

---

# Stochastic Optimization in Semi-Discrete Optimal Transport: Convergence Analysis and Minimax Rate

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We investigate the semi-discrete Optimal Transport (OT) problem, where a contin-  
2 uous source measure  $\mu$  is transported to a discrete target measure  $\nu$ , with particular  
3 attention to the OT map approximation. In this setting, Stochastic Gradient Descent  
4 (SGD) based solvers have demonstrated strong empirical performance in recent  
5 machine learning applications, yet their theoretical guarantee to approximate the  
6 OT map is an open question. In this work, we answer it positively by providing  
7 both computational and statistical convergence guarantees of SGD. Specifically,  
8 we show that SGD methods can estimate the OT map with a minimax convergence  
9 rate of  $\mathcal{O}(1/\sqrt{n})$ , where  $n$  is the number of samples drawn from  $\mu$ . To establish  
10 this result, we study the averaged projected SGD algorithm, and identify a suitable  
11 projection set that contains a minimizer of the objective, even when the source  
12 measure is not compactly supported. Our analysis holds under mild assumptions  
13 on the source measure and applies to MTW cost functions, which include  $\|\cdot\|^p$  for  
14  $p \in (1, \infty)$ . We finally provide numerical evidence for our theoretical results.

## 15 1 Introduction

16 Optimal Transport (OT) has become a central tool in machine learning for comparing and manip-  
17 ulating probability measures. A particularly important variant is the *semi-discrete* setting, where a  
18 continuous source distribution  $\mu$  is transported to a discrete target measure  $\nu$ . This formulation arises  
19 naturally in a wide range of applications, including image processing [15, 21], statistics [10, 17], and  
20 generative modeling [3, 8, 23]. Its hybrid structure bridges the gap between fully continuous and  
21 fully discrete formulations, allowing for expressive modeling while remaining more amenable to  
22 scalable numerical methods.

23 Despite its practical appeal, solving semi-discrete OT efficiently and reliably at scale remains  
24 challenging. Existing methods with convergence guarantees often require *full knowledge* of the  
25 source density and are typically confined to *low-dimensional settings*. For instance, Newton-type  
26 methods [25, 22, 20] and combinatorial [1] methods have been developed and are provided with  
27 convergence rates guarantees. However, these techniques become impractical in high-dimensional  
28 settings, since they need full knowledge of the source measure, and employ constructions that suffer  
29 from the curse of dimensionality such as meshes representation of the source measure.

30 In high dimensional and/or when the source measure can only be accessed through samples, SGD  
31 and its variants have become a popular choice for solving semi-discrete OT, especially in applications  
32 such as generative modeling [3, 8, 23]. These methods solve the semi-dual formulation of the OT  
33 problem using only i.i.d. samples from  $\mu$  and scale well to large data, not requiring to store samples.  
34 Yet, a *fundamental theoretical gap remains*: while SGD methods are widely used in practice, they  
35 lack convergence guarantees for approximating the OT map.

Providing convergence rates for SGD to approximate the semi-discrete OT map is both a computational and a statistical problem, since using more gradient steps in SGD for the semi-discrete setting is equivalent to using more samples. From a statistical point of view, recent work by Pooladian et al. [29] offers promising evidence for the convergence of SGD: they show that semi-discrete OT escapes the curse of dimensionality, unlike the continuous setting [14], and that a convergence rate of  $\mathcal{O}(1/\sqrt{n})$  with  $n$  samples is achievable when the cost to move mass is the quadratic cost  $\frac{1}{2}\|\cdot\|^2$  and the source measure is compactly supported. Moreover, this rate is minimax optimal for estimating the OT map. While their approach requires solving a discrete OT problem, i.e., first sampling points from  $\mu$  and then solving the corresponding empirical problem, their key result motivates studying SGD as a way to accurately estimate the OT map in an online setting, where the estimator is refined as more samples become available, without needing to store them, as is often required in practice, for instance in generative modeling tasks

Taken together, these observations highlight two major open questions that are answered positively in this work:

- (i) Can we establish convergence guarantees for SGD-based algorithms in the semi-discrete OT setting for OT map estimation, especially when only samples from  $\mu$  are available?
- (ii) Can we obtain statistical guarantees for the estimation of OT quantities (e.g., cost, potential, map) beyond the compact and quadratic setting of Pooladian et al. [29]?

**Contributions.** We establish convergence guarantees for SGD-based algorithms applied to the non-regularized semi-dual formulation of optimal transport. We focus our analysis on the convergence of the averaged Projected Stochastic Gradient Descent (PSGD) algorithm, relying on a key result: the existence of a compact projection set  $\mathcal{C}$  (Lemma 3.1) that contains a minimizer of the semi-dual objective, even when the source measure  $\mu$  has unbounded support. To the best of our knowledge, this is the first time such a projection set has been identified in the semi-discrete OT setting without assuming boundedness of  $\mu$ . This projection set allows us to derive key properties of the semi-dual functional, including a global weak form of strong convexity on  $\mathcal{C}$  in all our settings (Lemma 4.2), and to extend second-order regularity results that were previously known only in the compact setting (Prop 4.1), such as local strong convexity near the optimum.

The convergence rates for OT quantities, obtained when  $n$  samples are drawn from  $\mu$  (equivalently, when PSGD is run for  $n$  iterations), answer the two questions posed in the introduction. These results are summarized in Table 1.

Table 1: Summary of the convergence rate achieved by PSGD for the estimation of OT quantities

Costs	OT cost	OT potential	OT map	Non compact
MTW costs	$\mathcal{O}(1/n)$ [Cor 4.6]	$\mathcal{O}(1/n)$ [Th 4.5]	$\mathcal{O}(1/\sqrt{n})$ [Cor 5.2]	No
Quadratic	$\mathcal{O}(1/n)$ [Cor 4.6]	$\mathcal{O}(1/n)$ [Th 4.5]	$\mathcal{O}(1/\sqrt{n})$ [Cor 5.2]	Yes
All other costs	$\mathcal{O}(1/\sqrt{n})$ [Th 3.2]	No guarantees	No guarantees <sup>1</sup>	Yes

Moreover, we show our convergence rate to the OT map is minimax optimal (Theorem 5.3), giving the first statistical guarantees in semi-discrete OT for MTW and quadratic costs on unbounded support.

## 2 Background

### 2.1 Optimal Transport

Considering a source and target probability measures  $\mu \in \mathcal{P}(\mathbb{R}^d)$  and  $\nu \in \mathcal{P}(\mathbb{R}^d)$ , with a cost  $c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^+$  to transfer mass, the OT problem is defined as:

$$\text{OT}_c(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} c(x, y) d\pi(x, y), \quad (1)$$

where  $\Pi(\mu, \nu) := \{\pi \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d); \forall(x, y), \int_{\mathbb{R}^d} \pi(x, y) dy = \mu(x), \int_{\mathbb{R}^d} \pi(x, y) dx = \nu(y)\}$  is the set of joint probability measures on  $\mathbb{R}^d \times \mathbb{R}^d$  with marginals  $\mu$  and  $\nu$ .

<sup>1</sup>Not even necessarily defined.

In this article, we assume that  $\mu$  is continuous, while the target measure takes the form  $\nu = \sum_{i=1}^M w_i \delta_{\{y_i\}}$ , where  $\mathbf{w} = (w_1, \dots, w_M)$  are its probability weights and  $(y_1, \dots, y_M)$  its support. We are mainly interested in OT problems where (1) has a unique deterministic solution, which we refer to as the OT map. This holds under mild assumptions in our semi-discrete setting, as guaranteed by the generalized Brenier's Theorem.

**(Generalized) Brenier's Theorem** ([31], Th. 10.28). Suppose that  $\mu$  is an absolutely continuous probability measure with a second-order moment, and  $\nu$  is discrete. Then, if the cost  $c$  satisfies the MTW properties, as when  $c = \|\cdot\|^p$  with  $p \in (1, \infty)$ , (1) has, up to negligible sets, a unique solution of the form  $\gamma(dx, dy) = \mu(dx) \delta_{T_{\mu, \nu}(x)}(dy)$ ,  $\forall (x, y)$ , where:

$$T_{\mu, \nu}(x) = x - \nabla f^*(x).$$

$T_{\mu, \nu}$  is referred to as the **OT map**. Moreover,  $f^*$  is a solution of the dual problem of (1), given by:

$$\text{OT}_c(\mu, \nu) = \max_{f(x) + g(y) \leq c(x, y)} \int_{\mathbb{R}^d} f(x) d\mu(x) + \int_{\mathbb{R}^d} g(y) d\nu(y). \quad (2)$$

**The semi-dual problem.** The semi-dual formulation of (2) is particularly useful in the semi-discrete setting and can be expressed as a concave finite-dimensional problem:

$$\text{OT}_c(\mu, \nu) = \max_{\mathbf{g} \in \mathbb{R}^M} \left( H(\mathbf{g}) := \int_{\mathbb{R}^d} \mathbf{g}^c(x) d\mu(x) + \sum_{j=1}^M w_j g_j \right), \quad (3)$$

where  $\mathbf{g} = (g_1, \dots, g_M)$ , and for all  $x \in \mathbb{R}^d$ , the (vectorial)  $c$ -transform is defined as  $\mathbf{g}^c(x) := \min_{i \in \llbracket 1, M \rrbracket} \{c(x, y_i) - g_i\}$ . For any  $\mathbf{g}$ , the  $c$ -transform also defines the Laguerre cells  $\mathbb{L}_j^c(\mathbf{g})$  for  $j \in \llbracket 1, M \rrbracket$  as

$$\mathbb{L}_j^c(\mathbf{g}) := \{x \in \mathcal{X} \mid \mathbf{g}^c(x) = c(x, y_j) - g_j\}.$$

Using this formulation and under Brenier's generalized theorem, the OT map can be described as  $T_{\mu, \nu}(x) = x - \nabla(\mathbf{g}^*)^c(x) = x - y_i$  for  $x$  inside  $\mathbb{L}_i^c(\mathbf{g}^*)$ , where  $\mathbf{g}^*$ , that we refer to as the discrete **OT potential** solves (3).

**MTW costs.** For costs satisfying the Ma-Trudinger-Wang properties, referred to as MTW costs, such as  $\|\cdot\|^p$ ,  $p \in (1, \infty)$  or Bregman divergences  $\phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle$  for strictly convex  $\phi$ , we observe improved properties for  $H$ . Specifically, the differential is defined almost everywhere, except on a  $\mu$ -negligible set, and [20] proved that  $H$  is even locally smooth and strongly convex on the orthogonal complement of  $\mathbf{1}$ . Moreover, all these costs satisfy Brenier's generalized theorem. For a detailed definition of such costs in the semi-discrete setting, see [20] or Appendix A.

## 2.2 Stochastic approach for the semi-dual problem.

In the semi-discrete setting, the semi-dual formulation is particularly appealing because, even when  $\mu$  is a continuous measure, it reduces to a finite-dimensional problem. Efficient (quasi)-Newton schemes exist to solve this problem in low-dimensional settings when the density of  $\mu$  is known [25, 22, 20]. In the more general scenario, where the dimension can be high and/or only sampled points from  $\mu$  are available, we reformulate the semi-dual OT problem as a convex expected minimization problem, defined as:

$$\min \{H(\mathbf{g}) := \mathbb{E}_{X \sim \mu} [h(\mathbf{g}, X)] \mid \mathbf{g} \in \mathbb{R}^M\}, \quad (4)$$

where for all  $\mathbf{g} \in \mathbb{R}^M$ ,  $x \in \mathbb{R}^d$ ,  $h(\mathbf{g}, x) = -\mathbf{g}^c(x) - \sum_{j=1}^M w_j g_j$ . Note that we deliberately multiplied the semi-dual functional by  $-1$  to frame it as a convex minimization problem instead of a concave maximization problem; however, this is not a universal convention in the literature.

No matter the cost,  $H$  is subdifferentiable everywhere, and we consider the subdifferential  $\partial H(\mathbf{g}) = \mathbb{E}_{X \sim \mu} [\partial_{\mathbf{g}} h(\mathbf{g}, X)]$ , where for  $x \in \mathbb{R}^d$  and  $j \in \llbracket 1, M \rrbracket$ , we define:

$$\partial_{\mathbf{g}} h(\mathbf{g}, x)_j = \mathbb{1}_{x \in \mathbb{L}_j(\mathbf{g})} - w_j.$$

As long as  $x$  is in the interior of a Laguerre cell, this subdifferential  $\partial H$  is, in fact, a differential that we note  $\nabla H$ . Moreover,  $\partial_{\mathbf{g}} h(\mathbf{g}, X)$  is an unbiased estimator of  $\partial H(\mathbf{g})$ . Given access to samples from  $\mu$  naturally leads to the study of stochastic gradient descent schemes of the form

$$\mathbf{g}_n = \mathbf{g}_{n-1} - \gamma_n \partial_{\mathbf{g}} h(\mathbf{g}, X_n),$$

where we start from an initial point  $\mathbf{g}_0 \in \mathbb{R}^M$ , and at each iteration  $n$ , draw a sample  $X_n$  and take a gradient step with step size  $\gamma_n > 0$  (also referred to as the learning rate).

SGD algorithms are well-suited for this setting [3, 8, 23], as they adapt to the number of samples drawn, have linear  $\mathcal{O}(M)$  complexity per iteration, efficiently handle mini-batches through GPU parallelization, and do not require storing the drawn samples. Directly solving (4) also helps avoid discretization bias when estimating OT quantities [7, 16], which can be crucial in some applications of semi-discrete OT [8]. However, the specific structure of the OT semi-dual problem makes analyzing the convergence of SGD algorithms particularly challenging, especially regarding convergence to the optimizer  $\mathbf{g}^*$ . Unlike standard cases, it does not fall within the class of well-behaved problems, such as those that are globally strongly convex.

**Entropic regularization.** The idea of formulating the semi-dual OT problem as the minimization of an expectation and avoiding discretization by using SGD algorithms was introduced in [16]. However, possibly due to the lack of globally favorable properties of  $H$ , they propose using the entropy-regularized version of  $H$ , referred to as the entropic semi-dual  $H_\varepsilon$ , where  $\varepsilon$  is the regularization parameter. This results in a globally  $1/\varepsilon$ -smooth problem and as  $\varepsilon$  vanishes,  $H_\varepsilon$  converges to  $H$ . Unfortunately, the theoretical analysis of SGD algorithms for the regularized problem reveals prohibitive constants in  $\varepsilon^{-1}$  and higher [6], making the use of a small  $\varepsilon$  impractical for theoretical guarantees. Thus, avoiding both regularization and discretization bias to solve the semi-discrete problem highlights the relevance of studying SGD for the non-regularized OT problem.

