since  $0 < \beta < 1/2$  and

$$\|B_n\| = \mathcal{O}\left(\frac{1}{n^\beta}\right) \quad \text{a.s.}$$

where  $B_n$  is defined by (6.13). Therefore, we obtain from (6.34) that there exists a finite positive random variable  $C$  such that for all  $n \geq 1$

$$\|T_n\| \leq \frac{C}{n^\beta} \left(\frac{\log n}{n}\right)^{1/2} \quad \text{a.s.} \quad (6.35)$$

Consequently, it follows from (6.16) and (6.35) that for all  $n \geq 1$ ,

$$\|\Delta_n\| \leq \frac{C\kappa}{n^\rho} \sum_{k=1}^n \frac{k^\rho}{k^{1+\beta}} \left(\frac{\log k}{k}\right)^{1/2} \leq \frac{C\kappa\sqrt{\log n}}{n^\rho} \sum_{k=1}^n \frac{1}{k^a} \quad \text{a.s.}$$

where  $a = 3/2 + \beta - \rho$ , leading to

$$\|\Delta_n\| \leq \frac{D}{n^\beta} \left(\frac{\log n}{n}\right)^{1/2} \quad \text{a.s.} \quad (6.36)$$

with  $D = \kappa C/(1-a)$ . Hence, as  $\beta > 0$ , we infer from (6.36) that

$$\lim_{n \rightarrow \infty} \sqrt{n} \Delta_n = 0 \quad \text{a.s.}$$

which clearly implies that (6.33) is satisfied. On the other hand,  $Q_{n+1}$  is a sum of weighted martingale differences. We deduce from the first part of Proposition B.2 in [50] that

$$\sqrt{n} Q_{n+1} \xrightarrow{\mathcal{L}} \mathcal{N}(0, \Sigma) \quad (6.37)$$

where the asymptotic covariance matrix  $\Sigma$  is given by the integral form

$$\begin{aligned} \Sigma &= \int_0^\infty \left( \exp\left(-\left(\Gamma_\varepsilon(v^*) - \frac{1}{2}P_J\right)s\right) \right)^2 ds \\ &= \int_0^\infty \left( \exp\left(-2s\left(\Gamma_\varepsilon(v^*) - \frac{1}{2}P_J\right)\right) \right) ds \\ &= \frac{1}{2} \int_0^\infty \left( \exp\left(-t\left(\Gamma_\varepsilon(v^*) - \frac{1}{2}P_J\right)\right) \right) dt \\ &= \frac{1}{2} \left[ -\left(\Gamma_\varepsilon(v^*) - \frac{1}{2}P_J\right)^- \exp\left(-t\left(\Gamma_\varepsilon(v^*) - \frac{1}{2}P_J\right)\right) \right]_0^\infty \\ &= \left(2\Gamma_\varepsilon(v^*) - P_J\right)^-. \end{aligned}$$

Therefore, we obtain from (6.37) that

$$\sqrt{n} G_\varepsilon^{-1/2}(v^*) Q_{n+1} \xrightarrow{\mathcal{L}} \mathcal{N}\left(0, G_\varepsilon^{-1/2}(v^*) (2\Gamma_\varepsilon(v^*) - P_J)^- G_\varepsilon^{-1/2}(v^*)\right). \quad (6.38)$$

Finally, it follows from (6.32) together with (6.33) and (6.38) that

$$\sqrt{n}(\widehat{V}_n - v^*) \xrightarrow{\mathcal{L}} \mathcal{N}\left(0, G_\varepsilon^{-1/2}(v^*)(2\Gamma_\varepsilon(v^*) - P_J)^- G_\varepsilon^{-1/2}(v^*)\right),$$

which is exactly what we wanted to prove.

- It only remains to prove the asymptotic normality (3.14). We already saw from inequality (5.11) that for all  $v \in \mathbb{R}^J$ ,

$$H_\varepsilon(v) - H_\varepsilon(v^*) \leq \frac{1}{2\varepsilon} \|v - v^*\|^2.$$

Moreover, we have from (2.7) that  $\widehat{W}_n - W_\varepsilon(\mu, \nu)$  can be splitted into two terms,

$$\begin{aligned} \sqrt{n}(\widehat{W}_n - W_\varepsilon(\mu, \nu)) &= \frac{1}{\sqrt{n}} \sum_{k=1}^n (H_\varepsilon(v^*) - h_\varepsilon(X_k, \widehat{V}_{k-1})), \\ &= \frac{1}{\sqrt{n}} \sum_{k=1}^n (H_\varepsilon(\widehat{V}_{k-1}) - h_\varepsilon(X_k, \widehat{V}_{k-1})) - \frac{1}{\sqrt{n}} \sum_{k=1}^n (H_\varepsilon(\widehat{V}_{k-1}) - H_\varepsilon(v^*)). \end{aligned} \quad (6.39)$$

The second term in equation (6.39) goes to 0 a.s. thanks to the almost sure rate of convergence (3.8),

$$\frac{1}{\sqrt{n}} \sum_{k=1}^n (H_\varepsilon(\widehat{V}_{k-1}) - H_\varepsilon(v^*)) \leq \frac{1}{2\varepsilon\sqrt{n}} \sum_{k=1}^n \|\widehat{V}_{k-1} - v^*\|^2 = \mathcal{O}\left(\frac{\log^2(n)}{\sqrt{n}}\right) \quad \text{a.s.}$$

Finally, the first term in equation (6.39) is dealt using argument from the proof of Theorem 3.5 in [5], allowing to prove that it satisfies the asymptotic normality (3.14). This completes the proof of Theorem 3.3.  $\square$

## 6.5 Proofs of the non-asymptotic rates of convergence

We first detail some keystone results related to the use of the Kurdyka-Łojasiewicz functional inequality that is at the heart of the proof of Theorem 3.4.

### 6.5.1 Kurdyka-Łojasiewicz inequality.

The analysis that we carry out is essentially based on the so-called Kurdyka-Łojasiewicz functional inequality. We refer to the initial works [31, 32] and to [10, 11] for the use of such inequality in deterministic optimization. To this end, let  $\tilde{H}_\varepsilon$  be the positive and convex function defined, for all  $v \in \langle \mathbf{v}_J \rangle^\perp$ , by

$$\tilde{H}_\varepsilon(u) = H_\varepsilon(G_*^{-1/2}u) - H_\varepsilon(G_*^{-1/2}u^*),$$

where  $G_*^{-1/2}$  stands for the square root of the Moore-Penrose inverse of  $G_* = G_\varepsilon(v^*)$ , and  $u^* = G_*^{1/2}v^*$ . Since  $G_*^{-1/2}$  is symmetric, we notice that  $\nabla \tilde{H}_\varepsilon(u) = G_*^{-1/2} \nabla H_\varepsilon(G_*^{-1/2}u)$  and

$\nabla^2 \tilde{H}_\varepsilon(u) = G_*^{-1/2} \nabla^2 H_\varepsilon(G_*^{-1/2} u) G_*^{-1/2}$ . Hence, we obtain from the upper bounds (5.7) and (5.8) that for all  $u \in \langle \mathbf{v}_J \rangle^\perp$ ,

$$\|\nabla \tilde{H}_\varepsilon(u)\| \leq 2\lambda_{\max}(G_*^{-1/2}), \quad (6.40)$$

and

$$\lambda_{\max}(\nabla^2 \tilde{H}_\varepsilon(u)) \leq \frac{1}{\varepsilon} \lambda_{\max}(G_*^-). \quad (6.41)$$

Consequently, it follows from the Taylor-Lagrange formula that for all  $u \in \langle \mathbf{v}_J \rangle^\perp$ ,

$$\tilde{H}_\varepsilon(u) \leq \frac{1}{2\varepsilon} \lambda_{\max}(G_*^-) \|u - u^*\|^2. \quad (6.42)$$

First of all, we verify that the function  $\tilde{H}_\varepsilon$  satisfies a Kurdyka-Łojasiewicz inequality as stated in the next proposition. We refer the reader to [22] and the references therein for further details on this topic.

**Proposition 6.1.** *There exist two positive constants  $m_\varepsilon < M$  such that, for all  $u \in \langle \mathbf{v}_J \rangle^\perp$  with  $u \neq u^*$ ,*

$$0 < m_\varepsilon \leq \|\nabla \tilde{H}_\varepsilon(u)\|^2 + \frac{\|\nabla \tilde{H}_\varepsilon(u)\|^2}{\tilde{H}_\varepsilon(u)} \leq M < +\infty. \quad (6.43)$$

Moreover, the constant  $m_\varepsilon$  can be chosen as

$$m_\varepsilon = \varepsilon \lambda_{\min}(G_\varepsilon(v^*)) \min\left(1, \frac{\varepsilon}{4}\right). \quad (6.44)$$

The proof of this key inequality is postponed to Appendix B. From equation (6.44), we clearly observe that the magnitude of the constant  $m_\varepsilon$  depends on  $\varepsilon$ . Nevertheless, in the analysis carried out in this paper, the regularization parameter  $\varepsilon$  is held fixed, and we will not be interested in deriving sharp upper bounds depending on  $\varepsilon$  for the mean square error of  $\hat{V}_n$ . We believe that a careful analysis of the role of  $\varepsilon$  on the convergence of the SGN algorithm is a difficult issue that is left open for future investigation.