### 3 Projected Stochastic Gradient Descent on the Semi-Dual OT Problem

#### 3.1 Localizing a projection set

In convex optimization, particularly in an online or stochastic setting, localizing a set to restrict the optimization domain and using a projection step in the gradient descent scheme can be very useful, permitting straight-forward convergence proofs [19]. Our first lemma addresses this idea in the context of the OT semi-dual problem, showing that even when the support of  $\mu$  is not compact, it is still possible to localize a  $\|\cdot\|_\infty$ -ball within which a minimizer of the semi-dual function  $H$  exists. This projection set is formally defined in Lemma 3.1.

**Lemma 3.1** (Existence of a projection set). *Suppose that (4) admits a minimum. Then, there exists a minimizer  $\mathbf{g}^*$  contained in the set*

$$\mathcal{C} := \{\mathbf{g} \in \mathbb{R}^M \mid |g_j| \leq \|c\|_{K,\infty}\},$$

where  $\|c\|_{K,\infty} := \sup_{x \in K, j \in [1,M]} |c(x, y_j)|$ , for any compact  $K$  satisfying  $\mu(K) \geq 1 - \frac{1}{2} \min_j w_j$ .

While the existence of a  $\|\cdot\|_\infty$ -ball was previously established under the assumption that the cost function is uniformly bounded on  $\mathbb{R}^d \times \mathbb{R}^d$ , or when both measures have bounded support, we extend this result to the more general semi-discrete setting. On its own, this finding may enable further theoretical developments in semi-discrete OT, where a bounded potential is often required [2, 12].

**Example 1.** Consider  $\mu$  as the standard Gaussian on  $\mathbb{R}^3$ , the cost function  $c = \frac{1}{2}\|x - y\|^2$ , and  $\nu$  as a discrete measure with  $10^7$  points in  $[0, 1]^3$  and uniform weights. In this setting, it was previously believed that the Brenier potential could not be bounded, since  $c$  is not bounded on  $\mathbb{R}^3 \times [0, 1]^3$ . However, using Lemma 3.1, we can take the ball  $B(0, 6)$ , which allows us to restrict our search for the potential within a  $\|\cdot\|_\infty$ -ball of radius 18.

Incorporating this projection step into our SGD scheme has several advantages: (i) it significantly enhances both the practical performance and theoretical convergence of the algorithm; and (ii) the computational complexity of the projection step is  $\mathcal{O}(M)$ , as it simply involves clipping each coordinate of the vector. The projector is defined as

$$\text{Proj}_{\mathcal{C}}(\mathbf{g}) : \mathbf{g} \in \mathbb{R}^M \mapsto \arg\min\{\|\mathbf{g} - \mathbf{g}'\|; \mathbf{g}' \in \mathcal{C}\}.$$

In practice, this requires the knowledge of a compact set  $K$ , which we assume to be either given or previously estimated. The estimation is for instance trivial when the source is the standard Gaussian as in many applications (e.g. generative modeling). Based on this projection step, we derive the PSGD algorithm to minimize  $H$ , as presented in Algorithm 1.

### 3.2 A first convergence of PSGD in the general setting

As a first consequence of our projection step, we can directly establish a convergence rate for PSGD on the OT cost in a general setting, without assuming additional regularity of the cost function or the source measure. The proof follows a classical result for PSGD algorithms, and our projection step provides insights into the choice of the learning rate  $\gamma_1$ . Additionally, it allows us to recover the rate  $\mathcal{O}(1/\sqrt{n})$  from [16] for the averaged iterates of SGD when applied to the entropy-regularized semi-dual.

**Theorem 3.2** (PSGD in the general setting). *In*

*the general setting, choosing the learning rate  $\gamma_n = \gamma_1/n^b$  with  $\gamma_1 = \text{Diam}(\mathcal{C})/2\sqrt{2}$  and  $b = 1/2$ , we obtain*

$$\mathbb{E}[H(\bar{\mathbf{g}}_n) - H(\mathbf{g}^*)] \leq \frac{4\sqrt{2} \text{Diam}(\mathcal{C})}{\sqrt{n}}.$$

Although most of our results focus on MTW costs, as discussed in the next section, we also establish a convergence rate in a more general setting. This broader setting lies outside the scope of the generalized Brenier theorem and an OT map may not exist or be unique. However, estimating the OT cost can still be useful in certain applications, such as when  $c = \|\cdot\|$ , corresponding to the 1-Wasserstein distance. Thanks to our projection step, the same convergence rate carries over to other SGD-based methods, such as Adagrad [13], which is well-suited when the projection set is a hypercube [26, Section 4.2.4].

## 4 Convergence analysis of PSGD and minimax estimation for MTW costs

We now focus on costs satisfying the MTW properties, with particular attention to the cost  $c(x, y) = \frac{1}{2}\|x - y\|^2$ , as it is the most commonly used. For this cost, we extend our results to include non-compactly supported source measures. This setting is used, for instance, in [3, 23], where a standard Gaussian is mapped to a discrete distribution. Our objective is to establish the convergence rate of PSGD in approximating the true OT map and cost by estimating the Brenier potential  $\mathbf{g}^*$ . We make the following assumptions, distinguishing between the compact case, where we treat all MTW costs, and the non-compact case, where we focus solely on the quadratic cost  $\frac{1}{2}\|\cdot\|^2$ .

**Assumption A** (Compact case). The cost  $c$  satisfies the MTW condition,  $\text{Supp}(\mu)$  is bounded and  $c$ -convex (see Appendix A), and  $\mu$  satisfies a weighted  $(1, 1)$  Poincaré-Wirtinger inequality: there exists  $C_{\text{pw}} > 0$  such that for all  $f \in \mathcal{C}^1(\mathbb{R}^d)$ ,

$$\|f - \mathbb{E}_\mu[f]\|_{L^1(\mu)} \leq C_{\text{pw}} \|\nabla f\|_{L^1(\mu)}. \quad (\text{PW})$$

Note that (PW) relates to the local strong convexity of the semi-dual problem near the optimum and is a common assumption as in [20, 5]. Moreover, in the compact setting, the assumption that the density  $f_\mu$  is bounded from above and below by strictly positive values is a common assumption, as in [29]. This compact assumption implies (PW).

Regarding the non-compact case, our assumptions are satisfied notably, for non-degenerate Gaussians, finite mixtures of non-degenerate Gaussians, and heavy-tailed distributions such as Student distributions with degree of freedom larger than 2. While a broader class of MTW costs could potentially be covered, perhaps under stronger assumptions regarding the source measure, we deliberately omit these cases to avoid further technical complexity. Moreover, to the best of our knowledge, even for the quadratic cost, there are no existing theoretical results in semi-discrete OT when the support of  $\mu$  is unbounded.

**Assumption B** (Non-compact case). **(B1)** The cost is quadratic:  $c(x, y) = \frac{1}{2}\|x - y\|^2$  and the measure  $\mu$  has a finite second-order moment.

**(B2)** There exists a compact set  $K \subset \mathbb{R}^d$  with  $\mu(K) \geq 1 - \frac{1}{4}w_{\min}$ , such that the probability measure  $\mu_K$  with density  $f_{\mu_K}(x) := c_K f_\mu(x) \mathbf{1}_K(x)$  satisfies a (PW) inequality.

---

**Algorithm 1** Projected Stochastic Gradient Descent (PSGD)

---

**Parameters:**  $\gamma_1 > 0, b \in [\frac{1}{2}, 1)$

Initialize  $\mathbf{g}_0 \in \mathcal{C}$  and  $\bar{\mathbf{g}}_0 = \mathbf{g}_0$

**for**  $k = 1$  to  $n$  **do**

    Draw  $x_k \sim \mu$

$\mathbf{g}_k = \text{Proj}_{\mathcal{C}}(\mathbf{g}_{k-1} - \frac{\gamma_1}{k^b} \partial_{\mathbf{g}} h(\mathbf{g}_{k-1}, x_k))$

$\bar{\mathbf{g}}_k = \frac{1}{k+1} \mathbf{g}_k + \frac{k}{k+1} \bar{\mathbf{g}}_{k-1}$

**end for**

**return**  $\mathbf{g}_n$  and  $\bar{\mathbf{g}}_n$

---

(B3) The density  $f_\mu$  satisfies the following integrability and regularity condition: for  $R > 1$  and  $r \geq 1$ , define

$$f_\mu^R := f_\mu \cdot \mathbf{1}_{\|x\| \leq R}, \quad f_\mu^{R+r} := f_\mu \cdot \mathbf{1}_{R+(r-2) \leq \|x\| \leq R+r}.$$

Assume there exist  $C > 0$  and a modulus of continuity  $\omega$  such that for all  $\delta > 0$ ,

$$\sum_{r=0}^{\infty} (R+r)^{d-1} \omega_{f_\mu^{R+r}}(\delta) \leq C \omega(\delta), \quad \sum_{r=0}^{\infty} (R+r)^{d-1} C_{f_\mu^{R+r}} < \infty, \quad (5)$$

where  $C_{f_\mu^{R+r}} := \sup_{x \in \mathbb{R}^d} f_\mu^{R+r}(x)$  and  $\omega_{f_\mu^{R+r}}$  is the modulus of continuity of  $f_\mu^{R+r}$ .

#### 4.1 Properties of the semi-dual $H$

In our context, it is known that the discrete Brenier potential  $\mathbf{g}^* \in \mathbb{R}^M$  is unique only up to a transformation of the form  $\mathbf{g}^* + a\mathbb{1}_M$  with  $a \in \mathbb{R}$ . For clarity, we fix  $\mathbf{g}^*$  to be the Brenier potential such that  $\mathbf{g}^* \in \text{Vect}(\mathbb{1}_M)^\perp$ . Without losing information, we thus restrict our analysis to the orthogonal complement of the subspace spanned by the vector  $\mathbb{1}_M$ ,  $\text{Vect}(\mathbb{1}_M)^\perp$ . Therefore, for any  $\mathbf{g}, \mathbf{g}' \in \mathbb{R}^M$ , we define

$$\begin{aligned} \|\mathbf{g} - \mathbf{g}'\|_v &= \|\text{Proj}_{\text{Vect}(\mathbb{1}_M)^\perp}(\mathbf{g} - \mathbf{g}')\|, \\ \langle \mathbf{g}, \mathbf{g}' \rangle_v &= \langle \text{Proj}_{\text{Vect}(\mathbb{1}_M)^\perp}(\mathbf{g}), \text{Proj}_{\text{Vect}(\mathbb{1}_M)^\perp}(\mathbf{g}') \rangle. \end{aligned}$$

We start by stating the second-order regularity of  $H$  in our setting.

**Proposition 4.1.** *Under Assumption A or B, the function  $H$  is differentiable everywhere on  $\mathbb{R}^M$ , and we denote its gradient by  $\nabla H$ . Moreover, there exists a radius  $r > 0$  such that on the ball  $B(\mathbf{g}^*, r)$ ,  $H$  is  $C^2$  and strongly convex on  $B(\mathbf{g}^*, r)$ . If, in addition,  $f_\mu$  is  $\alpha$ -Hölder continuous with  $\alpha \in (0, 1]$ , then the Hessian of  $H$  is also  $\alpha$ -Hölder continuous.*

Naturally, the smallest eigenvalue of the Hessian is 0, with  $\mathbb{1}_M$  as its eigenvector. However, we still refer to the strong convexity of  $H$  since we focus on the orthogonal complement of  $\text{Vect}(\mathbb{1})$ . Notably, we extended the definition of the Hessian, originally provided in [20] for the compact case, to include the quadratic Euclidean cost in the non-compact setting. Furthermore, while the local strong convexity of  $H$  was established in [20] for the compact case, it was defined over the set of vectors  $\mathbf{g}$  such that the measures of all Laguerre cells are bounded by a positive constant. To formulate this result with respect to a ball, we also required a result on the quantitative stability of the measures of Laguerre cells. Additional details on this quantitative stability will be provided in section 5.1.

As a corollary of these results, we derive the following lemma, which stems from our projection step and can be viewed as a weak form of strong convexity of  $H$  on  $\mathcal{C}$ , also referred to as **Restricted Strong Convexity** (RSC) [34].

**Lemma 4.2.** *Under Assumptions A or B, there exists  $\eta > 0$  such that  $H$  satisfies a RSC property, uniformly for all  $\mathbf{g} \in \mathcal{C}$ :*

$$\langle \nabla H(\mathbf{g}), \mathbf{g} - \mathbf{g}^* \rangle_v \geq \eta \|\mathbf{g} - \mathbf{g}^*\|_v^2.$$

**Direct convergence guarantees under RSC.** As a direct consequence of the RSC property on  $\mathcal{C}$ , Projected SGD achieves linear convergence for our OT problem when using the step size  $\gamma_t = \frac{1}{\eta t}$ . However, our analysis does not provide a precise estimate of the parameter  $\eta$ , which is required to effectively implement this learning rate in practice, and is particularly difficult to obtain in the OT setting. This limitation motivates the following section, where we study PSGD with a learning rate of the form  $\gamma_t = \gamma_1/t^b$  for  $b < 1$ . This variant achieves optimal convergence rates without requiring prior knowledge of  $\eta$ .

**Remark 4.3.** Since RSC implies the Quadratic Growth (QG) condition  $H(\mathbf{g}) - H(\mathbf{g}^*) \geq \eta \|\mathbf{g} - \mathbf{g}^*\|_v^2$ , the same convergence guarantees also hold for other SGD variants that rely on either RSC or QG assumptions, such as S-Adagrad [9].

#### 4.2 Convergence rate of PSGD

**Convergence of the non-averaged iterates** Building on the RSC of  $H$  from Lemma 4.2, we derive the convergence rate of the non-averaged iterates of PSGD, which mirrors the convergence behavior observed in the strongly convex setting.

251 **Theorem 4.4** (Non-averaged iterates). *Under Assumptions A or B, and for any decay schedule of the*  
 252 *form  $\gamma_n = \gamma_1/n^b$  with  $\gamma_1 > 0$  and  $b \in (1/2, 1)$ , we have the convergence rate*

$$\mathbb{E}[\|\mathbf{g}_n - \mathbf{g}^*\|_v^2] = \mathcal{O}(1/n^b).$$

253 As  $b$  approaches 1, we observe a nearly linear convergence rate for the OT potential. In the next  
 254 section, we further show that a linear convergence rate is achievable by the averaged iterates sequence  
 255  $\bar{\mathbf{g}}_n$  and that it is minimax optimal, highlighting the strong performance of PSGD.

256 **Convergence of the averaged iterates** In convex stochastic optimization, averaging the iterates of  
 257 the SGD scheme is a widely used technique, as it enables achieving an optimal linear convergence  
 258 rate for strongly convex functions without requiring knowledge of the strong convexity parameter,  
 259 and regardless of the decay  $b \in (1/2, 1)$  for the gradient steps [28, 27]. Moreover, the averaged  
 260 scheme can adapt to the local strong convexity of the objective function, even when global strong  
 261 convexity does not hold [4].

262 Note that, as stated in Proposition 4.1,  $H$  is locally strongly convex, and as stated after Lemma 4.2,  
 263 the RSC parameter on  $\mathcal{C}$  is unknown. These observations motivate the study of the averaged iterates  
 264 of PSGD, and the next theorem confirms that these motivations hold true in our setting. To establish  
 265 this result, we also impose a mild regularity condition on  $f_\mu$ , requiring it to be  $\alpha$ -Hölder continuous  
 266 for some  $\alpha \in (0, 1]$ .

267 **Theorem 4.5** (Averaged iterates). *Under Assumptions A or B, and assuming that  $f_\mu$  is  $\alpha$ -Hölder*  
 268 *with  $\alpha \in (0, 1]$ , for any decay schedule of the form  $\gamma_n = \gamma_1/n^b$  with  $\gamma_1 > 0$  and  $b \in (\frac{1}{1+\alpha}, 1)$ , we*  
 269 *have the convergence rate*

$$\mathbb{E}[\|\bar{\mathbf{g}}_n - \mathbf{g}^*\|_v^2] = \mathcal{O}(1/n).$$

270 *Without assuming  $f_\mu$  to be  $\alpha$ -Hölder, and for  $b \in (1/2, 1)$ , we still obtain*

$$\mathbb{E}[\|\bar{\mathbf{g}}_n - \mathbf{g}^*\|_v^2] = \mathcal{O}(1/n^b).$$

271 As we can see, the averaged iterates achieve a linear convergence rate. Moreover, the constant term  
 272 hidden in the asymptotic rates, detailed in C, depends on the local strong convexity coefficient of  
 273  $H$  at the optimum. Notably, there is existing literature on  $\lambda_*^\varepsilon$  for the entropy-regularized semi-dual  
 274 problem when  $\mu$  has bounded support ([12], Theorem 3.2, [11], Proposition 5.1). These results can  
 275 be extended to estimate  $\lambda_*$  by letting the regularization parameter vanish.

## 276 4.2.1 Estimation of the OT cost

277 Building on the convergence rate of PSGD, we derive the corresponding rate for the OT cost  
 278 estimation.

279 **Corollary 4.6.** *Under Assumption A or B,  $H$  is  $C^1$ -smooth and uniformly bounded, so we have*

$$H(\mathbf{g}) - H(\mathbf{g}^*) = \mathcal{O}(\|\mathbf{g} - \mathbf{g}^*\|_v^2).$$

280 *Therefore, the OT cost exhibits the same convergence rate as the OT potential.*

281 As we can see, this result establishes a (nearly) linear convergence rate for the estimation of the  
 282 OT cost, matching the rate derived in Theorems 4.4 and 4.5. In particular, for MTW costs, a faster  
 283 convergence rate of  $\mathcal{O}(1/n)$  is achievable, in contrast to the  $\mathcal{O}(1/\sqrt{n})$  rate from Theorem 3.2 in the  
 284 general setting.

## 285 5 OT cost and map estimation with PSGD

### 286 5.1 Minimax estimation of the OT map and Brenier potential

287 Having the convergence rate of PSGD to the Brenier potential  $\mathbf{g}^*$ , we study here the convergence of  
 288 the map estimate  $T(\mathbf{g}) : x \mapsto x - \nabla \mathbf{g}^c(x)$ . Note that, as soon as there exists  $j \in \llbracket 1, M \rrbracket$  such that  $x$   
 289 is in the interior of  $\mathbb{L}_j(\mathbf{g}^*) \cap \mathbb{L}_j(\mathbf{g})$ , we have

$$T_{\mu,\nu}(x) = x - \nabla(\mathbf{g})^c(x) = y_j.$$

Therefore, a result on the quantitative stability of the measure of Laguerre cells is sufficient to establish a convergence rate for the map estimator obtained from PSGD. Such a result was previously established in the compact case in [5]. Here, we extend their result to the quadratic Euclidean cost in the non-compact setting, leading to the following theorem.

**Theorem 5.1.** *Under Assumption A or B, the function  $\mathbf{g} \mapsto \|T(\mathbf{g}) - T_{\mu,\nu}\|_{L^2(\mu)}^2$  is Lipschitz with respect to the infinity norm  $\|\cdot\|_\infty$ .*

As a corollary, we retrieve the convergence rate of our map estimator with PSGD.

**Corollary 5.2.** *Under the same assumptions as Theorem 4.4, taking  $\hat{\mathbf{g}}_n \in \{\mathbf{g}_n, \bar{\mathbf{g}}_n\}$*

$$\|T(\hat{\mathbf{g}}_n) - T_{\mu,\nu}\|_{L^p(\mu)} = \mathcal{O}\left(1/n^{b/2}\right).$$

*If in addition,  $f_\mu$  is  $\alpha$ -Hölder with  $\alpha \in (0, 1]$ , taking  $b \in (\frac{1}{1+\alpha}, 1)$ , we have*

$$\|T(\bar{\mathbf{g}}_n) - T_{\mu,\nu}\|_{L^p(\mu)} = \mathcal{O}(1/\sqrt{n}).$$

Furthermore, we prove that the rate  $\mathcal{O}(1/\sqrt{n})$  achieved by  $T(\bar{\mathbf{g}}_n)$  is minimax optimal.

**Theorem 5.3.** *Fixing  $c(x, y) = \frac{1}{2}\|x - y\|^2$  and  $\nu = \frac{1}{2}\delta_{\{0\}} + \frac{1}{2}\delta_{\{1\}}$  and noting  $\mathcal{P}_{Lip}(\mathbb{R})$  the set of probability measures on  $\mathbb{R}$  with Lipschitz densities, we have*

$$\inf_{T^{(n)}} \sup_{\mu \in \mathcal{P}_{Lip}(\mathbb{R})} \mathbb{E}_\mu \left[ \|T^{(n)} - T_{\mu,\nu}\|_{L^p(\mu)}^p \right] \gtrsim 1/\sqrt{n}.$$

*where the infimum is taken over all maps  $T^{(n)}$  constructed with the  $n$  i.i.d samples of  $\mu$ .*

We recover the same minimax lower bound as in the two-sample setting considered in [29], where the target measure  $\nu$  is also subsampled. This shows that, even though we have full information about the target measure, the asymptotic rates remain the same. However, we are able to achieve this rate in the non-batched setting, without the need to calibrate a regularization parameter as in [29], or to know the number of samples in advance. Note also that a direct corollary of Theorem 5.3 and Theorem 5.1 is that the convergence rate  $\mathcal{O}(1/n)$  for the estimation of the Brenier potential, achieved by the averaged iterates of PSGD, is also minimax optimal.

## 6 Numerical experiments

In this section, we numerically verify our convergence rate guarantees through various examples. All experiments demonstrating convergence rates were repeated 20 times, and the error plots represent the averaged errors. We set the learning rate to  $\gamma_1 = \text{Diam}(\mathcal{C})$ , as suggested by the analysis in Theorem 3.2. The step decay parameter  $b$  was set to  $3/4$ , unless stated otherwise. We find that this learning rate leads to robust results without requiring further tuning. For each example, we generate  $\mathbf{g}^*$  randomly, and approximate the associated Laguerre cell measures  $\mu(\mathbb{L}_i^c(\mathbf{g}))$ . We then fix  $w_i = \mu(\mathbb{L}_i^c(\mathbf{g}))$ , such that  $\mathbf{g}$  is optimal by the first-order condition. The Laguerre cells are estimated with  $10^9$  samples. All experiments were repeated 10 times, and the average performance was reported. We consider the following three settings to evaluate our method:

**Example 1: Non-quadratic cost.** The cost to move mass is set to  $\|\cdot\|^{1.5}$ . We set  $\mu$  as the uniform measure  $\mathcal{U}([0, 1]^{10})$  and take  $M = 50$  points  $y_1, \dots, y_M$  uniformly in  $[0, 1]^{50}$ . The projection set is then  $\mathcal{C} = [-10^{3/4}, 10^{3/4}]^{50}$ .

**Example 2: Non-compact case.** Here,  $\mu$  has full support on  $\mathbb{R}^{10}$  with cost  $c(x, y) = \frac{1}{2}\|x - y\|^2$ . We choose  $\mu$  with density  $f_\mu(x) \propto (1 + \|x\|)^{-d-3}$ , satisfying (B1-3). As in Example 1, we sample  $M = 50$  points in  $[0, 1]$ . The projection set  $\mathcal{C} = [-5, 5]^{50}$  since  $K = B(0, 1)$  satisfies Lemma 3.1.

**Example 3: Non-smooth source measure.** We define  $\mu$  with density  $f_\mu(x) = 1/(2\sqrt{x})\mathbf{1}_{x \in (0, 1]}$ , which satisfies a  $(1, 1)$ -Poincaré-Wirtinger inequality but is not  $\alpha$ -Hölder. We took  $M = 10$  points uniformly in  $[0, 1]$ . Since our results do not guarantee acceleration for non-Hölder densities, we set  $b = 0.9$  for PSGD, as our analysis recommends  $b$  close to 1 for the best rate. The projection set is  $\mathcal{C} = [-1, 1]^{10}$ .

As illustrated in Figure 1, our theoretical claims are well supported by empirical results. (i) Our convergence guarantees are matched: PSGD exhibits the expected convergence behavior across



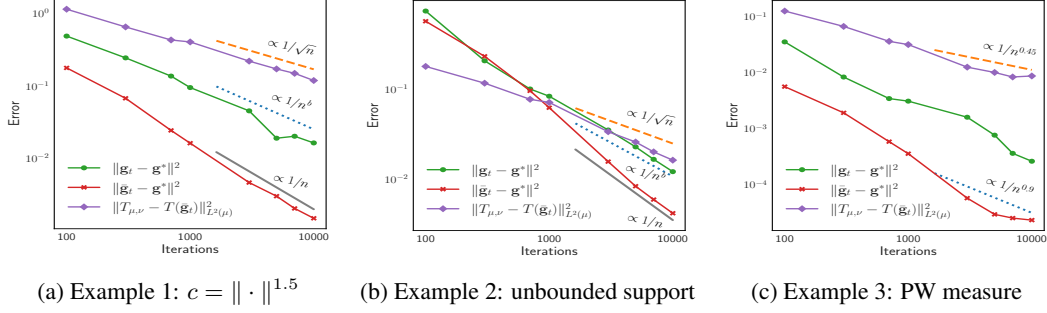


Figure 1: Convergence rates of our OT potential and OT map estimators across different settings.

all three settings. In particular, we observe a rate of  $\mathcal{O}(1/\sqrt{n})$  in Examples 1 and 2, and a rate of  $\mathcal{O}(1/n^{0.45})$  in Example 3 for estimating the OT map. The exponent 0.45 corresponds to  $b/2$  with  $b = 0.9$ , which aligns with our theoretical guidance to select  $b$  close to 1 when the source measure is not  $\alpha$ -Hölder regular but still satisfies a (PW) condition. (ii) Averaging yields optimal rates: In Examples 1 and 2, averaging leads to the optimal rate  $\mathcal{O}(1/\sqrt{n})$  for the OT map without requiring  $b = 1$ . This confirms our theory, which remains robust to the choice of  $\gamma_1 > 0$  and  $b \in (1/2, 1)$ , thanks to averaging, for achieving this minimax rate. (iii) We achieve minimax rates across our settings: Our results match the minimax rate  $\mathcal{O}(1/\sqrt{n})$  and extend the findings of [29], who established similar behavior in the compact case with quadratic cost. Importantly, we observe that the estimation of the OT map avoids the curse of dimensionality in both compact (MTW cost) and non-compact (quadratic cost) semi-discrete settings.

## 7 Conclusion and Discussion

We studied SGD-based solvers for the semi-discrete optimal transport (OT) problem, focusing on settings where only one or a few samples are available per iteration. These solvers are widely used in machine learning applications involving semi-discrete OT, yet their theoretical understanding remains incomplete. Our work bridges this gap by proving that such methods can consistently estimate both the OT cost and the OT map across a broad class of settings. Focusing on PSGD, we established minimax-optimal rates for estimating the OT map under MTW-type costs on compact domains, and under the quadratic Euclidean cost on both compact and non-compact domains. These results rely on novel convergence guarantees and structural properties of the semi-dual OT functional, stemming from the projection set we introduced and the enhanced properties of  $H$  we obtained thanks to the restriction of our minimization space.

**Future directions: exploiting RSC with adaptive methods.** Our analysis suggests a promising avenue for future work: leveraging RSC to improve the performance of adaptive SGD methods such as using S-Adam [33] with projection for OT. While Adam is commonly used in semi-discrete OT (especially in generative modeling), it often suffers from convergence plateaus due to its fixed step size. In contrast, S-Adam incorporates a decaying learning rate and is specifically tailored for strongly convex objectives. Despite lacking theoretical guarantees under RSC, our empirical results (Figure 2) show that Projected S-Adam better exploits the local geometry of the semi-dual problem, outperforming Projected Adam in practice. Formalizing these observations and extending our theory to include adaptive methods, while challenging, remains a compelling direction for future research.

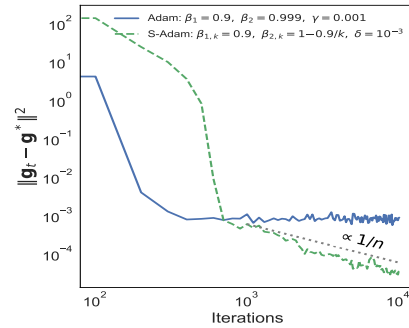


Figure 2: S-Adam outperforms Adam on Ex. 1, avoiding convergence plateau.