### 6.5.2 Choice of a Lyapunov function

Hereafter, a key step in our analysis is based on the Lyapunov function  $\Phi$  defined, for all  $u \in \langle \mathbf{v}_J \rangle^\perp$ , by

$$\Phi(u) = \tilde{H}_\varepsilon(u) \exp(\tilde{H}_\varepsilon(u)). \quad (6.45)$$

On the one hand, it follows from the elementary inequality  $\exp(x) \leq 1 + x \exp(x)$  that for all  $u \neq u^*$ ,

$$\frac{\Phi(u)}{\tilde{H}_\varepsilon(u)} \leq 1 + \Phi(u). \quad (6.46)$$

We shall repeatedly use inequality (6.46) in all the sequel. On the other hand, we can easily compute for all  $u \neq u^*$ ,

$$\nabla \Phi(u) = \left(1 + \frac{1}{\tilde{H}_\varepsilon(u)}\right) \Phi(u) \nabla \tilde{H}_\varepsilon(u),$$

which implies that

$$\langle \nabla \Phi(u), \nabla \tilde{H}_\varepsilon(u) \rangle = \left(1 + \frac{1}{\tilde{H}_\varepsilon(u)}\right) \Phi(u) \|\nabla \tilde{H}_\varepsilon(u)\|^2.$$

Consequently, we deduce from Proposition 6.1 that  $m_\varepsilon \Phi(u) \leq \langle \nabla \Phi(u), \nabla \tilde{H}_\varepsilon(u) \rangle \leq M \Phi(u)$ . Moreover, if  $\Sigma$  denotes a positive semi-definite matrix, by an application of Proposition 6.1, we also obtain the following *key* lower bound

$$\langle \nabla \Phi(u), \Sigma \nabla \tilde{H}_\varepsilon(u) \rangle \geq m_\varepsilon \lambda_{\min}(\Sigma) \Phi(u), \quad (6.47)$$

that will be useful to control the non-asymptotic rate of convergence of the SGN algorithm. Finally, as remarked in [22], a straightforward computation leads, for all  $u \neq u^*$ , to

$$\nabla^2 \Phi(u) = \Phi(u) \left( \left(1 + \frac{2}{\tilde{H}_\varepsilon(u)}\right) \nabla \tilde{H}_\varepsilon(u) \nabla \tilde{H}_\varepsilon(u)^T + \left(1 + \frac{1}{\tilde{H}_\varepsilon(u)}\right) \nabla^2 \tilde{H}_\varepsilon(u) \right).$$

Hence, using the fact that  $\lambda_{\max}(\nabla \tilde{H}_\varepsilon(u) \nabla \tilde{H}_\varepsilon(u)^T) = \|\nabla \tilde{H}_\varepsilon(u)\|^2$ , we obtain that for all  $u \neq u^*$ ,

$$\lambda_{\max}(\nabla^2 \Phi(u)) \leq \Phi(u) \left( \left(1 + \frac{2}{\tilde{H}_\varepsilon(u)}\right) \|\nabla \tilde{H}_\varepsilon(u)\|^2 + \left(1 + \frac{1}{\tilde{H}_\varepsilon(u)}\right) \lambda_{\max}(\nabla^2 \tilde{H}_\varepsilon(u)) \right).$$

Therefore, using inequality (6.46) together with the upper bounds (6.40) and (6.41), we obtain that for all  $u \in \langle \mathbf{v}_J \rangle^\perp$  with  $u \neq u^*$ ,

$$\lambda_{\max}(\nabla^2 \Phi(u)) \leq \delta_\varepsilon \lambda_{\max}(G_*^-) (1 + \Phi(u)) \quad (6.48)$$

where  $\delta_\varepsilon = 2(6 + \varepsilon^{-1})$ . Inequality (6.48) will also be crucial to derive the non-asymptotic rate of convergence of the SGN algorithm. Finally, thanks to the following result, we will be able to relate the study of the Lyapunov function  $\Phi$  to the quadratic risk of  $\hat{V}_n$  and  $S_n$ .

**Proposition 6.2.** *There exists a positive constant  $d_\varepsilon$  such that for all  $u \in \langle \mathbf{v}_J \rangle^\perp$ ,*

$$\|u - u^*\|^2 \leq d_\varepsilon \Phi(u). \quad (6.49)$$

*Proof.* First, one can verify that for all  $u \in \langle \mathbf{v}_J \rangle^\perp$  in a neighborhood of  $u^*$  with  $u \neq u^*$ , the function  $\|u - u^*\|^{-2} \tilde{H}_\varepsilon(u)$  is lower bounded. Moreover, using that  $\tilde{H}_\varepsilon$  is a convex function on  $\langle \mathbf{v}_J \rangle^\perp$  that attains its minimal value at  $u^*$  with a non-degenerate minimum, we also have

$$\liminf_{\|u\| \rightarrow +\infty} \frac{\tilde{H}_\varepsilon(u)}{\|u\|} > 0.$$

It implies that for any positive  $t$ ,

$$\lim_{\|u\| \rightarrow +\infty} \frac{\exp(t \tilde{H}_\varepsilon(u))}{\|u - u^*\|^2} = +\infty.$$

Since  $\tilde{H}_\varepsilon(u) \geq 0$ , one always has  $\Phi(u) \geq \tilde{H}_\varepsilon(u)$ . Consequently, for all  $u \neq u^*$ ,

$$\begin{aligned} \frac{\Phi(u)}{\|u - u^*\|^2} &= \frac{\Phi(u)}{\|u - u^*\|^2} \mathbb{1}_{\|u - u^*\| \leq 1} + \frac{\Phi(u)}{\|u - u^*\|^2} \mathbb{1}_{\|u - u^*\| \geq 1}, \\ &\geq \frac{\tilde{H}_\varepsilon(u)}{\|u - u^*\|^2} \mathbb{1}_{\|u - u^*\| \leq 1} + \min_{\|u - u^*\| \geq 1} \frac{\tilde{H}_\varepsilon(u) \exp(\tilde{H}_\varepsilon(u))}{\|u - u^*\|^2} \geq \frac{1}{d_\varepsilon} \end{aligned}$$

for some positive constant  $d_\varepsilon$ , where we used the local behavior around  $u^*$  of  $\tilde{H}_\varepsilon$  to derive a lower bound for the first term and the asymptotic behavior of  $\tilde{H}_\varepsilon$  for the second one with  $t = 1$ , which is exactly what we wanted to prove.  $\square$

### 6.5.3 A recursive inequality and proof of Theorem 3.4

We first describe the one-step evolution of the sequence  $(\Phi(G_*^{1/2}\hat{V}_n))$  where  $(\hat{V}_n)$  is the recursive sequence defined by (2.6) corresponding to the SGN algorithm. From this analysis, we shall also deduce the rate of convergence of the expected quadratic risk associated with  $\hat{V}_n$  and  $\bar{S}_n$ . Denote  $\tilde{U}_n = G_*^{1/2}\hat{V}_n$ .

**Proposition 6.3.** *Assume that  $\alpha \in [0, 1/2[$  and that  $\alpha + \beta < 1/2$ . Then, there exist an integer  $n_0 \geq J$  and a positive constant  $c_\varepsilon > 0$  such that, for all  $n \geq n_0$ ,*

$$\mathbb{E}[\Phi(\tilde{U}_{n+1})|\mathcal{F}_n] \leq \left(1 - \frac{m_\varepsilon n^\alpha \lambda_{\min}(G_*) \lambda_{\min}(S_n^{-1})}{2}\right) \Phi(\tilde{U}_n) + c_\varepsilon n^{2\alpha} \lambda_{\max}^2(S_n^{-1}). \quad (6.50)$$

The proof of Proposition 6.3 is postponed to Appendix B. We are now in position to establish the non-asymptotic rates of convergence for the SGN algorithm.

*Proof of Theorem 3.4.* In the proof, we use the notation  $c_\varepsilon$  to denote a positive constant (depending on  $\varepsilon$  and possibly on  $\alpha$  and  $\beta$ ) that is independent from  $n$  and whose value may change from line to line. We only consider the situation where  $\alpha \in ]0, 1/2[$  and our analysis is based on the function  $\Phi$ . We will establish that for all  $n \geq 1$ ,

$$\mathbb{E}[\Phi(\tilde{U}_n)] \leq \frac{c_\varepsilon}{n^{1-\alpha}}. \quad (6.51)$$

Step 1: Preliminary rate. Our starting point is inequality (6.50) that we combine with (B.6) and (B.7) to obtain that there exists an integer  $n_0$  such that, for all  $n \geq n_0$ ,

$$\mathbb{E}[\Phi(\tilde{U}_{n+1})|\mathcal{F}_n] \leq (1 - c_1(n)n^{-1+\alpha})\Phi(\tilde{U}_n) + c_2(n)n^{2(\alpha+\beta-1)} \quad (6.52)$$

where  $n_J = n/J$ ,

$$c_1(n) = \frac{m_\varepsilon \lambda_{\min}(G_*)}{8 + 4\gamma \max(\nu) + 2n^{-1}},$$

$$c_2(n) = \frac{c_\varepsilon}{J^{2(\beta-1)}(\gamma \min(\nu)(1 - n_J^{-1})^{1-\beta} + (1 - 2\gamma \min(\nu))n_J^{-1})^2}.$$

By taking the expectation on both sides of (6.52), we obtain that for all  $n \geq n_0$ ,

$$\mathbb{E}[\Phi(\tilde{U}_{n+1})] \leq (1 - c_1(n)n^{-1+\alpha}) \mathbb{E}[\Phi(\tilde{U}_n)] + c_2(n)n^{-2(1-\alpha-\beta)}.$$

Hereafter, it is not hard to see that for  $n$  large enough, there exist two positive constants  $c_1$  and  $c_2$ , depending on  $\varepsilon$ , such that  $c_1(n) \geq c_1$  and  $c_2(n) \leq c_2$ . Hence, there exists an integer  $n_0$  such that, for all  $n \geq n_0$ ,

$$\mathbb{E}[\Phi(\tilde{U}_{n+1})] \leq (1 - c_1 n^{-1+\alpha}) \mathbb{E}[\Phi(\tilde{U}_n)] + c_2 n^{-2(1-\alpha-\beta)}.$$

Therefore, it follows from the proof of Lemma A.3 in [5] that there exist a positive constant  $c_\varepsilon$  and an integer  $n_0$  such that, for all  $n \geq n_0$ ,

$$\mathbb{E}[\Phi(\tilde{U}_n)] \leq \frac{c_\varepsilon}{n^{1-2(\alpha+\beta)}}.$$

We emphasize that at this stage, we do not obtain the announced result that necessitates further work with a plug-in strategy. The rest of the proof details this additional step.