## References

- [1] P. K. Agarwal, S. Raghvendra, P. Shirzadian, and K. Yao. Fast and accurate approximations of the optimal transport in semi-discrete and discrete settings. In *Proceedings of the 2024 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 4514–4529. SIAM, 2024.
- [2] J. M. Altschuler, J. Niles-Weed, and A. J. Stromme. Asymptotics for semidiscrete entropic optimal transport. *SIAM Journal on Mathematical Analysis*, 54(2):1718–1741, 2022.
- [3] D. An, Y. Guo, N. Lei, Z. Luo, S.-T. Yau, and X. Gu. Ae-ot: A new generative model based on extended semi-discrete optimal transport. *ICLR 2020*, 2019.
- [4] F. Bach. Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic regression. *The Journal of Machine Learning Research*, 15(1):595–627, 2014.
- [5] M. Bansil and J. Kitagawa. Quantitative stability in the geometry of semi-discrete optimal transport. *International Mathematics Research Notices*, 2022(10):7354–7389, 2022.
- [6] B. Bercu and J. Bigot. Asymptotic distribution and convergence rates of stochastic algorithms for entropic optimal transportation between probability measures. *Annals of Statistics*, 49(2): 968–987, 2021.
- [7] L. Bottou and O. Bousquet. The tradeoffs of large scale learning. *Advances in neural information processing systems*, 20, 2007.
- [8] Y. Chen, M. Telgarsky, C. Zhang, B. Bailey, D. Hsu, and J. Peng. A gradual, semi-discrete approach to generative network training via explicit wasserstein minimization. In *International Conference on Machine Learning*, pages 1071–1080. PMLR, 2019.
- [9] Z. Chen, Y. Xu, E. Chen, and T. Yang. Sadagrad: Strongly adaptive stochastic gradient methods. In *International Conference on Machine Learning*, pages 913–921. PMLR, 2018.
- [10] V. Chernozhukov, A. Galichon, M. Hallin, and M. Henry. Monge–kantorovich depth, quantiles, ranks and signs. *The Annals of Statistics*, 45(1):223–256, 2017.
- [11] L. Chizat, A. Delalande, and T. Vaškevičius. Sharper exponential convergence rates for sinkhorn’s algorithm in continuous settings. *arXiv preprint arXiv:2407.01202*, 2024.
- [12] A. Delalande. Nearly tight convergence bounds for semi-discrete entropic optimal transport. In *International Conference On Artificial Intelligence And Statistics*, pages 1619–1642, 2022.
- [13] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- [14] N. Fournier and A. Guillin. On the rate of convergence in wasserstein distance of the empirical measure. *Probability theory and related fields*, 162(3):707–738, 2015.
- [15] B. Galerne, A. Leclaire, and J. Rabin. A texture synthesis model based on semi-discrete optimal transport in patch space. *SIAM Journal on Imaging Sciences*, 11(4):2456–2493, 2018.
- [16] A. Genevay, M. Cuturi, G. Peyré, and F. Bach. Stochastic optimization for large-scale optimal transport. In *Advances In Neural Information Processing Systems*, volume 29, 2016.
- [17] P. Ghosal and B. Sen. Multivariate ranks and quantiles using optimal transport: Consistency, rates and nonparametric testing. *The Annals of Statistics*, 50(2):1012–1037, 2022.
- [18] A. Godichon and B. Portier. An averaged projected robbins-monro algorithm for estimating the parameters of a truncated spherical distribution. *Electronic Journal of Statistics*, 11(1): 1890–1927, 2017.
- [19] E. Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- [20] J. Kitagawa, Q. Mérigot, and B. Thibert. Convergence of a newton algorithm for semi-discrete optimal transport. *Journal of the European Mathematical Society*, 21(9):2603–2651, 2019.
- [21] A. Leclaire and J. Rabin. A stochastic multi-layer algorithm for semi-discrete optimal transport with applications to texture synthesis and style transfer. *Journal of Mathematical Imaging and Vision*, 63(2):282–308, 2021.
- [22] B. Lévy. A numerical algorithm for  $l_2$  semi-discrete optimal transport in 3d. *ESAIM: Mathematical Modelling and Numerical Analysis*, 49(6):1693–1715, 2015.

- 421 [23] Z. Li, S. Li, Z. Wang, N. Lei, Z. Luo, and D. X. Gu. Dpm-ot: a new diffusion probabilistic  
422 model based on optimal transport. In *Proceedings of the IEEE/CVF International Conference on*  
423 *Computer Vision*, pages 22624–22633, 2023.
- 424 [24] X.-N. Ma, N. S. Trudinger, and X.-J. Wang. Regularity of potential functions of the optimal  
425 transportation problem. *Archive for rational mechanics and analysis*, 177:151–183, 2005.
- 426 [25] Q. Mérigot. A multiscale approach to optimal transport. In *Computer Graphics Forum*,  
427 volume 30, pages 1583–1592. Wiley Online Library, 2011.
- 428 [26] F. Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*, 2019.
- 429 [27] M. Pelletier. Asymptotic almost sure efficiency of averaged stochastic algorithms. *SIAM Journal*  
430 *on Control and Optimization*, 39(1):49–72, 2000.
- 431 [28] B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM*  
432 *journal on control and optimization*, 30(4):838–855, 1992.
- 433 [29] A.-A. Pooladian, V. Divol, and J. Niles-Weed. Minimax estimation of discontinuous optimal  
434 transport maps: The semi-discrete case. In *International Conference on Machine Learning*,  
435 pages 28128–28150. PMLR, 2023.
- 436 [30] F. Santambrogio. *Optimal transport for applied mathematicians*. Progress in nonlinear differen-  
437 tial equations and their applications. Birkhauser, 1 edition, Oct. 2015.
- 438 [31] C. Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.
- 439 [32] M. J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48.  
440 Cambridge university press, 2019.
- 441 [33] G. Wang, S. Lu, W. Tu, and L. Zhang. Sadam: A variant of adam for strongly convex functions.  
442 *arXiv preprint arXiv:1905.02957*, 2019.
- 443 [34] H. Zhang and L. Cheng. Restricted strong convexity and its applications to convergence analysis  
444 of gradient-type methods in convex optimization. *Optimization Letters*, 9:961–979, 2015.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: Proofs and numerical experiments are given for all the claims made in the abstract and the introduction.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: The paper clearly outlines its main limitations, including the underlying assumptions and the reliance on specific parameters, ensuring that the scope and constraints of the proposed approach are transparent to the reader.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: All assumptions are explicitly stated, and detailed proofs for each result are provided in the appendix.

### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: All experimental settings are clearly specified, and the parameters of our algorithm are provided.

### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [\[Yes\]](#)

Justification: All the code to reproduce our experiments is provided and attached in the supplementary materials.

### 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: All the parameters in each numerical experiment are explicitly provided.

### 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[Yes\]](#)

Justification: All the experiments are run several times, and the error plots represent the average error across the experiments.

### 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

494 Answer: [Yes]  
 495 Justification: All our experiments can be run on any modern computer, even without needing  
 496 a GPU.

497 **9. Code Of Ethics**  
 498 Question: Does the research conducted in the paper conform, in every respect, with the  
 499 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>  
 500 Answer: [Yes]  
 501 Justification: The authors have read the NeurIPS Code of Ethics and guarantee that the paper  
 502 conforms to it.

503 **10. Broader Impacts**  
 504 Question: Does the paper discuss both potential positive societal impacts and negative  
 505 societal impacts of the work performed?  
 506 Answer: [NA]

507 **11. Safeguards**  
 508 Question: Does the paper describe safeguards that have been put in place for responsible  
 509 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
 510 image generators, or scraped datasets)?  
 511 Answer: [NA]

512 **12. Licenses for existing assets**  
 513 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
 514 the paper, properly credited and are the license and terms of use explicitly mentioned and  
 515 properly respected?  
 516 Answer: [Yes]  
 517 Justification: The code provided is original and is included for reproducibility under the  
 518 correct license.

519 **13. New Assets**  
 520 Question: Are new assets introduced in the paper well documented and is the documentation  
 521 provided alongside the assets?  
 522 Answer: [Yes]  
 523 Justification: The algorithms and experiments are clearly explained in the paper, and the  
 524 code to reproduce the experiments is provided in an anonymized zip file.

525 **14. Crowdsourcing and Research with Human Subjects**  
 526 Question: For crowdsourcing experiments and research with human subjects, does the paper  
 527 include the full text of instructions given to participants and screenshots, if applicable, as  
 528 well as details about compensation (if any)?  
 529 Answer: [NA]

530 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**  
 531 **Subjects**  
 532 Question: Does the paper describe potential risks incurred by study participants, whether  
 533 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
 534 approvals (or an equivalent approval/review based on the requirements of your country or  
 535 institution) were obtained?  
 536 Answer: [NA]

537 **16. Declaration of LLM usage**  
 538 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
 539 non-standard component of the core methods in this research? Note that if the LLM is used  
 540 only for writing, editing, or formatting purposes and does not impact the core methodology,  
 541 scientific rigor, or originality of the research, declaration is not required.  
 542 Answer: [NA]

# Appendix

## Table of Contents

546	<b>A Further details on our Assumptions</b>	<b>16</b>
547	A.1 The Ma-Trudinger-Wang Properties . . . . .	16
548	A.2 The Poincaré-Wirtinger Inequality . . . . .	16
549	<b>B Properties of the semi-discrete OT problem</b>	<b>16</b>
550	B.1 Known results in the compact case . . . . .	17
551	B.2 New properties for the non-compact case with the quadratic Euclidean cost . . .	17
552	B.2.1 Definition and regularity of the Hessian . . . . .	18
553	B.2.2 Local strong convexity with respect to the mass of Laguerre cells . . . .	22
554	B.2.3 Quantitative stability of the Laguerre cells . . . . .	22
555	B.3 New properties in both our compact and non compact settings . . . . .	24
556	<b>C Proofs of the convergence rates of PSGD</b>	<b>25</b>
557	C.1 Fast convergence rates for MTW costs . . . . .	25
558	C.1.1 Convergence of the non-averaged iterates. . . . .	25
559	C.1.2 Convergence rate for higher order moments of the non-averaged iterates.	25
560	C.2 Convergence of the averaged iterates. . . . .	27
561	C.3 Convergence rate of PSGD in the general setting . . . . .	30
562	<b>D Proof of Lemma 3.1: Localisation of a projection set</b>	<b>31</b>
563	<b>E Minimax estimation of OT quantities</b>	<b>32</b>
564	E.1 Kantorovich potential . . . . .	32
565	E.2 OT map . . . . .	33
566	<b>F Technical Lemmas</b>	<b>34</b>
567	F.1 Technical Lemmas for Appendix B . . . . .	34
568	F.2 Technical Lemmas for Appendix C . . . . .	36

## Notations

- $\mathbb{R}^*$ : The set  $\mathbb{R} \setminus \{0\}$  (i.e., real numbers excluding zero).
- $\|\cdot\|$ : The Euclidean norm. item  $\lambda_{\mathbb{R}^d}$ : The Lebesgue measure on  $\mathbb{R}^d$ .
- **Diameter**: For  $\mathcal{C} \subset \mathbb{R}^d$ , we define its diameter as:

$$D_{\mathcal{C}} := \sup\{\|x - y\| \mid x, y \in \mathcal{C}\}.$$

- **Indicator function**: For a set  $A \subset \mathbb{R}^d$ ,  $\mathbb{1}_A(x)$  is defined as  $\mathbb{1}_A(x) = 1$  if  $x \in A$ , and  $\mathbb{1}_A(x) = 0$  otherwise.
- **Component-wise Minimum**: For  $v \in \mathbb{R}^d$ :

$$v_{\min} := \min_{1 \leq j \leq d} v_j.$$

- **Special Vectors**:

- $\mathbf{1}_M = (1, \dots, 1) \in \mathbb{R}^M$ .
- $\mathbf{0}_M = (0, \dots, 0) \in \mathbb{R}^M$ .
- $\mathbf{e}_j \in \mathbb{R}^M$ , for any  $1 \leq j \leq M$ , is the vector with zeros except for the  $j$ -th entry, which is equal to 1.

- **Probability Measures**:

- $\mathcal{P}(\mathbb{R}^d)$ : The set of probability measures on  $\mathbb{R}^d$ .
- For  $\rho \in \mathcal{P}(\mathbb{R}^d)$ ,  $\text{Supp}(\rho)$  denotes its support.

- **Asymptotic Orders**:

- $\mathcal{O}(\cdot)$  and  $o(\cdot)$ : Standard approximation orders.
- $f \lesssim g$  means there exists a constant  $C > 0$  such that  $f(\cdot) \leq Cg(\cdot)$ .
- $a \asymp b$  means both  $a \lesssim b$  and  $b \lesssim a$ .

- **Filtration**: We denote by  $\mathcal{F}_n$  the filtration generated by the sample  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mu$ , i.e.,

$$\mathcal{F}_n = \sigma(X_1, \dots, X_n), \quad n \geq 1.$$

- **Sets**:

- For any  $\varepsilon > 0$ , define:

$$K_{\varepsilon} := \{\mathbf{g} \in \mathbb{R}^M \mid \forall i \in \llbracket 1, M \rrbracket, \mu(\mathbb{L}_i(\mathbf{g})) \geq \varepsilon\}.$$

- Define:

$$K_+ := \{\mathbf{g} \in \mathbb{R}^M \mid \forall i \in \llbracket 1, M \rrbracket, \mu(\mathbb{L}_i(\mathbf{g})) > 0\}.$$

- **Density of an Absolutely Continuous Measure**: For an absolutely continuous measure  $\rho$  on  $\mathbb{R}^d$ , we denote its density w.r.t. the Lebesgue measure by  $f_{\rho}$ .

- **Orthogonal of Vect(1)**:

- For any  $\mathbf{g}, \mathbf{g}' \in \mathbb{R}^M$ , we define:

$$\|\mathbf{g} - \mathbf{g}'\|_v = \|\text{Proj}_{\mathbf{1}^{\perp}}(\mathbf{g} - \mathbf{g}')\|^2.$$

- Inner product in this space:

$$\langle \mathbf{g}, \mathbf{g}' \rangle_v = \langle \text{Proj}_{\mathbf{1}^{\perp}}(\mathbf{g}), \text{Proj}_{\mathbf{1}^{\perp}}(\mathbf{g}') \rangle.$$

- **Strong Convexity**: We discuss the strong convexity of the semi-dual function  $H$  when the strong convexity holds on the orthogonal complement of  $\mathbf{1}$ .