Step 2: Plug-in. Let us now explain how one may improve the above result from  $n^{-(1-2(\alpha+\beta))}$  to  $n^{-(1-2\alpha)}$ . First of all, thanks to Proposition 6.2, we obtain via Step 1 that

$$\mathbb{E}[\|\tilde{U}_n - u^*\|^2] \leq \frac{c_\varepsilon}{n^{1-2(\alpha+\beta)}}. \quad (6.53)$$

Next, we shall consider the study of the convergence rate of  $\bar{S}_n$  to improve the pessimistic bounds (B.6) and (B.7). If  $p_n$  denotes the largest integer such that  $p_n J \leq n$ , we already saw from (4.3) and (6.4) that for all  $n \geq 1$ ,

$$\bar{S}_n = \bar{\Sigma}_n + \bar{R}_n$$

where

$$\begin{aligned} \bar{\Sigma}_n &= \frac{1}{n} \sum_{k=1}^n \nabla_v h_\varepsilon(X_k, \hat{V}_{k-1}) \nabla_v h_\varepsilon(X_k, \hat{V}_{k-1})^T, \\ \bar{R}_n &= \frac{1}{n} I_J + \frac{1}{n} \left( \sum_{m=1}^{p_n} m^{-\beta} \right) \gamma \operatorname{diag}(\nu) + \frac{1}{n} \sum_{k=p_n+1}^n \gamma \left( 1 + \lfloor \frac{k}{J} \rfloor \right)^{-\beta} Z_k Z_k^T. \end{aligned}$$

On the one hand, it is not hard to see that it exists a positive constant  $c_3$  such that

$$\|\bar{R}_n\|_F^2 \leq \frac{c_3}{n^{2\beta}}. \quad (6.54)$$

On the other hand, starting from the fact that

$$\bar{\Sigma}_{n+1} = \bar{\Sigma}_n + \frac{1}{n+1} \left( \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n)^T - \bar{\Sigma}_n \right),$$

we obtain that

$$\begin{aligned} \|\bar{\Sigma}_{n+1} - G_*\|_F^2 &= \|\bar{\Sigma}_n - G_*\|_F^2 + \frac{1}{(n+1)^2} \|\nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n)^T - \bar{\Sigma}_n\|_F^2 \\ &\quad + \frac{2}{n+1} \langle \bar{\Sigma}_n - G_*, \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n)^T - \bar{\Sigma}_n \rangle_F. \end{aligned}$$

Therefore, by taking the conditional expectation on both sides of the above equality, we obtain that

$$\begin{aligned} \mathbb{E}[\|\bar{\Sigma}_{n+1} - G_*\|_F^2 | \mathcal{F}_n] &= \|\bar{\Sigma}_n - G_*\|_F^2 + \frac{2}{n+1} \langle \bar{\Sigma}_n - G_*, G_\varepsilon(\hat{V}_n) - \bar{\Sigma}_n \rangle_F \\ &\quad + \frac{1}{(n+1)^2} \mathbb{E}[\|\nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n)^T - \bar{\Sigma}_n\|_F^2 | \mathcal{F}_n], \\ &= \|\bar{\Sigma}_n - G_*\|_F^2 - \frac{2}{n+1} \|\bar{\Sigma}_n - G_*\|_F^2 + \frac{2}{n+1} \langle \bar{\Sigma}_n - G_*, G_\varepsilon(\hat{V}_n) - G_* \rangle_F \\ &\quad + \frac{1}{(n+1)^2} \mathbb{E}[\|\nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n)^T - \bar{\Sigma}_n\|_F^2 | \mathcal{F}_n], \\ &\leq \|\bar{\Sigma}_n - G_*\|_F^2 \left( 1 - \frac{1}{n+1} \right) + \frac{1}{n+1} \|G_\varepsilon(\hat{V}_n) - G_*\|_F^2 \\ &\quad + \frac{1}{(n+1)^2} \mathbb{E}[\|\nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n)^T - \bar{\Sigma}_n\|_F^2 | \mathcal{F}_n] \quad (6.55) \end{aligned}$$

where the last line follows from Cauchy-Schwarz and Young inequalities. Moreover, we deduce from inequality (5.6) that

$$\begin{aligned} & \mathbb{E} \left[ \left\| \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n) \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n)^T - \bar{\Sigma}_n \right\|_F^2 | \mathcal{F}_n \right] \\ & \leq 2\mathbb{E} \left[ \left\| \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n) \nabla_v h_\varepsilon(X_{n+1}, \widehat{V}_n)^T \right\|_F^2 | \mathcal{F}_n \right] + 2\|\bar{\Sigma}_n\|_F^2 \leq 2 \times 4^2 + 2 \times 4 = 40. \end{aligned}$$

Furthermore, we obtain from (5.15) that

$$\|G_\varepsilon(\widehat{V}_n) - G_*\|_F^2 \leq \frac{16J}{\varepsilon^2} \|\widehat{V}_n - v^*\|^2.$$

Using the previous bounds in (6.55) leads to the recursive inequality

$$\mathbb{E} \left[ \|\bar{\Sigma}_{n+1} - G_*\|_F^2 \right] \leq \left(1 - \frac{1}{n+1}\right) \mathbb{E} \left[ \|\bar{\Sigma}_n - G_*\|_F^2 \right] + \frac{40}{(n+1)^2} + \frac{16J}{(n+1)\varepsilon^2} \mathbb{E} \left[ \|\widehat{V}_n - v^*\|^2 \right].$$

Since  $\tilde{U}_n - u^* = G_*^{1/2}(\widehat{V}_n - v^*)$ , inequality (6.53) ensures that

$$\mathbb{E} \left[ \|\bar{\Sigma}_{n+1} - G_*\|_F^2 \right] \leq \left(1 - \frac{1}{n+1}\right) \mathbb{E} \left[ \|\bar{\Sigma}_n - G_*\|_F^2 \right] + \frac{40}{(n+1)^2} + \frac{16Jc_\varepsilon \lambda_{\max}(G_*^-)}{n^{2(1-\alpha-\beta)}\varepsilon^2}.$$

By applying Lemma A.3 in [5], we obtain that for  $n$  large enough,

$$\mathbb{E} \left[ \|\bar{\Sigma}_n - G_*\|_F^2 \right] \leq \frac{c_\varepsilon}{n^{(1-2(\alpha+\beta))}}.$$

From this last inequality and (6.54), we conclude that for  $n$  large enough,

$$\mathbb{E} \left[ \|\bar{S}_n - G_*\|_F^2 \right] \leq \frac{c_\varepsilon}{n^{(1-2(\alpha+\beta))}} + \frac{c_3}{n^{2\beta}}. \quad (6.56)$$

Since  $4\beta < 1 - 2\alpha$ , it follows that  $1 - 2(\alpha + \beta) > 2\beta$ , which means that  $n^{-(1-2(\alpha+\beta))}$  decays faster than  $n^{-2\beta}$  and the second inequality in (3.18) holds true.

• From now on, we focus our attention on the first inequality in (3.18). It follows from (6.50) and the previous calculation that for  $n$  large enough,

$$\mathbb{E}[\Phi(\tilde{U}_{n+1})] \leq (1 - c_1(n)n^{-1+\alpha})\mathbb{E}[\Phi(\tilde{U}_n)] + c_\varepsilon n^{2\alpha-2}\mathbb{E}[\lambda_{\max}^2(\bar{S}_n^{-1})], \quad (6.57)$$

using that  $\bar{S}_n^{-1} = nS_n^{-1}$ . The identity  $\bar{S}_n^{-1} - G_*^- = G_*^-(G_* - \bar{S}_n)\bar{S}_n^{-1}$  implies that

$$\lambda_{\max}(\bar{S}_n^{-1}) \leq \|\bar{S}_n^{-1} - G_*^-\|_2 + \|G_*^-\|_2 \leq \lambda_{\max}(G_*^-)\lambda_{\max}(\bar{S}_n^{-1})\|\bar{S}_n - G_*\|_2 + \|G_*^-\|_2.$$

It follows from inequality (B.6) that for  $n$  large enough,  $\lambda_{\max}(\bar{S}_n^{-1}) = n\lambda_{\max}(S_n^{-1}) \leq c_4 n^\beta$  where  $c_4$  is a positive constant. It ensures that for  $n$  large enough,

$$\lambda_{\max}^2(\bar{S}_n^{-1}) \leq 2c_4^2 \lambda_{\max}^2(G_*^-) n^{2\beta} \|\bar{S}_n - G_*\|_2^2 + 2\|G_*^-\|_2^2.$$

Since  $\|\bar{S}_n - G_*\|_2^2 \leq \|\bar{S}_n - G_*\|_F^2$ , we obtain from inequality (6.56) that for  $n$  large enough,

$$\mathbb{E} \left[ \lambda_{\max}^2(\bar{S}_n^{-1}) \right] \leq c_\varepsilon \left( 1 + \frac{n^{4\beta}}{n^{1-2\alpha}} \right) + 2\|G_*^-\|_2^2.$$

Consequently, we deduce from the condition  $4\beta < 1 - 2\alpha$  that there exists a positive constant  $C_\varepsilon$  such that for all  $n \geq 1$ ,

$$\mathbb{E}[\lambda_{\max}^2(\bar{S}_n^{-1})] \leq C_\varepsilon. \quad (6.58)$$

It follows from (6.57) and (6.58) that there exist two positive constants  $c_1$  and  $c_\varepsilon$  such that for  $n$  large enough,

$$\mathbb{E}[\Phi(\tilde{U}_{n+1})] \leq (1 - c_1 n^{-1+\alpha})\mathbb{E}[\Phi(\tilde{U}_n)] + c_\varepsilon n^{2\alpha-2}.$$

Hereafter, using once again Lemma A.3 in [5], we obtain that for all  $n \geq 1$ ,

$$\mathbb{E}[\Phi(\tilde{U}_n)] \leq \frac{c_\varepsilon}{n^{1-\alpha}}$$

which is exactly the announced inequality (6.51). Finally, as  $\hat{V}_n - v_* = G_*^{-1/2}(\tilde{U}_n - u^*)$ , the first inequality in (3.18) clearly follows from (6.51) together with Proposition 6.2.

• It only remains to prove the inequalities (3.19) and (3.20). We already saw the decomposition

$$\begin{aligned} \widehat{W}_n - W_\varepsilon(\mu, \nu) &= \frac{1}{n} \sum_{k=1}^n (H_\varepsilon(v^*) - h_\varepsilon(X_k, \hat{V}_{k-1})), \\ &= \frac{1}{n} \sum_{k=1}^n \xi_k - \frac{1}{n} \sum_{k=1}^n (H_\varepsilon(\hat{V}_{k-1}) - H_\varepsilon(v^*)), \end{aligned} \quad (6.59)$$

where the martingale increment  $\xi_k = -h_\varepsilon(X_k, \hat{V}_{k-1}) + H_\varepsilon(\hat{V}_{k-1})$ . As  $\mathbb{E}[\xi_{k+1} | \mathcal{F}_k] = 0$ , it follows from (6.59) together with inequality (5.11) and the first inequality in (3.18) that for all  $n \geq 1$ ,

$$\left| \mathbb{E}[\widehat{W}_n] - W_\varepsilon(\mu, \nu) \right| = \left| \frac{1}{n} \sum_{k=1}^n \mathbb{E}[H_\varepsilon(\hat{V}_{k-1}) - H_\varepsilon(v^*)] \right| \leq \frac{c_\varepsilon}{n} \sum_{k=1}^n \frac{1}{k^{1-\alpha}} \leq \frac{c_{\varepsilon, \alpha}}{n^{1-\alpha}},$$

which proves inequality (3.19).

Regarding the  $\mathbb{L}^1$  risk  $\mathbb{E}[|\widehat{W}_n - W_\varepsilon(\mu, \nu)|]$ , we still use the decomposition (6.59). The triangle inequality and inequality (5.11) imply that

$$\begin{aligned} \mathbb{E}[|\widehat{W}_n - W_\varepsilon(\mu, \nu)|] &\leq \frac{1}{n} \mathbb{E}\left[\left|\sum_{k=1}^n \xi_k\right|\right] + \frac{1}{n} \sum_{k=1}^n \mathbb{E}[|H_\varepsilon(\hat{V}_{k-1}) - H_\varepsilon(v^*)|], \\ &\leq \frac{1}{n} \left( \mathbb{E}\left(\sum_{k=1}^n \xi_k^2\right)^{1/2} + \frac{1}{2\varepsilon n} \sum_{k=1}^n \mathbb{E}[\|\hat{V}_{k-1} - v^*\|^2] \right), \end{aligned}$$

where the last line comes from the Cauchy-Schwarz inequality. Let us now prove that  $\sup_{k \geq 1} \mathbb{E}[\xi_k^2] < +\infty$ . To this end, we observe that

$$\mathbb{E}[\xi_k^2] = \mathbb{E}\left(\mathbb{E}\left[\left(h_\varepsilon(X_k, \hat{V}_{k-1}) - H_\varepsilon(\hat{V}_{k-1})\right)^2 \middle| \mathcal{F}_k\right]\right) \leq \mathbb{E}[h_\varepsilon^2(X_k, \hat{V}_{k-1})],$$



thanks to the property that  $\mathbb{E}[h_\varepsilon(X_k, \widehat{V}_{k-1}) | \mathcal{F}_k] = H_\varepsilon(\widehat{V}_{k-1})$ . Using that  $\|\partial_v h_\varepsilon(x, v)\| = \|\pi(x, v) - \nu\| \leq 2$ , it follows by integration that  $h_\varepsilon(x, v) \leq h_\varepsilon(x, v^*) + 2\|v - v^*\|$ . We then deduce that

$$\mathbb{E}[\xi_k^2] \leq 2\mathbb{E}[h_\varepsilon^2(X, v^*)] + 4\mathbb{E}[\|\widehat{V}_{k-1} - v^*\|^2]$$

Using the first inequality in (3.18), it follows that  $(\mathbb{E}[\|\widehat{V}_{k-1} - v^*\|^2])_{k \geq 1}$  is a bounded sequence. Moreover, arguing as in the proof of [5, Theorem 3.5], the condition  $\int_{\mathcal{X}} c^2(x, y_j) d\mu(x) < +\infty$ , for any  $1 \leq j \leq J$ , implies that  $\mathbb{E}[h_\varepsilon^2(X, v^*)]$  is finite. Therefore, we conclude that  $\sup_{k \geq 1} \mathbb{E}[\xi_k^2] < +\infty$ . Hence, using once again the first inequality in (3.18) together with a conditional expectation argument, we obtain that

$$\mathbb{E}[\|\widehat{W}_n - W_\varepsilon(\mu, \nu)\|] \leq \frac{c_\varepsilon}{n} \sqrt{n} + \frac{c_\varepsilon}{2\varepsilon n} \sum_{k=1}^n \frac{1}{k^{1-\alpha}} \leq \frac{c_\varepsilon}{\sqrt{n}},$$

which proves inequality (3.20). This achieves the proof of Theorem 3.4.  $\square$

## A Appendix - Proofs of auxiliary results

This appendix contains the proofs of some auxiliary results of the paper.

### A.1 Proof of Lemma 5.2

Since  $\nabla_v h_\varepsilon(X, v) = \pi(X, v) - \nu$ , we first remark that

$$\begin{aligned} G_\varepsilon(v) &= L(v) + \nu \nu^T - \nu \mathbb{E}[\pi(X, v)]^T - \mathbb{E}[\pi(X, v)] \nu^T, \\ &= L(v) - \nu \nu^T - \nu \nabla H_\varepsilon(v)^T - \nabla H_\varepsilon(v) \nu^T, \end{aligned} \quad (\text{A.1})$$

where  $L(v) = \mathbb{E}[\pi(X, v) \pi(X, v)^T]$ . For  $u \in \langle \nu_J \rangle^\perp$  and  $t \in [0, 1]$ , we define the real-valued function  $\phi_u(t) = u^T G_\varepsilon(v_t) u$  with  $v_t = v^* + t(v - v^*)$ . It follows from the decomposition (A.1) that

$$\phi_u(t) = u^T L(v_t) u - \langle u, \nu \rangle^2 - 2\langle u, \nu \rangle \langle \Phi(t), u \rangle$$

where  $\Phi(t) = \nabla H_\varepsilon(v_t)$  is a vector-valued function satisfying

$$\Phi'(t) = \nabla^2 H_\varepsilon(v_t)(v - v^*) \quad \text{and} \quad \Phi''(t) = \nabla^3 H_\varepsilon(v_t)[v - v^*, v - v^*],$$

where  $\nabla^3 H_\varepsilon$  stands for the third-order tensor derivative of  $H_\varepsilon$ . Now, since  $\Phi(1) = \nabla H_\varepsilon(v)$  and  $\Phi(0) = \nabla H_\varepsilon(v^*) = 0$ , and using the property that  $\phi_u(1) - \phi_u(0) = \int_0^1 \phi'_u(t) dt$ , we obtain that

$$u^T (L(v) - L(v^*)) u - 2\langle u, \nu \rangle \langle \nabla H_\varepsilon(v), u \rangle = \int_0^1 \phi'_u(t) dt. \quad (\text{A.2})$$

We clearly have

$$\phi'_u(t) = \frac{\partial}{\partial t} u^T L(v_t) u - 2\langle u, \nu \rangle \langle \Phi'(t), u \rangle$$

where

$$\frac{\partial}{\partial t} u^T L(v_t) u = \langle \nabla_v \psi_u(v_t), v - v^* \rangle \quad \text{with} \quad \psi_u(v) = u^T L(v) u.$$

In addition, we also have  $\langle \Phi'(t), u \rangle = u^T \nabla^2 H_\varepsilon(v_t)(v - v^*)$ . Since

$$u^T L(v) u = \mathbb{E}[\langle u, \pi(X, v) \rangle^2]$$

and

$$\nabla_v \pi(x, v) = \frac{1}{\varepsilon} \left( \text{diag}(\pi(x, v)) - \pi(x, v) \pi(x, v)^T \right),$$

one obtains that

$$\nabla_v \psi_u(v_t) = \frac{2}{\varepsilon} \mathbb{E}[\langle u, \pi(X, v_t) \rangle (\text{diag}(\pi(X, v_t)) u - \langle u, \pi(X, v_t) \rangle \pi(X, v_t))].$$

Consequently,

$$\begin{aligned} \phi'_u(t) &= \frac{2}{\varepsilon} \mathbb{E}[\langle u, \pi(X, v_t) \rangle (u^T \text{diag}(\pi(X, v_t))(v - v^*) - \langle u, \pi(X, v_t) \rangle \langle \pi(X, v_t), v - v^* \rangle)] \\ &\quad - 2\langle u, \nu \rangle u^T \nabla^2 H_\varepsilon(v_t)(v - v^*). \end{aligned}$$

Then, we deduce from equality (5.4) that

$$\phi'_u(t) = \frac{2}{\varepsilon} \mathbb{E}[\langle u, \pi(X, v_t) - \nu \rangle u^T A_\varepsilon(X, v_t)(v - v^*)],$$

with  $A_\varepsilon(x, v) = \text{diag}(\pi(x, v)) - \pi(x, v) \pi(x, v)^T$ . It follows from Cauchy-Schwarz inequality and the upper bound (5.6) that

$$|\langle u, \pi(X, v_t) - \nu \rangle| \leq 2\|u\|,$$

which together with the fact that  $\lambda_{\max}(A_\varepsilon(x, v)) \leq 1$  yields

$$|\phi'_u(t)| \leq \frac{4}{\varepsilon} \|v - v^*\| \|u\|^2.$$

Hence, inserting the above upper bound in (A.2), we obtain that

$$|u^T (L(v) - L(v^*)) u - 2\langle u, \nu \rangle \langle \nabla H_\varepsilon(v), u \rangle| \leq \frac{4}{\varepsilon} \|v - v^*\| \|u\|^2.$$

Therefore, in the sense of partial ordering between positive semi-definite matrices, we have shown that

$$-\frac{4}{\varepsilon} \|v - v^*\| I_J \leq L(v) - L(v^*) - \nabla H_\varepsilon(v) \nu^T - \nu \nabla H_\varepsilon(v)^T \leq \frac{4}{\varepsilon} \|v - v^*\| I_J.$$

Inequality (5.15) thus follows from the decomposition (A.1) since

$$G_\varepsilon(v) - G_\varepsilon(v^*) = L(v) - L(v^*) - \nu \nabla H_\varepsilon(v)^T - \nabla H_\varepsilon(v) \nu^T,$$

which completes the proof of Lemma 5.2.