## A Further details on our Assumptions

### A.1 The Ma-Trudinger-Wang Properties

In the semi-discrete setting, the class of cost functions verifying the Ma–Trudinger–Wang (MTW) properties [24] is defined as the set of cost functions satisfying the following conditions: (Reg), (Twist), and Loeper’s condition (QC), detailed below.

$$\begin{aligned} c(\cdot, y_i) &\in C^2(\text{Supp}(\mu)), \forall i \in \{1, \dots, M\} & (\text{Reg}) \\ \nabla_x c(x, y_i) &\neq \nabla_x c(x, y_k), \forall x \in \text{Supp}(\mu), i \neq k & (\text{Twist}) \end{aligned}$$

**Definition A.1** (Loeper’s condition). We say  $c$  satisfies *Loeper’s condition* if for each  $i \in \{1, \dots, M\}$  there exists a convex set  $X_i \subset \mathbb{R}^d$  and a  $C^2$  diffeomorphism  $\exp_i^c(\cdot) : X_i \rightarrow \text{Supp}(\mu)$  such that

$$\forall t \in \mathbb{R}, 1 \leq k, i \leq M, \{p \in X_i \mid -c(\exp_i^c(p), y_k) + c(\exp_i^c(p), y_i) \leq t\} \text{ is convex.} \quad (\text{QC})$$

**Definition A.2** ( $c$ -convexity). We say that  $X \subset \mathbb{R}^d$  is  $c$ -convex if  $(\exp_i^c)^{-1}(X)$  is a convex set for every  $i \in \{1, \dots, N\}$

For a detailed discussion on this class of cost functions and their implications in the semi-discrete optimal transport framework, we refer the reader to Section 1.5 of [20]

### A.2 The Poincaré–Wirtinger Inequality

A probability measure  $\rho = w(x) dx$  on a domain  $\Omega \subset \mathbb{R}^d$  is said to satisfy the *weighted Poincaré–Wirtinger inequality* (PW) if

$$\int_{\Omega} |f - \mathbb{E}_{\rho}[f]| d\rho \leq C_{\text{PW}} \int_{\Omega} |\nabla f| d\rho, \quad \forall f \in C^1(\Omega).$$

The existence of a finite constant  $C_{\text{PW}}$  provides a quantitative connectedness of the source measure and is necessary for  $H$  to be locally strongly convex. We provide here two examples of measures satisfying (PW). Notably, Example 1 shows that non-degenerate Gaussians, mixture of non-degenerate Gaussians and Student distributions satisfy (PW), when we take their restrictions on any ball  $B(0, R)$ ,  $R > 0$ .

**Example 2.** (bounded support, density bounded above and below) Let  $\Omega$  be bounded, connected,  $\alpha$ -Holder with  $\alpha \in (0, 1]$ , and assume  $0 < m \leq w(x) \leq M < \infty$  almost everywhere.

Note that the assumption of the support being bounded from above and below is classical in the semi-discrete OT literature, as in [12, 29].

**Example 3.** (Annular support with radial concave profile, ([20], Proposition A.1)) Let  $0 < r < R$ , and let  $\bar{\rho} \in C^0([0, R])$  be a nonnegative function such that  $\bar{\rho}(s) = 0$  for  $s \in [0, r]$ , and  $\bar{\rho}$  is concave on  $[r, R]$ , with

$$\int_r^R \bar{\rho}(s) ds = 1.$$

Define the probability measure  $\rho$  on the annulus  $X := B(0, R) \subset \mathbb{R}^d$  by

$$\rho(x) = \frac{1}{\|x\|^{d-1} \omega_{d-1}} \bar{\rho}(\|x\|),$$

where  $\omega_{d-1}$  denotes the surface volume of the unit sphere  $\mathbb{S}^{d-1}$ . Then  $\rho$  satisfies the weighted Poincaré–Wirtinger inequality for some positive constant.

## B Properties of the semi-discrete OT problem

### Regularity properties of $H$

In the main article, we concisely presented the regularity properties of the function  $H$ . In this section, we provide a more detailed breakdown of these properties, organizing them into sub-properties and referring to the corresponding proofs in the quadratic case with unbounded support (assumptions B1–B3).



- 637 • **Differentiability:** The function  $H$  is differentiable on the entire space  $\mathbb{R}^M$ , and we denote  
638 its gradient by  $\nabla H$  (see Proposition B.5).
- 639 • **Local  $C^2$  regularity:** There exists a radius  $r > 0$  such that  $H$  is twice continuously  
640 differentiable ( $C^2$ ) on the ball  $B(\mathbf{g}^*, r)$  (see Proposition B.5).
- 641 • **Local strong convexity:** The function  $H$  is strongly convex on the ball  $B(\mathbf{g}^*, r)$  (see  
642 Proposition B.12).
- 643 • **Hölder continuity of the Hessian:** If the density  $f_\mu$  is  $\alpha$ -Hölder continuous for some  
644  $\alpha \in (0, 1]$ , then the Hessian of  $H$  inherits this regularity and is also  $\alpha$ -Hölder continuous  
645 (see Corollary B.11).

## 646 B.1 Known results in the compact case

647 In this section, we recall known properties of the semi-dual semi-discrete problem when the support  
648 of the source measure  $\mu$  is  $c$ -convex and contained within a compact set and  $c$  is a cost satisfying the  
649 MTW properties. We will then extend these results to the non-compact case for the quadratic cost.

650 Here, we fix  $\varepsilon > 0$  and recall that  $K_\varepsilon := \{\mathbf{g} \in \mathbb{R}^M : \forall i \in \llbracket 1, M \rrbracket, \mu(\mathcal{L}_i(\mathbf{g})) \geq \varepsilon\}$ . The two  
651 theorems presented below are taken from [20] and have been adapted to our notation. We emphasize  
652 that the authors of [20] considered the semi-dual OT problem as a concave problem, studying the  
653 objective function  $-H$  instead of  $H$  under their notation. For a better understanding of the constants  
654 in their theorems, we refer the reader to their article.

655 **Proposition B.1** (Theorem 1.1 in [20]). *Let  $\mu$  be an absolutely continuous density with bounded  
656 support included in  $\mathbb{R}^d$ , then the functional  $H$  is  $C^1$  smooth, its gradient is given by*

$$\nabla H(\mathbf{g})_i = -\mu(\mathbb{L}_i^c(\mathbf{g})) + w_i ,$$

657 *and its Hessian by*

$$\begin{aligned} (i \neq j) \quad \nabla^2 H(\mathbf{g})_{ij} &= - \int_{\mathbb{L}_i(\mathbf{g}) \cap \mathbb{L}_j(\mathbf{g})} \frac{f_\mu(x)}{\|y_i - y_j\|} d\mathcal{H}^{d-1}(x), \\ \nabla^2 H(\mathbf{g})_{ii} &= - \sum_{j \neq i} \nabla^2 H(\mathbf{g})_{ij} . \end{aligned}$$

658 **Proposition B.2** (Theorem 5.1 in [20]). *Under the assumption (A1), that  $\mu$  satisfies a weighted  
659  $(1,1)$ -Poincaré–Wirtinger inequality, there exists a constant  $\lambda$  such that for any  $\mathbf{g} \in K_\varepsilon$ , the second  
660 smallest eigenvalue of  $\nabla^2 H(\mathbf{g})$ , denoted  $\lambda_2(\nabla^2 H(\mathbf{g}))$ , satisfies*

$$\lambda_2(\nabla^2 H(\mathbf{g})) > \lambda .$$

661 *That is,  $H$  is strongly convex on  $K_\varepsilon$ , considering the problem on  $\mathbf{1}^\perp$ .*

662 **Theorem B.3** (Theorem 1.3 in [20]). *If  $\mu$  has its density  $f_\mu$  in  $C^{0,\alpha}(\text{Supp}(\mu))$ . Then, the functional  
663  $H$  is  $C^{2,\alpha}$  on the set*

$$K_\varepsilon := \{\mathbf{g} \in \mathbb{R}^M, \forall i, \mu(\mathbb{L}_i(\mathbf{g})) > \varepsilon\} ,$$

664 Lastly, we state a result concerning the quantitative stability of Laguerre cells, as presented in [5].

665 **Lemma B.4** (Lemma 5.5 in [5]). *Under the same assumptions as in Proposition B.2, for  $\mathbf{g}, \mathbf{g}' \in \mathbb{R}^M$ ,  
666 we have*

$$\mu(\mathbb{L}_i^c(\mathbf{g}) \setminus \mathbb{L}_i^c(\mathbf{g}')) \lesssim M \|\mathbf{g} - \mathbf{g}'\|_\infty, \quad \forall i \in \llbracket 1, M \rrbracket .$$

667 Once again, we refer to [5] for a more detailed understanding of the constant involved in this lemma.

## 668 B.2 New properties for the non-compact case with the quadratic Euclidean cost

669 In this section, we give second order properties of the semi-dual when the source measure is not  
670 supported on a compact. In the compact case, this properties are already known as discussed in  
671 Appendix B.1.

For the reader's convenience, we recall the Assumptions that we made for the non-compact case.

**Assumption (Non-compact case) (B1)** The cost is quadratic:  $c(x, y) = \frac{1}{2}\|x - y\|^2$  and the measure  $\mu$  has a finite second-order moment.

**(B2)** There exists a compact set  $K \subset \mathbb{R}^d$  with  $\mu(K) \geq 1 - \frac{1}{4}w_{\min}$ , such that the probability measure  $\mu_K$  with density  $f_{\mu_K}(x) := c_K f_\mu(x) \mathbf{1}_K(x)$  satisfies a Poincaré-Wirtinger inequality.

**(B3)** The density  $f_\mu$  satisfies the following integrability and regularity condition: for  $R > 1$  and  $r \geq 1$ , define

$$f_\mu^R := f_\mu \cdot \mathbf{1}_{\|x\| \leq R}, \quad f_\mu^{R+r} := f_\mu \cdot \mathbf{1}_{R+(r-2) \leq \|x\| \leq R+r}.$$

Assume there exist  $C > 0$  and a modulus of continuity  $\omega$  such that for all  $\delta > 0$ ,

$$\sum_{r=0}^{\infty} (R+r)^{d-1} \omega_{f_\mu^{R+r}}(\delta) \leq C \omega(\delta), \quad \sum_{r=0}^{\infty} (R+r)^{d-1} C_{f_\mu^{R+r}} < \infty, \quad (6)$$

where  $C_{f_\mu^{R+r}} := \sup_{x \in \mathbb{R}^d} f_\mu^{R+r}(x)$  and  $\omega_{f_\mu^{R+r}}$  is the modulus of continuity of  $f_\mu^{R+r}$ .

**Additional notation.** Since here the cost is fixed, we define the Laguerre cells by  $\mathbb{L}_i(\mathbf{g})$  instead of  $\mathbb{L}_i^c(\mathbf{g})$ .

### B.2.1 Definition and regularity of the Hessian

**Proposition B.5.** Under Assumptions (B1) and (B3), the semi-dual  $H$  is differentiable everywhere, and its gradient is given by

$$\nabla H(\mathbf{g})_i = \mu(\mathbb{L}_i(\mathbf{g})) - w_i.$$

Moreover  $H$  is  $C^2$  smooth on  $K_{w_{\min}/2} \cap \mathcal{C}$  and its Hessian is given by

$$(i \neq j) \quad \nabla^2 H(\mathbf{g})_{ij} = - \int_{\mathbb{L}_i(\mathbf{g}) \cap \mathbb{L}_j(\mathbf{g})} \frac{f_\mu(x)}{\|y_i - y_j\|} d\mathcal{H}^{d-1}(x),$$

$$\nabla^2 H(\mathbf{g})_{ii} = - \sum_{j \neq i} \nabla^2 H(\mathbf{g})_{ij}.$$

*Proof. Definition of the gradient.* The proof follows the lines of the proof of Theorem 4.1 of [20] and extend this result to the non-compact case. If  $x$  is in the interior of Laguerre cell  $\mathbb{L}_j(\mathbf{g})$ , we have

$$\nabla_{\mathbf{g}} h(\mathbf{g}, x) = \mathbb{1}_{i=j} - w_i.$$

where we recall that  $h(\mathbf{g}, x) = -\mathbf{g}^c(x) - \sum_{i=1}^M w_i g_i$ . Since the boundaries of the Laguerre cells are defined by the intersections of  $M$  hyperplanes, they form a negligible set with respect to the measure  $\mu$ . As a result, the gradient definition of  $H$  follows immediately.

**Definition of the Hessian.** Fix  $i \in \llbracket 1, M \rrbracket$ . We aim to prove the differentiability of the measure  $\mu(\mathbb{L}_i(\mathbf{g}))$  with respect to  $g_j$  and that its differential is defined by:

$$\frac{\partial \mu(\mathbb{L}_i(\mathbf{g}))}{\partial g_j} = - \int_{\mathbb{L}_i(\mathbf{g}) \cap \mathbb{L}_j(\mathbf{g})} \frac{f_\mu(x)}{\|y_i - y_j\|} d\mathcal{H}^{d-1}(x), \quad j \neq i,$$

$$\frac{\partial \mu(\mathbb{L}_j(\mathbf{g}))}{\partial g_j} = - \sum_{i \neq j} \frac{\partial \mu(\mathbb{L}_i(\mathbf{g}))}{\partial g_j}.$$

This gives us exactly the line  $(\nabla^2 H(\mathbf{g})_{1i}, \dots, \nabla^2 H(\mathbf{g})_{Mi})$  of the Hessian.

**Suppose  $i \neq j$ .**

Suppose  $\delta \geq 0$ . Defining  $h_{ij}(x) := \frac{1}{2}\|x - y_i\|^2 - \frac{1}{2}\|x - y_j\|^2 - g_i + g_j$ , note that we have  $\mathbb{L}_i(\mathbf{g}) = \cap_j h_{ij}^{-1}([-\infty, 0])$ .