## A.2 A recursive formula to compute the inverse of $S_n$ for the stochastic Newton algorithm

In this section, we discuss the construction of a recursive formula to compute, from the knowledge of  $S_{n-1}^{-1}$ , the inverse of the matrix  $S_n$  defined by the recursive equation (2.8) that corresponds to the use of the stochastic Newton (SN) algorithm. To this end, let us first recall the following matrix inversion lemma classically referred to as the Sherman-Morrison-Woodbury (SMW) formula [27], also known as Woodbury's formula or Riccati's matrix identity.

**Lemma A.1** (SMW formula). *Suppose that  $A$  and  $C$  are invertible matrices of size  $d \times d$  and  $q \times q$  respectively. Let  $U$  and  $V$  be  $d \times q$  and  $q \times d$  matrices. Then,  $A + UCV$  is invertible iff  $C^{-1} + VA^{-1}U$  is invertible. In that case, we have*

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}. \quad (\text{A.3})$$

A repeated use of the SMW formula allows to prove the following result.

**Proposition A.1.** *Let  $P_J = I_J - \mathbf{v}_J \mathbf{v}_J^T$  be the projection matrix onto  $\langle \mathbf{v}_J \rangle^\perp$ . Suppose that*

$$\mathbb{S}_n = I_J + \sum_{k=1}^n \nabla_v^2 h_\varepsilon(X_k, \hat{V}_{k-1}) = \mathbb{H}_n + \mathbf{v}_J \mathbf{v}_J^T,$$

where

$$\mathbb{H}_n = P_J + \sum_{k=1}^n \nabla_v^2 h_\varepsilon(X_k, \hat{V}_{k-1}),$$

with  $\mathbb{H}_0 = P_J$ . Then, for all  $n \geq 1$ , one has  $\mathbb{S}_n^{-1} = \mathbb{H}_n^- + \mathbf{v}_J \mathbf{v}_J^T$  with  $\mathbb{H}_n^-$  that satisfies the recursive formula

$$\mathbb{H}_n^- = P_J \left( \mathbb{H}_{n-1}^- + \frac{1}{\varepsilon} \text{diag}(\pi_n) \right)^{-1} P_J, \quad (\text{A.4})$$

$$= \mathbb{H}_{n-1}^- - \mathbb{H}_{n-1}^- (\mathbb{H}_{n-1}^- + \varepsilon \text{diag}(\pi_n^{-1}))^{-1} \mathbb{H}_{n-1}^-, \quad (\text{A.5})$$

where  $\pi_n^{-1}$  stands for the vector whose entries are the inverse of those of  $\pi_n = \pi(X_n, \hat{V}_{n-1})$ .

*Proof.* For  $n \geq 1$ , we define

$$\tilde{S}_n = I_J + \sum_{k=1}^n \left( \nabla_v^2 h_\varepsilon(X_k, \hat{V}_{k-1}) + \mathbf{v}_J \mathbf{v}_J^T \right) = \tilde{S}_{n-1} + \Sigma_n = \mathbb{H}_n + (n+1) \mathbf{v}_J \mathbf{v}_J^T, \quad (\text{A.6})$$

where  $\Sigma_n = \nabla_v^2 h_\varepsilon(X_n, \hat{V}_{n-1}) + \mathbf{v}_J \mathbf{v}_J^T$  and  $\tilde{S}_0 = I_J = P_J + \mathbf{v}_J \mathbf{v}_J^T$ . As discussed in Section 5, the eigenvectors of the matrix  $\nabla_v^2 h_\varepsilon(X_k, \hat{V}_{k-1})$  associated to non-zero eigenvalues belong to  $\langle \mathbf{v}_J \rangle^\perp$  for any  $k \geq 1$ , which implies that  $\mathbb{H}_n$  is a matrix of rank  $J - 1$  and that  $\mathbf{v}_J$  is its only eigenvector associated to the eigenvalue  $\lambda_1 = 0$ . Hence, for all  $n \geq 0$ ,  $\mathbb{H}_n$  is also a matrix such that all its eigenvectors associated to non-zero eigenvalues belong to  $\langle \mathbf{v}_J \rangle^\perp$ . Therefore, for any  $n \geq 0$ , the inverse of the matrix  $\tilde{S}_n$  (which is of full rank  $J$ ) satisfies the identity

$$\tilde{S}_n^{-1} = \mathbb{H}_n^- + \frac{1}{n+1} \mathbf{v}_J \mathbf{v}_J^T. \quad (\text{A.7})$$

Moreover, given that  $\mathbb{S}_n = \mathbb{H}_n + \mathbf{v}_J \mathbf{v}_J^T$  one has that

$$\mathbb{S}_n^{-1} = \mathbb{H}_n^{-1} + \mathbf{v}_J \mathbf{v}_J^T = \tilde{S}_n^{-1} + \frac{n}{n+1} \mathbf{v}_J \mathbf{v}_J^T. \quad (\text{A.8})$$

The computation of  $\tilde{S}_n^{-1}$  can be done recursively as follows. By applying the SMW formula (A.3) with  $U = V = I_J$  we obtain that

$$\tilde{S}_n^{-1} = \tilde{S}_{n-1}^{-1} - \tilde{S}_{n-1}^{-1} \left( \Sigma_n^{-1} + \tilde{S}_{n-1}^{-1} \right)^{-1} \tilde{S}_{n-1}^{-1}. \quad (\text{A.9})$$

Now, introducing the notation  $\pi_n = \pi(X_n, \hat{V}_{n-1})$ , we remark that

$$\nabla_v^2 h_\varepsilon(X_n, \hat{V}_{n-1}) = \frac{1}{\varepsilon} (\text{diag}(\pi_n) - \pi_n \pi_n^T)$$

is proportional to a multinomial matrix (up to a minus sign and the multiplicative factor  $\varepsilon^{-1}$ ). Consequently, from the pseudo inverse formula of multinomial matrices [47], we obtain that the inverse of the matrix  $\Sigma_n$  is given by

$$\Sigma_n^{-1} = \varepsilon P_J \text{diag}(\pi_n^{-1}) P_J + \mathbf{v}_J \mathbf{v}_J^T, \quad (\text{A.10})$$

where  $\pi_n^{-1}$  stands for the vector whose entries are the inverse of those of  $\pi_n$ . Now, introducing the notation  $T_{n-1} = \tilde{S}_{n-1}^{-1} + \mathbf{v}_J \mathbf{v}_J^T$  and  $Q_{n-1} = \Sigma_n^{-1} + \tilde{S}_{n-1}^{-1}$ , we deduce from equation (A.10) that

$$Q_{n-1} = T_{n-1} + \varepsilon P_J \text{diag}(\pi_n^{-1}) P_J.$$

Consequently, by the SMW formula (A.3), it follows that

$$Q_{n-1}^{-1} = T_{n-1}^{-1} - T_{n-1}^{-1} P_J \left( P_J T_{n-1}^{-1} P_J + \frac{1}{\varepsilon} \text{diag}(\pi_n) \right)^{-1} P_J T_{n-1}^{-1}.$$

Then, applying once again the SMW formula and equality (A.6), one has that

$$\begin{aligned} T_{n-1}^{-1} &= \tilde{S}_{n-1} - \frac{1}{\mathbf{v}_J^T \tilde{S}_{n-1} \mathbf{v}_J + 1} \tilde{S}_{n-1} \mathbf{v}_J \mathbf{v}_J^T \tilde{S}_{n-1} = \tilde{S}_{n-1} - \frac{n^2}{n+1} \mathbf{v}_J \mathbf{v}_J^T, \\ &= \mathbb{H}_{n-1} + \frac{n}{n+1} \mathbf{v}_J \mathbf{v}_J^T. \end{aligned}$$

Hence, by the fact that  $\mathbb{H}_n$  maps the subspace  $\langle \mathbf{v}_J \rangle^\perp$  onto itself and since  $P_J$  is the projection matrix onto  $\langle \mathbf{v}_J \rangle^\perp$ , we thus obtain that

$$Q_{n-1}^{-1} = T_{n-1}^{-1} - \mathbb{H}_{n-1} \left( \mathbb{H}_{n-1} + \frac{1}{\varepsilon} \text{diag}(\pi_n) \right)^{-1} \mathbb{H}_{n-1}.$$

Consequently, we have shown that

$$Q_{n-1}^{-1} = \left( \Sigma_n^{-1} + \tilde{S}_{n-1}^{-1} \right)^{-1} = \mathbb{H}_{n-1} + \frac{n}{n+1} \mathbf{v}_J \mathbf{v}_J^T - \mathbb{H}_{n-1} \left( \mathbb{H}_{n-1} + \frac{1}{\varepsilon} \text{diag}(\pi_n) \right)^{-1} \mathbb{H}_{n-1}.$$

Therefore, combining the above equality with (A.7), one obtains that

$$\begin{aligned} \tilde{S}_{n-1}^{-1} \left( \Sigma_n^{-1} + \tilde{S}_{n-1}^{-1} \right)^{-1} \tilde{S}_{n-1}^{-1} &= \mathbb{H}_{n-1}^{-1} + \frac{1}{n(n+1)} \mathbf{v}_J \mathbf{v}_J^T \\ &\quad - \mathbb{H}_{n-1}^{-1} \mathbb{H}_{n-1} \left( \mathbb{H}_{n-1} + \frac{1}{\varepsilon} \text{diag}(\pi_n) \right)^{-1} \mathbb{H}_{n-1} \mathbb{H}_{n-1}^{-1}. \end{aligned}$$