We also have  $\mathbb{L}_i(\mathbf{g} + \delta \mathbf{e}_j) = \mathbb{L}_i(\mathbf{g}) \setminus (\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}([-\delta, 0]))$ . Note that  $h$  is Lipschitz and for all  $x$ ,  $\|\nabla h(x)\| = \|y_j - y_i\|$ . Moreover, under assumption (B3), we can use Lemma F.2,

698 which states that for all hyperplane  $H$ ,  $\int_H f_\mu d\mathcal{H}^{d-1} \lesssim 1$ . Therefore, we can apply the coarea formula  
 699 to pass from the second to the third equality below,

$$\begin{aligned} \mu(\mathbb{L}_i(\mathbf{g} + \delta \mathbf{e}_j)) &= \mu(\mathbb{L}_i(\mathbf{g})) - \mu(\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}([-\delta, 0])) \\ &= \mu(\mathbb{L}_i(\mathbf{g})) - \int_{\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}([-\delta, 0])} f_\mu(x) dx \\ &= \mu(\mathbb{L}_i(\mathbf{g})) - \int_{-\delta}^0 \int_{\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}(\{t\})} \frac{f_\mu(x)}{\|y_j - y_i\|} d\mathcal{H}^{d-1}(x) dt. \end{aligned}$$

700 By analogy, for  $\delta \leq 0$  we have:

$$\mathbb{L}_i(\mathbf{g} + \delta \mathbf{e}_j) = \mathbb{L}_i(\mathbf{g}) \cup (\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}([-\delta, 0])),$$

701 Using the coarea formula gives:

$$\mu(\mathbb{L}_i(\mathbf{g} + \delta \mathbf{e}_j)) = \mu(\mathbb{L}_i(\mathbf{g})) + \int_0^{-\delta} \int_{\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}(\{t\})} \frac{f_\mu(x)}{\|y_j - y_i\|} d\mathcal{H}^{d-1}(x) dt.$$

702 Applying Lemma B.6, which is stated and proved later in the appendice, the integrand defined above  
 703 is continuous on  $K_{w_{\min}/2} \cap \mathcal{C}$ . As a consequence, we can apply the Fundamental Theorem of Calculus  
 704 and justify the limit in

$$\lim_{\delta \rightarrow 0^-} \frac{\mu(\mathbb{L}_i(\mathbf{g} + \delta \mathbf{e}_j)) - \mu(\mathbb{L}_i(\mathbf{g}))}{\delta} = - \int_{\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}(\{0\})} f_\mu(x) \frac{1}{\|y_j - y_i\|} d\mathcal{H}^{d-1}(x).$$

705 By symmetry

$$\lim_{\delta \rightarrow 0^+} \frac{\mu(\mathbb{L}_i(\mathbf{g} + \delta \mathbf{e}_j)) - \mu(\mathbb{L}_i(\mathbf{g}))}{\delta} = - \int_{\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}(\{0\})} f_\mu(x) \frac{1}{\|y_j - y_i\|} d\mathcal{H}^{d-1}(x).$$

706 Since  $\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}(\{0\}) = \mathbb{L}_i(\mathbf{g}) \cap \mathbb{L}_j(\mathbf{g})$ , we thus obtain

$$F_{ij}(\mathbf{g}) := \frac{\partial \mu(\mathbb{L}_i(\mathbf{g}))}{\partial g_j} = - \int_{\mathbb{L}_i(\mathbf{g}) \cap \mathbb{L}_j(\mathbf{g})} \frac{f_\mu(x)}{\|y_i - y_j\|} d\mathcal{H}^{d-1}(x), \quad j \neq i.$$

707 **Suppose**  $i = j$ . No matter  $\mathbf{g} \in \mathbb{R}^M$ , the Laguerre cells verifies  $\mu(\cup_i \mathbb{L}_i(\mathbf{g})) = \mu(\mathbb{R}^d) = 1$ .  
 708 Therefore,

$$\frac{\partial}{\partial g_j} \sum_{i=1}^M \mu(\mathbb{L}_i(\mathbf{g})) = \frac{\partial}{\partial g_j} 1 = 0.$$

709 This equality gives  $\frac{\partial}{\partial g_j} \mu(\mathbb{L}_j(\mathbf{g})) = - \sum_{i \neq j} \frac{\partial}{\partial g_j} \mu(\mathbb{L}_i(\mathbf{g}))$ .

710 □

711 **Lemma B.6.** Under assumption (B1) and (B3), for any  $\mathbf{g} \in K_{w_{\min}/2} \cap \mathcal{C}$  the function  $F_{f_\mu}$  defined  
 712 for all  $t \in \mathbb{R}$  as

$$F_{f_\mu}(t) = \int_{\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}(\{t\})} \frac{f_\mu(x)}{\|y_j - y_i\|} d\mathcal{H}^{d-1}(x).$$

713 admits  $\omega$  as modulus of continuity in some neighborhood of  $\{0\}$ .

714 *Proof of Lemma B.6.* We keep the same notation as in [20]. Restricting ourselves to the quadratic  
 715 cost, we get

$$\begin{aligned} \varepsilon_{nd} = \varepsilon_{tw} &= \min_{i \neq j} \|y_i - y_j\|, \\ C_{\nabla}^R &= O(R), \\ C_{exp} &= O(1), \\ C_{cond} &= O(1), \\ C_{det} &= O(1). \end{aligned}$$

716 In order to apply the results in [20], we need to extend Proposition 4.5 to measures with unbounded  
 717 support in the case of the quadratic cost. We prove the following result:

718 **Proposition B.7.** Consider  $\mathcal{K}_\varepsilon := \{\mathbf{g} \in \mathbb{R}^{N-1} \mid \mu(\mathbb{L}_i(\mathbf{g})) > \varepsilon, \forall i = 1, \dots, M\}$  for some  $\varepsilon > 0$   
719 sufficiently small and let  $\mathcal{C}$  be the compact set of the projection. Then, there exists a positive  
720 constant  $\delta_1$ , such that for all  $\mathbf{g} \in \mathcal{K}_\varepsilon \cap \mathcal{C}$  and all  $p \in \mathbb{R}^d$  such that there exist  $i \neq j$ , for which  
721  $c(p, y_i) - c(p, y_0) = g_i - g_0$  and  $c(p, y_j) - c(p, y_0) = g_j - g_0$ , then

$$\left( \frac{\langle y_i - y_0, y_j - y_0 \rangle}{\|y_i - y_0\| \|y_j - y_0\|} \right)^2 \leq 1 - \delta_1^2. \quad (7)$$

722 *Remark B.8.* Note that the transversality equation (7) is independent of  $p \in \mathbb{R}^d$ . However,  $p$  is still  
723 involved in the transversality condition by the value of the difference of the costs.

724 *Proof.* Fix an index  $i$ . The set of points  $p \in \mathbb{R}^d$  such that  $c(p, y_i) - c(p, y_0) = g_i - g_0$  defines  
725 a hyperplane in  $\mathbb{R}^d$ . The intersection of two such hyperplanes, denoted  $H_{ij}(g)$ , is, unless the  
726 hyperplanes are parallel, a codimension-2 affine subspace of  $\mathbb{R}^d$ .

727 For each such  $H_{ij}(g)$ , let  $d_{ij}(g) \in \mathbb{R}^d$  be the orthogonal projection of the origin  $0_d$  onto  $H_{ij}(g)$ .  
728 Since the vector of potentials  $(g_i)_{i=0, \dots, M-1}$  is bounded due to the projection set  $\mathcal{C}$ , the set

$$\{d_{ij}(g) : i \neq j \text{ admissible, and } g \in \mathcal{K}_\varepsilon \cap \mathcal{C}\}$$

729 is contained in a ball  $B(0, R - 1)$  for some sufficiently large  $R > 0$ .

730 Hence, for every  $g \in \mathcal{K}_\varepsilon \cap \mathcal{C}$  and every admissible pair  $(i, j)$ , there exists a point  $p \in H_{ij}(g)$  lying in  
731 the interior of the ball  $B(0, R)$ .

732 We now apply the transversality result from [20] to the measure

$$\mu_R := \mathbf{1}_{B(0, R)} \cdot \frac{\mu}{\mu(B(0, R))},$$

733 which is the normalized restriction of  $\mu$  to  $B(0, R)$ . This result guarantees transversality of the  
734 hyperplanes  $H_{ij}$  associated with potentials in the set  $\mathcal{K}'_{\varepsilon'}$  (defined analogously for  $\mu_R$  and some  
735  $\varepsilon' > 0$ ) over the whole space  $\mathbb{R}^d$ .

736 Choose  $R$  such that  $\mu(B(0, R)) \geq 1 - \varepsilon/2$ . Then it is sufficient to take

$$\varepsilon' = \frac{\varepsilon/2}{1 - \varepsilon/2},$$

737 since one can check that  $\mathcal{K}_\varepsilon \subset \mathcal{K}'_{\varepsilon'}$ . This completes the proof.  $\square$

738 We aim to use the transversality result of Kitagawa et al. on a decomposition of  $\mathbb{R}^d$  into balls centered  
739 at 0. To this goal, we now prove a transversality result at the boundary of  $B(0, R)$  uniform for  $R$   
740 sufficiently large.

741 **Proposition B.9.** Consider  $\mathcal{K}_\varepsilon := \{\mathbf{g} \in \mathbb{R}^{N-1} \mid \mu(\mathbb{L}_i(\mathbf{g})) > \varepsilon, \forall i = 1, \dots, M\}$  for some  $\varepsilon > 0$   
742 sufficiently small and let  $\mathcal{C}$  be the compact set of the projection. Then, for any positive constant  $\delta_2$   
743 there exists  $R$  sufficiently large such that for all  $\mathbf{g} \in \mathcal{K}_\varepsilon \cap \mathcal{C}$  and all  $p \in \mathbb{R}^d$ : if there exists  $i$  such that  
744  $c(p, y_i) - c(p, y_0) = g_i - g_0$  for some  $p \in \partial B(0, R)$ , then

$$\left\langle \frac{p}{\|p\|}, \frac{y_i - y_0}{\|y_i - y_0\|} \right\rangle \leq 1 - \delta_2^2. \quad (8)$$

745 *Proof.* Rewrite the condition  $c(p, y_i) - c(p, y_0) = g_i - g_0$  as  $\langle p, y_i - y_0 \rangle = g_i - g_0$ , dividing by  
746  $\|p\|$  we get

$$\left\langle \frac{p}{\|p\|}, y_i - y_0 \right\rangle = \frac{g_i - g_0}{\|p\|}. \quad (9)$$

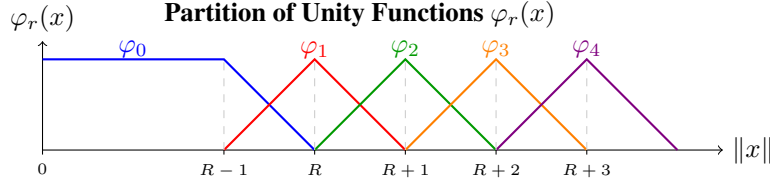
747 The right hand side tends to 0 uniformly with  $R$  since  $g_i - g_0$  lies in a compact set. The conclusion  
748 follows directly.  $\square$

749 The rest of the proof follows the lines of the proofs Appendix B of [20] with  $\varepsilon_{tr} = \delta_1$  as in Proposition  
 750 B.7. From Proposition B.9 applied to  $\delta_2 = \delta_1$  there exists  $R > 1$  such as the transversality condition  
 751 on the boundary (8) holds for every  $B(0, R + r)$  and the same  $\delta_1$  for every  $r \geq 0$ . Let us fix such  
 752  $R > 1$  so that Assumption (5) is also satisfied.

753 The next step is to decompose the integral using a partition of unity. We define the sequence of  
 754 functions  $(\varphi_r)_{r \geq 0}$  by

$$\varphi_r(x) = \begin{cases} (R - \|x\|)_+ \wedge 1, & \text{if } r = 0, \\ ((\|x\| - (R + r - 2))_+ \wedge (R + r - \|x\|)_+), & \text{if } r \geq 1, \end{cases} \quad x \in \mathbb{R}^d.$$

755 An illustration of the functions  $\varphi_r(x)$  defined above is shown below:



756 By definition,  $\sum_{r \geq 0} \varphi_r(x) = 1$ , for all  $x \in \mathbb{R}^d$ , and every  $\varphi_r$  is supported on  $\{R + r - 2 \leq \|x\| \leq$   
 757  $R + r\}$  for every  $r \geq 1$ ,  $\varphi_0$  being supported by  $B(0, R)$ . Moreover  $\varphi_r$  is Lipschitz continuous and we  
 758 denote  $\omega_{\varphi_r}$  its modulus of continuity on its support. The moduli of continuity satisfy  $\omega_{\varphi_r}(\delta) \leq |\delta|$ ,  
 759  $\delta > 0$   $r \geq 0$ . Then we decompose the integral

$$\int_{\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}(\{t\})} \frac{f_\mu(x)}{\|y_j - y_i\|} d\mathcal{H}^{d-1}(x) = \sum_{r=0}^{\infty} \int_{\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}(\{t\})} \frac{f_\mu(x) \varphi_r(x)}{\|y_j - y_i\|} d\mathcal{H}^{d-1}(x). \quad (10)$$

761 We apply Proposition B.1 of [20] to each term of the decomposition. We recall below his crucial  
 762 result, tracking the order of the constants established by [20].

**Proposition B.10.** *Let  $\sigma$  be a continuous non-negative function on  $B(0, R)$  bounded by  $\sigma_\infty$  and with modulus of continuity  $\omega_\sigma$ . Let the functions  $h_{ij}$  satisfy the transversality conditions (7) and (8) with the same constant  $\varepsilon_{tr} > 0$ . Then*

$$F_\sigma(t) := \int_{\cap_{k \neq j} h_{ik}^{-1}([-\infty, 0]) \cap h_{ij}^{-1}(\{t\})} \frac{\sigma(x)}{\|y_j - y_i\|} d\mathcal{H}^{d-1}(x)$$

has modulus of continuity

$$\omega_{h_\sigma}(\delta) = C_1 \omega_\sigma(C_2 \delta) + C_3 |\delta|$$

763 where  $C_1 = O(\mathcal{H}^{d-1}(\partial B(0, R)))$ ,  $\mathcal{H}^{d-1}(\partial B(0, R)) = O(R^{d-1})$ ,  $C_2 = O(\varepsilon_{tr}^{-1}) = O(1)$  and  
 764  $C_3 = O(\sigma_\infty C(d, 2R) \varepsilon_{tr}^{-4} + \mathcal{H}^{d-1}(\partial B(0, R)))$ , where  $C(d, 2R)$  defined in (3.5) of [20] satisfies  
 765  $C(d, 2R) = O(R^{d-1})$ .