Inserting the above equality into (A.9) and using again (A.7), one infers that

$$\begin{aligned}\tilde{S}_n^{-1} &= \tilde{S}_{n-1}^{-1} - \mathbb{H}_{n-1}^{-1} - \frac{1}{n(n+1)} \mathbf{v}_J \mathbf{v}_J^T + \mathbb{H}_{n-1}^{-1} \mathbb{H}_{n-1} \left( \mathbb{H}_{n-1} + \frac{1}{\varepsilon} \text{diag}(\pi_n) \right)^{-1} \mathbb{H}_{n-1}^{-1} \mathbb{H}_{n-1}, \\ &= \frac{1}{n+1} \mathbf{v}_J \mathbf{v}_J^T + \mathbb{H}_{n-1}^{-1} \mathbb{H}_{n-1} \left( \mathbb{H}_{n-1} + \frac{1}{\varepsilon} \text{diag}(\pi_n) \right)^{-1} \mathbb{H}_{n-1}^{-1} \mathbb{H}_{n-1}, \\ &= \frac{1}{n+1} \mathbf{v}_J \mathbf{v}_J^T + P_J \left( \mathbb{H}_{n-1} + \frac{1}{\varepsilon} \text{diag}(\pi_n) \right)^{-1} P_J,\end{aligned}$$

by noticing that  $\mathbb{H}_{n-1}^{-1} \mathbb{H}_{n-1} = \mathbb{H}_{n-1} \mathbb{H}_{n-1}^{-1} = P_J$ . Hereafter, we immediately deduce (A.4) from the above identity together with (A.7). Finally, (A.5) follows from an application of a generalization of the SMW formula [17] to the setting of the Moore-Penrose inverse to handle the situation where the matrix  $A$  in equation (A.3) is not invertible, which completes the proof of Proposition A.1.  $\square$

## B Proofs of auxiliary results related to the KL inequality

*Proof of Proposition 6.1.* The proof consists in a study of  $\tilde{H}_\varepsilon(u)$  when a vector  $u \in \langle \mathbf{v}_J \rangle^\perp$  is either near  $u^*$  or such that  $\|u\| \rightarrow +\infty$ . First, we observe that  $u \mapsto \|\nabla \tilde{H}_\varepsilon(u)\|^2 + \frac{\|\nabla \tilde{H}_\varepsilon(u)\|^2}{\tilde{H}_\varepsilon(u)}$  is a continuous function except at  $u^*$ . Then, since  $\tilde{H}_\varepsilon(u^*) = 0$ , a local approximation of  $\tilde{H}_\varepsilon$  using a Taylor expansion shows that, for all  $h \in \langle \mathbf{v}_J \rangle^\perp$ ,

$$\tilde{H}_\varepsilon(u^* + h) = \frac{1}{2} h^T \nabla^2 \tilde{H}_\varepsilon(u^*) h + o(\|h\|^2),$$

and

$$\nabla \tilde{H}_\varepsilon(u^* + h) = \nabla^2 \tilde{H}_\varepsilon(u^*) h + o(\|h\|).$$

with  $\nabla^2 \tilde{H}_\varepsilon(u^*) = G_*^{-1/2} \nabla^2 H_\varepsilon(v^*) G_*^{-1/2}$ . Since the matrices  $G_*^{-1/2}$  and  $\nabla^2 H_\varepsilon(v^*)$  are of rank  $J-1$  with all eigenvectors corresponding to positive eigenvalues that belong to  $\langle \mathbf{v}_J \rangle^\perp$ , the Courant-Fischer minmax Theorem yields:

$$\begin{aligned}0 < \lambda_{\min}(\nabla^2 \tilde{H}_\varepsilon(u^*)) &= \liminf_{u \rightarrow u^*, u \in \langle \mathbf{v}_J \rangle^\perp} \frac{\|\nabla \tilde{H}_\varepsilon(u)\|^2}{\tilde{H}_\varepsilon(u)} \\ &\leq \limsup_{u \rightarrow u^*, u \in \langle \mathbf{v}_J \rangle^\perp} \frac{\|\nabla \tilde{H}_\varepsilon(u)\|^2}{\tilde{H}_\varepsilon(u)} = \lambda_{\max}(\nabla^2 \tilde{H}_\varepsilon(u^*)),\end{aligned}\tag{B.1}$$

where  $\lambda_{\min}(\nabla^2 \tilde{H}_\varepsilon(u^*))$  denotes the second smallest eigenvalue of the matrix  $\nabla^2 \tilde{H}_\varepsilon(u^*)$ , and the notation  $u \rightarrow u^*$  corresponds to the convergence of  $u \in \langle \mathbf{v}_J \rangle^\perp$ . Inequality (B.1) thus implies that the continuous function  $u \mapsto \|\nabla \tilde{H}_\varepsilon(u)\|^2 + \frac{\|\nabla \tilde{H}_\varepsilon(u)\|^2}{\tilde{H}_\varepsilon(u)}$  is upper and lower bounded by positive constants in a neighborhood of  $u^*$ . Finally, we observe that  $\|\nabla \tilde{H}_\varepsilon\|$  is bounded thanks to inequality (5.7), and it can be checked that the function  $\tilde{H}_\varepsilon$  is coercive (over the finite dimensional vector space  $\langle \mathbf{v}_J \rangle^\perp$ ) since it has a unique minimizer at  $u^*$ . These facts together with the boundedness of  $u \mapsto \|\nabla \tilde{H}_\varepsilon(u)\|^2 + \frac{\|\nabla \tilde{H}_\varepsilon(u)\|^2}{\tilde{H}_\varepsilon(u)}$  in a neighborhood of  $u^*$  implies that inequality (6.43) holds.

Now, let us show that the constant  $m$  appearing in inequality (6.43) can be made more explicit thanks to Lemma 5.1. Indeed, since  $\nabla \tilde{H}_\varepsilon(u) = G_*^{-1/2} \nabla H_\varepsilon(G_*^{-1/2}u)$ , we immediately obtain from inequality (5.14) that

$$\begin{aligned} \langle \nabla \tilde{H}_\varepsilon(u), u - u^* \rangle &\geq \frac{1 - \exp(-\delta(u))}{\delta(u)} (u - u^*)^T G_*^{-1/2} \nabla^2 H_\varepsilon(v^*) G_*^{-1/2} (u - u^*), \\ &\geq \frac{1 - \exp(-\delta(u))}{\delta(u)} \|u - u^*\|^2, \end{aligned} \quad (\text{B.2})$$

where  $\delta(u) = \frac{\sqrt{2}}{\varepsilon} \lambda_{\max}(G_*^{-1/2}) \|u - u^*\|$  and the second inequality above follows from the fact that

$$\lambda_{\min}^{\langle v_J \rangle^\perp} \left( G_*^{-1/2} \nabla^2 H_\varepsilon(v^*) G_*^{-1/2} \right) \geq 1,$$

by inequality (3.4). Note that inequality (B.2) corresponds to a local strong convex property of the function  $\tilde{H}_\varepsilon$  in the neighborhood of  $u^*$ . Then, by the Cauchy-Schwarz inequality, one has that

$$\langle \nabla \tilde{H}_\varepsilon(u), u - u^* \rangle \leq \|\nabla \tilde{H}_\varepsilon(u)\| \|u - u^*\|.$$

Thus, we obtain by inequality (B.2) that, for any  $u \in \langle v_J \rangle^\perp$ ,

$$\|\nabla \tilde{H}_\varepsilon(u)\| \geq \frac{1 - \exp(-\delta(u))}{\delta(u)} \|u - u^*\|. \quad (\text{B.3})$$

We then consider two cases. If  $\|u - u^*\| \leq \frac{\varepsilon}{\lambda_{\max}(G_*^{-1/2})}$ , then, using that the function  $\delta \mapsto \frac{1}{\delta} (1 - \exp(-\delta))$  is decreasing, it follows from (6.42) and (B.3) that,

$$\|\nabla \tilde{H}_\varepsilon(u)\|^2 + \frac{\|\nabla \tilde{H}_\varepsilon(u)\|^2}{\tilde{H}_\varepsilon(u)} \geq \frac{\|\nabla \tilde{H}_\varepsilon(u)\|^2}{\tilde{H}_\varepsilon(u)} \geq \frac{2\varepsilon}{\lambda_{\max}(G_*^-)} \left(1 - \exp(-\sqrt{2})\right)^2. \quad (\text{B.4})$$

To the contrary, if  $\|u - u^*\| \geq \frac{\varepsilon}{\lambda_{\max}(G_*^{-1/2})}$ , then

$$\|\nabla \tilde{H}_\varepsilon(u)\|^2 + \frac{\|\nabla \tilde{H}_\varepsilon(u)\|^2}{\tilde{H}_\varepsilon(u)} \geq \|\nabla \tilde{H}_\varepsilon(u)\|^2 \geq \frac{\varepsilon^2}{2\lambda_{\max}(G_*^-)} \left(1 - \exp(-\sqrt{2})\right)^2. \quad (\text{B.5})$$

using the inequality  $1 - \exp(-\delta(u)) \geq 1 - \exp(-\sqrt{2})$  that holds for  $\delta(u) \geq \sqrt{2}$ . Consequently, since  $\left(1 - \exp(-\sqrt{2})\right)^2 > 1/2$  and combining inequalities (B.4) and (B.5), it follows that the constant  $m$  appearing in inequality (6.43) can be chosen as  $m = m_\varepsilon$  with  $m_\varepsilon$  defined by (6.44). This concludes the proof of Proposition 6.1.  $\square$

We then show the proof of the one-step evolution of the SGN algorithm.

*Proof of Proposition 6.3.* First, by using the arguments from the proof of (i) of Theorem 3.1, we remark that the matrix  $S_n$  defined by (2.9) satisfies:

$$\lambda_{\min}(S_n) \geq 1 + \gamma \min(\nu) \left( \sum_{m=1}^{p_n} m^{-\beta} \right) \quad \text{and} \quad \lambda_{\max}(S_n) \leq 1 + 4n + \gamma \max(\nu) \sum_{m=1}^{p_n+1} m^{-\beta},$$

where  $p_n \geq 1$  denotes the largest integer such that  $p_n J \leq n$ . Consequently, using the fact that  $\frac{1}{1-\beta}(p_n^{1-\beta} - 1) \leq \sum_{m=1}^{p_n} m^{-\beta} \leq \frac{1}{1-\beta} p_n^{1-\beta}$  the above inequalities imply that:

$$\begin{aligned} \lambda_{\max}(S_n^{-1}) &\leq \frac{1-\beta}{1-\beta+\gamma \min(\nu) \left( (n_J-1)^{1-\beta} - 1 \right)} \\ &\leq \frac{1}{1-2\gamma \min(\nu) + \gamma \min(\nu) (n_J-1)^{1-\beta}}, \end{aligned} \quad (\text{B.6})$$

where  $n_J = n/J$ , and

$$\lambda_{\min}(S_n^{-1}) \geq \frac{1}{1-\beta+4(1-\beta)n+\gamma \max(\nu) (n/J)^{1-\beta}} \geq \frac{1}{1+(4+2\gamma \max(\nu))n}. \quad (\text{B.7})$$

Step 1: Taylor expansion. First, we introduce the notation  $\tilde{U}_n = G_*^{1/2} \hat{V}_n$ , and in the proof we repeatedly use the property that the eigenvectors of  $G_*$  associated to non-zero eigenvalues belong to  $\langle \mathbf{v}_J \rangle^\perp$ . Using equation (2.6) and the fact that  $P_J G_* = G_*$ , a second order Taylor expansion yields

$$\begin{aligned} \Phi(\tilde{U}_{n+1}) &= \Phi\left(\tilde{U}_n - n^\alpha G_*^{1/2} S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n)\right) \\ &= \Phi(\tilde{U}_n) - n^\alpha \left\langle \nabla \Phi(\tilde{U}_n), G_*^{1/2} S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \right\rangle \\ &\quad + \frac{n^{2\alpha}}{2} \nabla^2 \Phi(\xi_{n+1}) \left( G_*^{1/2} S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \right)^{\otimes 2}, \end{aligned} \quad (\text{B.8})$$

where  $\xi_{n+1}$  is such that  $\xi_{n+1} = \tilde{U}_n + t_{n+1}(\tilde{U}_{n+1} - \tilde{U}_n)$  with  $t_{n+1} \in (0, 1)$ . Now, applying inequalities (5.6) and (6.48), the second order term in equation (B.8) can be bounded as follows:

$$\begin{aligned} \left\| \nabla^2 \Phi(\xi_{n+1}) \left( G_*^{1/2} S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \right)^{\otimes 2} \right\| &\leq \delta_\varepsilon \lambda_{\max}(G_*^-) (1 + \Phi(\xi_{n+1})) \times \\ &\quad \left\| G_*^{1/2} S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \right\|^2 \\ &\leq 4\delta_\varepsilon \lambda_{\max}(G_*^-) \lambda_{\max}^2(G_*^{1/2} S_n^{-1}) (1 + \Phi(\xi_{n+1})), \end{aligned}$$

which yields the inequality:

$$\begin{aligned} \Phi(\tilde{U}_{n+1}) &\leq \Phi(\tilde{U}_n) - n^\alpha \left\langle \nabla \Phi(\tilde{U}_n), G_*^{1/2} S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n) \right\rangle \\ &\quad + 2\delta_\varepsilon \frac{\lambda_{\max}(G_*)}{\lambda_{\min}(G_*)} n^{2\alpha} \lambda_{\max}^2(S_n^{-1}) (1 + \Phi(\xi_{n+1})). \end{aligned} \quad (\text{B.9})$$

Step 2: An auxiliary inequality. We now establish a technical bound to relate  $\Phi(\xi_{n+1})$  and  $\Phi(\tilde{U}_n)$ . To this end, by a first order Taylor expansion of the function  $s \mapsto \tilde{H}_\varepsilon(\tilde{U}_n + s(\tilde{U}_{n+1} - \tilde{U}_n))$  where  $\tilde{U}_n^s = \tilde{U}_n + s(\tilde{U}_{n+1} - \tilde{U}_n)$ , one has that:

$$\tilde{H}_\varepsilon(\tilde{U}_{n+1}) = \tilde{H}_\varepsilon(\tilde{U}_n) + \int_0^1 \nabla \tilde{H}_\varepsilon(\tilde{U}_n^s) (\tilde{U}_{n+1} - \tilde{U}_n) ds,$$

and thus, combining the Cauchy-Schwarz inequality with the upper bounds (5.6) and (6.40), we obtain that:

$$\begin{aligned}
\tilde{H}_\varepsilon(\xi_n) &\leq \tilde{H}_\varepsilon(\tilde{U}_n) + \sup_{t \in [0,1]} \|\nabla \tilde{H}_\varepsilon(\tilde{U}_n^T)\| \|\tilde{U}_{n+1} - \tilde{U}_n\| \\
&\leq \tilde{H}_\varepsilon(\tilde{U}_n) + \sup_{t \in [0,1]} \|\nabla \tilde{H}_\varepsilon(\tilde{U}_n^T)\| \|n^\alpha G_*^{1/2} S_n^{-1} \nabla_v h_\varepsilon(X_{n+1}, \hat{V}_n)\| \\
&\leq \tilde{H}_\varepsilon(\tilde{U}_n) + 4\lambda_{\max}(G_*^{-1/2})\lambda_{\max}(G_*^{1/2})n^\alpha \lambda_{\max}(S_n^{-1}).
\end{aligned} \tag{B.10}$$

Note that, under the condition  $\alpha + \beta < 1/2$ , it follows from inequality (B.6) that  $n^\alpha \lambda_{\max}(S_n^{-1}) \leq c_0$  for some constant  $c_0 \geq 1$  for all  $n \geq J$ . Hence, inserting the upper bound (B.10) into the definition of  $\Phi$  and using inequality (6.46), we obtain that

$$\begin{aligned}
\Phi(\xi_{n+1}) &\leq \left( \tilde{H}_\varepsilon(\tilde{U}_n) + 4 \frac{\lambda_{\max}(G_*^{1/2})}{\lambda_{\min}(G_*^{1/2})} n^\alpha \lambda_{\max}(S_n^{-1}) \right) \times \\
&\quad \exp \left( \tilde{H}_\varepsilon(\tilde{U}_n) + 4 \frac{\lambda_{\max}(G_*^{1/2})}{\lambda_{\min}(G_*^{1/2})} n^\alpha \lambda_{\max}(S_n^{-1}) \right) \\
&\leq \exp \left( 4 \frac{\lambda_{\max}(G_*^{1/2})}{\lambda_{\min}(G_*^{1/2})} n^\alpha \lambda_{\max}(S_n^{-1}) \right) \times \\
&\quad \left( \tilde{H}_\varepsilon(\tilde{U}_n) + 4 \frac{\lambda_{\max}(G_*^{1/2})}{\lambda_{\min}(G_*^{1/2})} n^\alpha \lambda_{\max}(S_n^{-1}) \right) \exp(\tilde{H}_\varepsilon(\tilde{U}_n)) \\
&\leq \exp \left( 4c_0 \frac{\lambda_{\max}(G_*^{1/2})}{\lambda_{\min}(G_*^{1/2})} \right) \left( \Phi(\tilde{U}_n) + 4c_0 \frac{\lambda_{\max}(G_*^{1/2})}{\lambda_{\min}(G_*^{1/2})} (1 + \Phi(\tilde{U}_n)) \right) \\
&\leq \tilde{c}_0 (1 + \Phi(\tilde{U}_n)),
\end{aligned} \tag{B.11}$$

where  $\tilde{c}_0 = \max \left( 1, 4c_0 \frac{\lambda_{\max}(G_*^{1/2})}{\lambda_{\min}(G_*^{1/2})} \right) \exp \left( 4c_0 \frac{\lambda_{\max}(G_*^{1/2})}{\lambda_{\min}(G_*^{1/2})} \right)$ .

Step 3: Derivation of a recursive inequality. Inserting inequality (B.11) into (B.9), and taking the conditional expectation with respect to  $\mathcal{F}_n$ , we obtain that:

$$\begin{aligned}
\mathbb{E} [\Phi(\tilde{U}_{n+1}) | \mathcal{F}_n] &\leq \Phi(\tilde{U}_n) - n^\alpha \left\langle \nabla \Phi(\tilde{U}_n), G_*^{1/2} S_n^{-1} G_*^{1/2} \nabla \tilde{H}_\varepsilon(\tilde{U}_n) \right\rangle \\
&\quad + 2\delta_\varepsilon \frac{\lambda_{\max}(G_*)}{\lambda_{\min}(G_*)} n^{2\alpha} \lambda_{\max}^2(S_n^{-1}) \left( 1 + \tilde{c}_0 (1 + \Phi(\tilde{U}_n)) \right),
\end{aligned}$$

where we used the property that  $\nabla H_\varepsilon(\hat{V}_n) = G_*^{1/2} \nabla \tilde{H}_\varepsilon(\tilde{U}_n)$ .

Consequently, using inequality (6.47) and introducing  $c_\varepsilon = 2\tilde{c}_0 \delta_\varepsilon \frac{\lambda_{\max}(G_*)}{\lambda_{\min}(G_*)}$ , we have:

$$\begin{aligned}
\mathbb{E} [\Phi(\tilde{U}_{n+1}) | \mathcal{F}_n] &\leq (1 - m_\varepsilon n^\alpha \lambda_{\min}(G_*) \lambda_{\min}(S_n^{-1}) + c_\varepsilon n^{2\alpha} \lambda_{\max}^2(S_n^{-1})) \Phi(\tilde{U}_n) \\
&\quad + c_\varepsilon n^{2\alpha} \lambda_{\max}^2(S_n^{-1})
\end{aligned} \tag{B.12}$$

Now, thanks to inequalities (B.6) and (B.7), and the condition  $0 < \alpha + \beta < 1/2$ , it follows that there exists an integer  $n_0$  such that, for all  $n \geq n_0$ ,

$$m_\varepsilon n^\alpha \lambda_{\min}(G_*) \lambda_{\min}(S_n^{-1}) \geq 2c_\varepsilon n^{2\alpha} \lambda_{\max}^2(S_n^{-1}). \tag{B.13}$$



Hence, by combining inequalities (B.12) and (B.13) we obtain inequality (6.50) which concludes the proof of *i*.  $\square$

## Acknowledgments

The authors gratefully acknowledge financial support from the Agence Nationale de la Recherche (MaSDOL grant ANR-19-CE23-0017). J. Bigot and S. Gadat are members of the Institut Universitaire de France (IUF), and part of this work has been carried out with financial support from the IUF.