When applying Proposition B.10 to the continuous function  $\sigma^{R+r}(x) := f_\mu(x) \varphi_r(x)$  we easily estimate  $\sigma_\infty^{R+r} \leq C_{f_\mu}^{R+r}$  and  $w_{\sigma^{R+r}} \leq \omega_{f_\mu}^{R+r} + C_{f_\mu}^{R+r} \omega_{\varphi_r}$ . Using that  $\omega_{\varphi_r}(\delta) \leq |\delta|$  we obtain:

$$\omega_{h_{\sigma^{R+r}}}(\delta) = O((R+r)^{d-1}) \left( \omega_{f_\mu}^{R+r}(O(1)\delta) + O(C_{f_\mu}^{R+r})\delta \right),$$

uniformly for every  $\mathbf{g} \in \mathcal{K}_\varepsilon \cap \mathcal{C}$  on a neighborhood of  $\{0\}$ . Noticing that this neighborhood depends solely on the transversality properties that are common to every  $r \geq 0$  and since

$$\omega_{F_{f_\mu}} = \sum_{r \geq 0} \omega_{h_{\sigma^{R+r}}}$$

766 from the decomposition in (10) the desired result follows under the assumption in (5).  $\square$

767 **Corollary B.11.** *Under Assumptions (B1) and (B3), if moreover  $f_\mu$  is  $\alpha$ -Hölder with  $\alpha \in (0, 1]$ , then*  
 768 *the Hessian of the semi-dual  $H$  is also  $\alpha$ -Hölder on  $K_{w_{\min}/2} \cap \mathcal{C}$ .*

769 *Proof.* Since  $f_\mu$  is  $\alpha$ -Hölder, the function  $f_\mu^{R+r}$  with  $r \geq 0$  are also  $\alpha$ -Hölder and we note  $\omega_{f_\mu}^{R+r} =$   
 770  $\kappa_{f_\mu}^{R+r} \delta^\alpha$ , applying Proposition B.10,

$$\begin{aligned} & \omega_{h_{\sigma R+r}}(\delta) O((R+r)^{d-1}) \left( \omega_{f_\mu}^{R+r}(O(1)\delta) + O(C_{f_\mu}^{R+r})\delta \right) \\ &= O((R+r)^{d-1}) \kappa_{f_\mu}^{R+r} O(1)\delta^\alpha + ((R+r)^{d-1}) O(C_{f_\mu}^{R+r})\delta. \end{aligned}$$

771 By the summability conditions of assumption (B3), for some constants  $C_1, C_2 > 0$ , we have

$$\begin{aligned} \omega_{F_{f_\mu}} &= \sum_{r \geq 0} \omega_{h_{\sigma R+r}} \\ &= O(1)\delta^\alpha \sum_{r \geq 0} O((R+r)^{d-1}) \kappa_{f_\mu}^{R+r} + \delta \sum_{r \geq 0} O(C_{f_\mu}^{R+r}) ((R+r)^{d-1}) \\ &\leq C_1 \delta^\alpha + C_2 \delta, \end{aligned}$$

772 Applying Proposition B.3 in [20] under our hypothesis shows that for  $\mathbf{g}, \mathbf{g}' \in K_{w_{\min}/2} \cap \mathcal{C}$ , there  
 773 exists a constant  $C$  depending on  $\varepsilon_{tr}$  and  $f_\mu$  such that we have

$$|\nabla^2 H(\mathbf{g}) - \nabla^2 H(\mathbf{g}')| \leq \omega_{F_{f_\mu}}(\|\mathbf{g} - \mathbf{g}'\|_\infty) + C\|\mathbf{g} - \mathbf{g}'\|$$

774 which gives the  $\alpha$ -Hölder regularity since  $\omega_{F_{f_\mu}} \leq C_1 \delta^\alpha + C_2 \delta$ . □

## 775 B.2.2 Local strong convexity with respect to the mass of Laguerre cells

776 **Proposition B.12.** *Under Assumptions (B1-B3)  $H$  is strongly convex on  $K_{\frac{1}{2}w_{\min}}$ .*

777 *Proof.* Recall  $K \subset \mathbb{R}^d$  the compact set such that  $\mu(K) > 1 - \frac{1}{4}w_{\min}$  and  $\mu_K$ , the probability  
 778 measure with density defined by  $f_{\mu_K}(x) = c_K f_\mu(x) \mathbf{1}_K(x)$ , with  $c_K \in [1, 2]$ , satisfies a weighted  
 779 Poincaré-Wirtinger inequality.

We thus can use Proposition B.2, which states that the semi-dual between  $\mu_R$  and  $\nu$  is strongly-convex on  $\text{Vect}_1^\perp$  on the set

$$K_{\frac{1}{4}w_{\min}}^R := \left\{ \mathbf{g} \in \mathbb{R}^M : \forall i \in \llbracket 1, M \rrbracket, \mu_R(\mathbb{L}_i^R(\mathbf{g})) \geq \frac{1}{4}w_{\min} \right\}.$$

780 This gives us that, for any  $\mathbf{g} \in K_{\frac{1}{4}w_{\min}}^R$ , there exists  $\lambda_K > 0$ , lower bounding the second smallest  
 781 value of the semi-dual function between  $\mu_R$  and  $\nu$ . This also gives us that for any  $\mathbf{g} \in K_{\frac{1}{4}w_{\min}}^R$ , the  
 782 second smallest eigenvalue of the matrix  $B$  defined by

$$\begin{aligned} (i \neq j) \quad B_{ij} &= - \int_{\mathbb{L}_i(\mathbf{g}) \cap \mathbb{L}_j(\mathbf{g})} \frac{f_\mu(x) \mathbf{1}_{B(0,R)}(x)}{\|y_i - y_j\|} d\mathcal{H}^{d-1}(x), \\ B_{ii} &= - \sum_{j \neq i} B_{ij}, \end{aligned}$$

783 is lower bounded by  $\lambda_K/c_R \geq \lambda_K/2 > 0$ .

784 Since for any  $\mathbf{g}$ , and  $i \neq j$ ,  $\nabla^2 H(\mathbf{g})_{ij} \leq B_{ij}$  and that both  $\nabla^2 H(\mathbf{g})$  and  $B$  are Laplacian matrices,  
 785 we apply Lemma F.1 to obtain that the second smallest eigenvalue of the hessian  $\nabla^2 H$  is lower  
 786 bounded by  $\lambda_R/c_R$  on  $K_{\frac{1}{4}w_{\min}}^R$ .

787 For any  $\mathbf{g} \in K_{w_{\min}/2}$ , we thus have for all  $i \in \llbracket 1, M \rrbracket$ ,  $\mu(\mathbb{L}_i(\mathbf{g}) \cap B(0, R)) \geq \frac{1}{4}w_{\min}$ . That is,  
 788  $K_{\frac{1}{2}w_{\min}} \subset K_{\frac{1}{4}w_{\min}}^R$ , which completes the proof. □

## 789 B.2.3 Quantitative stability of the Laguerre cells

790 **Lemma B.13.** *Under Assumptions (B1-B3), we have*

$$\mu(\mathbb{L}_i^c(\mathbf{g}) \setminus \mathbb{L}_i^c(\mathbf{g}')) \lesssim M\|\mathbf{g} - \mathbf{g}'\|_\infty, \quad \forall i \in \llbracket 1, M \rrbracket.$$

791 *Proof.* For all  $j \in \llbracket 1, M \rrbracket$ , if  $x$  is in the interior of  $\mathbb{L}_j(\mathbf{g})$ , we have

$$\nabla \mathbf{g}^c(x) = x - y_j. \quad (11)$$

792 Therefore, given  $\mathbf{g}, \mathbf{g}' \in \mathbb{R}^M$ , if there exists  $j \in \llbracket 1, M \rrbracket$  such that  $x$  is the interior of  $\mathbb{L}_j(\mathbf{g}) \cap \mathbb{L}_j(\mathbf{g}')$   
793 we have

$$\nabla \mathbf{g}^c(x) = \nabla (\mathbf{g}')^c(x).$$

794 Moreover, since the support of  $\nu$  is finite, we have

$$\sup_{x \in \mathbb{R}^d} \|\nabla \mathbf{g}^c(x) - \nabla (\mathbf{g}')^c(x)\| = \max_{i \neq j} \|y_i - y_j\|.$$

795 Hence, to bound the error of  $T(\mathbf{g})(x) = x - \nabla \mathbf{g}^c(x)$  and  $T_{\mu, \nu} = x - \nabla (\mathbf{g}^*)^c(x)$ , we just need to  
796 bound the difference of measure of Laguerre cells made by  $\mathbf{g}$  and  $\mathbf{g}^*$ . More generally, we now proceed  
797 here to bound the difference of measure of Laguerre cells between  $\mathbf{g}, \mathbf{g}' \in \mathbb{R}^M$  fixed arbitrarily.

798 Our proof will follow some arguments from [20]. Let us fix  $i \in \llbracket 1, M \rrbracket$  and suppose  $x \in \mathbb{L}_i(\mathbf{g}) \setminus$   
799  $\mathbb{L}_i(\mathbf{g}')$ . The definition of the Laguerre cells implies that there is a  $k \neq i$  such that  $c(x, y_k) + \mathbf{g}'_k <$   
800  $c(x, y_i) + \mathbf{g}'_i$  while  $c(x, y_i) + \mathbf{g}_i \leq c(x, y_k) + \mathbf{g}_k$ . Combining these two inequalities yields to

$$\mathbf{g}'_k - \mathbf{g}'_i < c(x, y_i) - c(x, y_k) \leq \mathbf{g}_k - \mathbf{g}_i.$$

801 Hence, writing  $f_{ik}(x) = c(x, y_i) - c(x, y_k)$ , we have

$$\mathbb{L}_i(\mathbf{g}) \setminus \mathbb{L}_i(\mathbf{g}') \subset \bigcup_{k \neq i} f_{ik}^{-1}([\mathbf{g}'_k - \mathbf{g}'_i, \mathbf{g}_k - \mathbf{g}_i]). \quad (12)$$

802 We now now bound  $\mu(f_{ik}^{-1}([\mathbf{g}'_k - \mathbf{g}'_i, \mathbf{g}_k - \mathbf{g}_i]))$  using the coarea formula, using that for all  $x$ ,  
803  $\|\nabla f_{ik}(x)\| = \|y_i - y_k\|$ :

$$\begin{aligned} \mu(f_{ik}^{-1}([a, b])) &= \int_{f_{ik}^{-1}([a, b])} d\mu(x) \\ &= \int_a^b \int_{f_{ik}^{-1}(\{t\})} \frac{1}{\|\nabla f_{ik}(x)\|} \mu(x) d\mathcal{H}^{d-1}(x) dt \\ &= \int_a^b \int_{f_{ik}^{-1}(\{t\})} \frac{1}{\|y_i - y_k\|} \mu(x) d\mathcal{H}^{d-1}(x) dt. \end{aligned}$$

804 Observe that for the quadratic cost,  $f_{ik}(x) = \langle x, y_k - y_i \rangle + \|y_i\|^2 - \|y_k\|^2$  and so  $f_{ik}^{-1}(\{t\}) =$   
805  $\{x \in \mathbb{R}^d, \langle x, y_k - y_i \rangle = t - \frac{1}{2}(\|y_i\|^2 + \|y_k\|^2)\}$  and is therefore a hyperplane. Applying Lemma  
806 F.2, there exists a constant  $C$  such that, for any  $i, k$  and  $t$ , we have

$$\int_{f_{ik}^{-1}(\{t\})} \frac{1}{\|y_i - y_k\|} \mu(x) d\mathcal{H}^{d-1}(x) \leq C.$$

807 Therefore, we have

$$\mu(f_{ik}^{-1}([a, b])) \leq \int_a^b \frac{C}{\|y_i - y_k\|} \leq (b - a) \frac{C}{\|y_i - y_k\|}.$$

808 Since  $\mathbf{g}_k - \mathbf{g}_i - (\mathbf{g}'_k - \mathbf{g}'_i) \leq 2\|\mathbf{g} - \mathbf{g}'\|_\infty$ , by combining the above with (12) we conclude

$$\mu(\mathbb{L}_i(\mathbf{g}) \setminus \mathbb{L}_i(\mathbf{g}')) \leq \sum_{k \neq i} \mu(f_{ik}^{-1}([\mathbf{g}'_k - \mathbf{g}'_i, \mathbf{g}_k - \mathbf{g}_i])) \lesssim M\|\mathbf{g} - \mathbf{g}'\|_\infty. \quad (13)$$

809 □

### 810 B.3 New properties in both our compact and non compact settings

811 **Theorem 5.1.** Under Assumptions (A1) or (B1-3), the function  $\mathbf{g} \mapsto \|T(\mathbf{g}) - T_{\mu,\nu}\|_{L^2(\mu)}^2$  is Lipschitz  
 812 with respect to the infinity norm  $\|\cdot\|_\infty$ . Moreover, the Lipschitz constant grows at most quadratically  
 813 in  $M$ .

814 *Proof.* Using that no matter  $\mathbf{g} \in \mathbb{R}^M$ , for any  $x \in \mathbb{R}^d$

$$\|T_{\mu,\nu}(\mathbf{g}^*) - T_{\mu,\nu}(\mathbf{g})\| \leq \begin{cases} 0 & \text{if } x \in \mathbb{L}_i(\mathbf{g}^*) \cap \mathbb{L}_i(\mathbf{g}) \text{ for a certain } i \in \llbracket 1, M \rrbracket, \\ \max_{i \neq j} \|y_i - y_j\| & \text{else,} \end{cases}$$

815 and that for any  $\mathbf{g}, \mathbf{g}' \in \mathbb{R}^M$ , we have  $\mu(\mathbb{L}_i(\mathbf{g}) \setminus \mathbb{L}_i(\mathbf{g}')) \lesssim M\|\mathbf{g} - \mathbf{g}'\|$  using Proposition B.4 in the  
 816 compact case, or Proposition B.13, we have for any  $p \in [1, \infty)$

$$\begin{aligned} \|T_{\mu,\nu}(\mathbf{g}^*) - T_{\mu,\nu}(\mathbf{g})\|_{L^p(\mu)}^p &= \int_{\mathbb{R}^d} \|T_{\mu,\nu}(\mathbf{g}^*)(x) - T_{\mu,\nu}(\mathbf{g})(x)\|^p d\mu(x) \\ &= \sum_{i=1}^M \int_{\mathbb{L}_i(\mathbf{g}^*)} \|T_{\mu,\nu}(\mathbf{g}^*)(x) - T_{\mu,\nu}(\mathbf{g})(x)\|^p d\mu(x) \\ &\leq \sum_{i=1}^M \int_{\mathbb{L}_i(\mathbf{g}^*) \setminus \mathbb{L}_i(\mathbf{g})} \max_{i \neq j} \|y_i - y_j\|^p d\mu(x) \\ &\leq \sum_{i=1}^M \max_{i \neq j} \|y_i - y_j\|^p \mu(\mathbb{L}_i(\mathbf{g}^*) \setminus \mathbb{L}_i(\mathbf{g})) \\ &\lesssim M^2 \|\mathbf{g}^* - \mathbf{g}\|_\infty, \end{aligned}$$

817 where we used  $\mu(\mathbb{L}_i(\mathbf{g}) \setminus \mathbb{L}_i(\mathbf{g}^*)) \lesssim M\|\mathbf{g} - \mathbf{g}^*\|$  for the last line.