## References

- [1] ALTSCHULER, J., WEED, J., AND RIGOLLET, P. Near-linear time approximation algorithms for optimal transport via sinkhorn iteration. In Advances in Neural Information Processing Systems 30. 2017, pp. 1964–1974.
- [2] BACH, F. R. Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic regression. Journal of Machine Learning Research 15, 1 (2014), 595–627.
- [3] BENAÏM, M., AND HIRSCH, M. Asymptotic pseudotrajectories and chain recurrent flows, with applications. Journal of Dynamics and Differential Equations 8, 1 (1996), 141–176.
- [4] BENAMOU, J.-D., CARLIER, G., CUTURI, M., NENNA, L., AND PEYRÉ, G. Iterative Bregman projections for regularized transportation problems. SIAM Journal on Scientific Computing 37, 2 (2015), A1111–A1138.
- [5] BERCU, B., AND BIGOT, J. Asymptotic distribution and convergence rates of stochastic algorithms for entropic optimal transportation between probability measures. Annals of Statistics 49, 2 (2021), 968–987.
- [6] BERCU, B., GODICHON, A., AND PORTIER, B. An efficient stochastic newton algorithm for parameter estimation in logistic regressions. SIAM Journal on Control and Optimization 58, 1 (2020), 348–367.
- [7] BIGOT, J., CAZELLES, E., AND PAPADAKIS, N. Data-driven regularization of Wasserstein barycenters with an application to multivariate density registration. Information and Inference 8, 4 (2019), 719–755.
- [8] BIGOT, J., GOUET, R., KLEIN, T., LÓPEZ, A., ET AL. Geodesic PCA in the Wasserstein space by convex PCA. Annales de l’Institut Henri Poincaré, Probabilités et Statistiques 53, 1 (2017), 1–26.
- [9] BIGOT, JÉRÉMIE. Statistical data analysis in the wasserstein space. ESAIM: ProcS 68 (2020), 1–19.
- [10] BOLTE, J., DANIILIDIS, A., AND LEWIS, A. The Łojasiewicz inequality for non-smooth subanalytic functions with applications to subgradient dynamical systems. SIAM J. Optim. 17, 4 (2006), 1205–1223.

- [11] BOLTE, J., NGUYEN, P., PEYPOUQUET, J., AND SUTER, B. W. From error bounds to the complexity of first-order descent methods for convex functions. Math. Program. (A), 165 (2017), 471–507.
- [12] BONNEEL, N., RABIN, J., PEYRÉ, G., AND PFISTER, H. Sliced and radon Wasserstein barycenters of measures. Journal of Mathematical Imaging and Vision 51, 1 (2015), 22–45.
- [13] CAZELLES, E., SEGUY, V., BIGOT, J., CUTURI, M., AND PAPADAKIS, N. Log-PCA versus Geodesic PCA of histograms in the Wasserstein space. SIAM Journal on Scientific Computing 40, 2 (2018), B429–B456.
- [14] CÉNAC, P., GODICHON-BAGGIONI, A., AND PORTIER, B. An efficient averaged stochastic gauss-newton algorithm for estimating parameters of non linear regressions models. Preprint - arXiv2006.12920, 2020.
- [15] CUTURI, M. Sinkhorn distances: Lightspeed computation of optimal transport. In Advances in Neural Information Processing Systems 26, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2013, pp. 2292–2300.
- [16] CUTURI, M., AND PEYRÉ, G. Semidual regularized optimal transport. SIAM Review 60, 4 (2018), 941–965.
- [17] DENG, C. Y. A generalization of the Sherman-Morrison-Woodbury formula. Applied Mathematics Letters 24, 9 (2011), 1561 – 1564.
- [18] DUFLO, M. Random iterative models, vol. 34 of Applications of Mathematics, New York. Springer-Verlag, Berlin, 1997.
- [19] FERRADANS, S., PAPADAKIS, N., PEYRÉ, G., AND AUJOL, J.-F. Regularized discrete optimal transport. SIAM Journal on Imaging Sciences 7, 3 (2014), 1853–1882.
- [20] FLAMARY, R., CUTURI, M., COURTY, N., AND RAKOTOMAMONJY, A. Wasserstein discriminant analysis. Machine Learning 107, 12 (Dec 2018), 1923–1945.
- [21] FROGNER, C., ZHANG, C., MOBAHI, H., ARAYA-POLO, M., AND POGGIO, T. Learning with a wasserstein loss. In Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2 (Cambridge, MA, USA, 2015), NIPS’15, MIT Press, pp. 2053–2061.
- [22] GADAT, S., AND PANLOUP, F. Optimal non-asymptotic bound of the Ruppert-Polyak averaging without strong convexity. Preprint - arXiv:1709.03342, May 2022.
- [23] GENEVAY, A., CUTURI, M., PEYRÉ, G., AND BACH, F. Stochastic optimization for large-scale optimal transport. In Advances in Neural Information Processing Systems 29. 2016, pp. 3440–3448.
- [24] GENEVAY, A., PEYRE, G., AND CUTURI, M. Learning generative models with sinkhorn divergences. In Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics (Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018), A. Storkey and F. Perez-Cruz, Eds., vol. 84 of Proceedings of Machine Learning Research, PMLR, pp. 1608–1617.

- [25] GORDALIZA, P., DEL BARRIO, E., GAMBOA, F., AND LOUBES, J.-M. Obtaining fairness using optimal transport theory. In Proceedings of the 36th International Conference on Machine Learning (2019), pp. 2357–2365.
- [26] GRAMFORT, A., PEYRÉ, G., AND CUTURI, M. Fast optimal transport averaging of neuroimaging data. In International Conference on Information Processing in Medical Imaging (2015), Springer, pp. 261–272.
- [27] HAGER, W. W. Updating the inverse of a matrix. SIAM Review 31, 2 (1989), 221–239.
- [28] KINGMA, D. P., AND BA, J. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings (2015), Y. Bengio and Y. LeCun, Eds.
- [29] KITAGAWA, J., MÉRIGOT, Q., AND B., T. Convergence of a newton algorithm for semi-discrete optimal transport. Journal of the European Math Society 21, 9 (2019), 2603–2651.
- [30] KLATT, M., TAMELING, C., AND MUNK, A. Empirical regularized optimal transport: Statistical theory and applications. SIAM J. Math. Data Sci. 2, 2 (2020), 419–443.
- [31] KURDYKA, K. On gradients of functions definable in o-minimal structures. Ann. Inst. Fourier (Grenoble) 48, 3 (1998), 769–783.
- [32] LOJASIEWICZ, S. Une propriété topologique des sous-ensembles analytiques réels. Editions du centre National de la Recherche Scientifique, Paris, Les Équations aux Dérivées Partielles (1963), 87–89.
- [33] MÉRIGOT, Q. A multiscale approach to optimal transport. Computer Graphics Forum 30, 5 (2011), 1583–1592.
- [34] MÉRIGOT, Q., MEYRON, J., AND THIBERT, B. An algorithm for optimal transport between a simplex soup and a point cloud. SIAM Journal on Imaging Sciences 11, 2 (2018), 1363–1389.
- [35] PANARETOS, V. M., AND ZEMEL, Y. Amplitude and phase variation of point processes. Annals of Statistics 44, 2 (2016), 771–812.
- [36] PANARETOS, V. M., AND ZEMEL, Y. Statistical aspects of wasserstein distances. Annual Reviews of Statistics and its Applications 6 (2018), 405–431.
- [37] PELLETIER, M. Asymptotic almost sure efficiency of averaged stochastic algorithms. SIAM J. Control and Optimization 39 (08 2000), 49–72.
- [38] PEYRÉ, G., AND CUTURI, M. Computational optimal transport. Foundations and Trends in Machine Learning 11, 5-6 (2019), 355–607.
- [39] RABIN, J., AND PAPADAKIS, N. Convex color image segmentation with optimal transport distances. In International Conference on Scale Space and Variational Methods in Computer Vision (2015), Springer, pp. 256–269.

- [40] RIGOLLET, P., AND WEED, J. Entropic optimal transport is maximum-likelihood deconvolution. Comptes Rendus Mathématique 356, 11 (2018), 1228 – 1235.
- [41] ROBBINS, H., AND SIEGMUND, D. A convergence theorem for non negative almost supermartingales and some applications. In Optimizing methods in statistics. Elsevier, 1971, pp. 233–257.
- [42] ROLET, A., CUTURI, M., AND PEYRÉ, G. Fast dictionary learning with a smoothed Wasserstein loss. In Proc. International Conference on Artificial Intelligence and Statistics (AISTATS) (2016).
- [43] SANJABI, M., BA, J., RAZAVIYAYN, M., AND LEE, J. D. On the convergence and robustness of training gans with regularized optimal transport. In Advances in Neural Information Processing Systems (2018), S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31, Curran Associates, Inc.
- [44] SEGUY, V., AND CUTURI, M. Principal geodesic analysis for probability measures under the optimal transport metric. In Advances in Neural Information Processing Systems 28. 2015, pp. 3294–3302.
- [45] SOLOMON, J., DE GOES, F., PEYRÉ, G., CUTURI, M., BUTSCHER, A., NGUYEN, A., DU, T., AND GUIBAS, L. Convolutional wasserstein distances: Efficient optimal transportation on geometric domains. ACM Trans. Graph. 34, 4 (2015), 66:1–66:11.
- [46] SOMMERFELD, M., AND MUNK, A. Inference for empirical Wasserstein distances on finite spaces. Journal of the Royal Statistical Society: Series B (Statistical Methodology) (2016).
- [47] STEERNEMAN, A., AND VAN PERLO -TEN KLEIJ, F. Properties of the matrix  $a$ -xy. Linear Algebra and Its Applications 410 (2005), 70–86.
- [48] VILLANI, C. Optimal transport: old and new, vol. 338. Springer Science & Business Media, 2008.
- [49] ZEMEL, Y., AND PANARETOS, V. M. Fréchet means and procrustes analysis in wasserstein space. Bernoulli 25, 2 (05 2019), 932–976.
- [50] ZHANG, L.-X. Central limit theorems of a recursive stochastic algorithm with applications to adaptive designs. Ann. Appl. Probab. 26, 6 (2016), 3630–3658.