818

□

819 **Lemma B.14.** For any  $\varepsilon > 0$ , there exists  $\mathbf{g}^* \in K_\varepsilon$  and  $d_\varepsilon > 0$  such that  $B(\mathbf{g}^*, d_\varepsilon) \subset K_\varepsilon$ .

820 *Proof.* This result is a simple application of the results on the quantitative stability of Laguerre cells  
 821 stated in Proposition B.4 and Proposition B.13. □

822 **Proposition B.15.** (Hessian and Local strong convexity). Under Assumptions (A1) or (B1-3),  $H$  is  
 823 twice differentiable. Moreover, there exists a constant  $\lambda > 0$  and a radius  $r > 0$ , such that for any  
 824  $\mathbf{g} \in \mathbb{R}^M$ ,  $\|\mathbf{g}^* - \mathbf{g}\|_v \leq r$  implies that the second smallest value of the Hessian  $\nabla^2 H(\mathbf{g})$  is lower  
 825 bounded by  $\lambda$ .

826 *Proof.* The fact that  $H$  is twice differentiable was already known in the compact case and is proven  
 827 in Proposition B.5 for the non compact case under (B1-3).

828 Under assumptions (A1) and using Proposition B.2 or under assumptions (B1-3) and using Proposition  
 829 B.12,  $H$  is strongly convex on  $K_{\frac{1}{2}w_{\min}}$ . Applying Lemma B.14 concludes the proof. □

830 **Proposition B.16.** There exists  $\eta > 0$ , such that uniformly in  $\mathbf{g} \in \mathcal{C}$ , we have

$$\langle \nabla H(\mathbf{g}), \mathbf{g} - \mathbf{g}^* \rangle_v \geq \eta \|\mathbf{g} - \mathbf{g}^*\|_v^2.$$

831 *Proof.* Observe that, by Proposition B.12,  $H$  is locally strongly convex on the orthogonal of  $\mathbf{1}$ , on  
 832 the set

$$K_{\frac{1}{2}w_{\min}} = \left\{ \mathbf{g} : \forall i \in \llbracket 1, M \rrbracket, \mu(\mathbb{L}_i(\mathbf{g})) \geq \frac{1}{2}w_{\min} \right\}.$$

833 Therefore, for any  $\mathbf{g} \in K_{\frac{1}{2}w_{\min}}$ , we have  $\langle \nabla H(\mathbf{g}), \mathbf{g} - \mathbf{g}^* \rangle_v \geq \lambda \|\mathbf{g} - \mathbf{g}^*\|_v^2$ .

834 Now, suppose  $\mathbf{g} \in \mathcal{C} \setminus K_{\frac{1}{2}w_{\min}}$ . Using Lemma B.14, we know that there exists  $d_0 > 0$  such that  
 835  $B(\mathbf{g}^*, d_0) \subset K_{\frac{1}{2}w_{\min}}$ . Therefore, by the convexity of  $H$ , there exists  $c > 0$  such that

$$\langle \nabla H(\mathbf{g}), \mathbf{g} - \mathbf{g}^* \rangle_v \geq H(\mathbf{g}) - H(\mathbf{g}^*) > c.$$



836 Defining

$$\lambda' := \inf_{\mathbf{g} \in \mathcal{C} \setminus K_{\frac{1}{2}w_{\min}}} \left\{ \frac{H(\mathbf{g}) - H(\mathbf{g}^*)}{\|\mathbf{g} - \mathbf{g}^*\|_v^2} \right\} > c/\text{Diam}_2(\mathcal{C})^2,$$

837 we have for any  $\mathbf{g} \in \mathcal{C} \setminus K_{\frac{1}{2}w_{\min}}$

$$\langle \nabla H(\mathbf{g}), \mathbf{g} - \mathbf{g}^* \rangle_v \geq \lambda' \|\mathbf{g} - \mathbf{g}^*\|_v^2.$$

838 Taking  $\eta = \min \{\lambda, \lambda'\}$  concludes the proof.  $\square$

## 839 C Proofs of the convergence rates of PSGD

### 840 C.1 Fast convergence rates for MTW costs

#### 841 C.1.1 Convergence of the non-averaged iterates.

842 **Theorem 4.5** (Non averaged iterates). Under Assumptions (A1) or (B1-3) and for any decay step of  
843 the form  $\gamma_n = \gamma_1/n^b$  with  $\gamma_1 > 0, b \in (1/2, 1)$ , we have the convergence rate

$$\mathbb{E}[\|\mathbf{g}_n - \mathbf{g}^*\|_v^2] = \mathcal{O}\left(\frac{\gamma_n}{\eta}\right).$$

844 *Proof.* By definition of the gradient step at time  $n + 1$ , sampling  $X_{n+1} \sim \mu$  and since  $\mathbf{g}^* \in \mathcal{C}$ , we  
845 have  $\|\nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1})\|_v \leq 2$  a.s. and

$$\begin{aligned} \|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^2 &= \|\text{Proj}_{\mathcal{C}}(\mathbf{g}_n - \gamma_{n+1} \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1})) - \mathbf{g}^*\|_v^2 \\ &\leq \|\mathbf{g}_n - \gamma_{n+1} \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}) - \mathbf{g}^*\|_v^2 \\ &\leq (\|\mathbf{g}_n - \mathbf{g}^*\|_v^2 - 2\gamma_{n+1} \langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v + \gamma_{n+1}^2 \|\nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1})\|_v^2) \\ &\leq (\|\mathbf{g}_n - \mathbf{g}^*\|_v^2 - 2\gamma_{n+1} \langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v + 4\gamma_{n+1}^2). \end{aligned} \quad (14)$$

846 Using Proposition B.16, we have for any  $\mathbf{g} \in \mathcal{C}$ ,  $\langle \nabla H(\mathbf{g}), \mathbf{g} - \mathbf{g}^* \rangle_v \geq \eta \|\mathbf{g} - \mathbf{g}^*\|_v^2$ . Therefore,  
847 taking the conditional expectation, we obtain

$$\mathbb{E}[\|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^2 \mid \mathcal{F}_n] \leq \|\mathbf{g}_n - \mathbf{g}^*\|_v^2 (1 - 2\eta\gamma_{n+1}) + 4\gamma_{n+1}^2.$$

848 Taking the expectation and applying Lemma F.3 with  $\delta_n = \mathbb{E}[\|\mathbf{g}_n - \mathbf{g}^*\|_v^2]$  and  $m_n = 4\gamma_n$  gives

$$\mathbb{E}[\|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^2] \leq \exp\left(-2\eta \sum_{k=\lceil n/2 \rceil}^n \gamma_k\right) \left(\sum_{n=n_0}^n 4\gamma_n^2 + \mathbb{E}[\|\mathbf{g}_{n_0} - \mathbf{g}^*\|_v^2]\right) + \frac{4}{\eta} \gamma_{\lceil n/2 \rceil - 1}$$

849 where  $n_0 = \min\{n \in \mathbb{N}, \eta\gamma_{n+1} \leq 1\}$ . Remark that the exponential term converges exponentially  
850 fast. Indeed, we have  $\sum_{k=\lceil n/2 \rceil}^n \gamma_k \gtrsim n^{1-b}$  with  $1 - b > 0$ . Moreover,  $\|\mathbf{g}_{n_0} - \mathbf{g}^*\|_v^2 \leq \text{Diam}_2(\mathcal{C})^2$ .

851 Therefore, since for all  $n \geq 2$ ,  $\gamma_{\lceil n/2 \rceil - 1} \leq 2^b \gamma_n$ , we have the desired convergence rate

$$\mathbb{E}[\|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^2] \lesssim \frac{\gamma_n}{\eta}$$

852 which concludes the proof.  $\square$

#### 854 C.1.2 Convergence rate for higher order moments of the non-averaged iterates.

855 We prove here the convergence rate of higher order moments of the error  $\|\mathbf{g}_n - \mathbf{g}^*\|_v$ . This conver-  
856 gence will be useful for the convergence rate of the averaged iterates of PSGD. While this proposition  
857 directly proves Theorem 4.4 by the use of Jensen's inequality, the proof is slightly more cumbersome  
858 so we decided to make a separate case.

859 **Proposition C.1.** *Under Assumptions (A1) or (B1-3) and for any decay step of the form  $\gamma_n = \gamma_1/n^b$*   
 860 *with  $\gamma_1 > 0, b \in (1/2, 1)$  and  $p \in \{1, 2, 3\}$ , we have the convergence rate*

$$\mathbb{E} [\|\mathbf{g}_n - \mathbf{g}^*\|^p] \lesssim \frac{\gamma_n^p}{\eta^p}$$

861 *where  $(\mathbf{g}_n)$  is the sequence of non-averaged iterates of PSGD.*

862 *Proof.* By definition of the gradient step at time  $n + 1$ , sampling  $X_{n+1} \sim \mu$  and using inequality  
 863 (14),

$$\|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^6 \leq (\|\mathbf{g}_n - \mathbf{g}^*\|_v^2 - 2\gamma_{n+1} \langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v + 4\gamma_{n+1}^2)^3.$$

864 Using that  $(A + B + C)^3 = \sum_{a+b+c=3} \frac{3!}{a!b!c!} A^a B^b C^c$ , we obtain

$$\begin{aligned} \|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^6 &\leq \|\mathbf{g}_n - \mathbf{g}^*\|_v^6 \\ &\quad - 6\|\mathbf{g}_n - \mathbf{g}^*\|_v^4 \gamma_{n+1} \langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v \\ &\quad + 12\|\mathbf{g}_n - \mathbf{g}^*\|_v^4 \gamma_{n+1}^2 \\ &\quad + 12\|\mathbf{g}_n - \mathbf{g}^*\|_v^2 \gamma_{n+1}^2 \langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v^2 \\ &\quad - 48\|\mathbf{g}_n - \mathbf{g}^*\|_v^2 \gamma_{n+1}^3 \langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v \\ &\quad + 48\|\mathbf{g}_n - \mathbf{g}^*\|_v^2 \gamma_{n+1}^4 \\ &\quad - 8\gamma_{n+1}^3 \langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v^3 \\ &\quad + 48\gamma_{n+1}^4 \langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v^2 \\ &\quad - 96\gamma_{n+1}^5 \langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v \\ &\quad + 2^6 \gamma_{n+1}^6. \end{aligned}$$

865 Taking the conditional expectation and already omitting some negative terms thanks to  
 866  $\langle \nabla H(\mathbf{g}_n), \mathbf{g}_n - \mathbf{g}^* \rangle_v \geq 0$ , which follows from the fact that  $H$  is convex and  $\mathbf{g}^*$  is a minimizer, gives  
 867 the simplification

$$\begin{aligned} \mathbb{E} [\|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^6 \mid \mathcal{F}_n] &\leq \|\mathbf{g}_n - \mathbf{g}^*\|_v^6 \\ &\quad - 6\|\mathbf{g}_n - \mathbf{g}^*\|_v^4 \gamma_{n+1} \langle \nabla H(\mathbf{g}_n), \mathbf{g}_n - \mathbf{g}^* \rangle_v \\ &\quad + 12\|\mathbf{g}_n - \mathbf{g}^*\|_v^4 \gamma_{n+1}^2 \\ &\quad + 12\|\mathbf{g}_n - \mathbf{g}^*\|_v^2 \gamma_{n+1}^2 \mathbb{E} [\langle \nabla h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v^2 \mid \mathcal{F}_n] \\ &\quad + 48\|\mathbf{g}_n - \mathbf{g}^*\|_v^2 \gamma_{n+1}^4 \\ &\quad + 48\gamma_{n+1}^4 \mathbb{E} [\langle \nabla h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v^2 \mid \mathcal{F}_n] \\ &\quad + 2^6 \gamma_{n+1}^6 \\ &\quad - 8\gamma_{n+1}^3 \mathbb{E} [\langle \nabla_{\mathbf{g}} h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v^3 \mid \mathcal{F}_n]. \end{aligned}$$

868 Using Proposition B.16, we have for any  $\mathbf{g} \in \mathcal{C}$ ,  $\langle \nabla H(\mathbf{g}_n), \mathbf{g}_n - \mathbf{g}^* \rangle_v \geq \eta \|\mathbf{g}_n - \mathbf{g}^*\|_v^2$ . The  
 869 Cauchy-Schwarz inequality gives  $|\langle \nabla h(\mathbf{g}_n, X_{n+1}), \mathbf{g}_n - \mathbf{g}^* \rangle_v| \leq 2\|\mathbf{g}_n - \mathbf{g}^*\|_v$ . Combining these  
 870 two inequalities we obtain

$$\begin{aligned} \mathbb{E} [\|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^6 \mid \mathcal{F}_n] &\leq \|\mathbf{g}_n - \mathbf{g}^*\|_v^6 (1 - 6\eta\gamma_{n+1}) + 60\gamma_{n+1}^2 \|\mathbf{g}_n - \mathbf{g}^*\|_v^4 + 240\gamma_{n+1}^4 \|\mathbf{g}_n - \mathbf{g}^*\|_v^2 \\ &\quad + 2^6 \gamma_{n+1}^6 + 64\gamma_{n+1}^3 \|\mathbf{g}_n - \mathbf{g}^*\|_v^3. \end{aligned}$$

871 Using Young's (generalized) inequality  $ab = ac \frac{b}{c} \leq \frac{(ac)^p}{p} + \frac{b^q}{c^q q}$  for  $c \neq 0, \frac{1}{p} + \frac{1}{q} = 1$  and applying  
 872 it to  $60\gamma_{n+1} \|\mathbf{g}_n - \mathbf{g}^*\|_v^4$  with  $c = \left(\frac{2}{3\eta}\right)^{2/3}, p = 3, q = \frac{3}{2}$  gives  $60\gamma_{n+1} \|\mathbf{g}_n - \mathbf{g}^*\|_v^4 \leq \frac{60^3 \gamma_{n+1}^3}{3}$ .

873  $\left(\frac{2}{3\eta}\right)^2 + \eta\|\mathbf{g}_n - \mathbf{g}^*\|_v^6$ . Analogously, one has  $64\gamma_{n+1}^2\|\mathbf{g}_n - \mathbf{g}^*\|_v^3 \leq \frac{2^{10}}{\eta}\gamma_{n+1}^4 + \eta\|\mathbf{g}_n - \mathbf{g}^*\|_v^6$ .  
 874 Thus, taking the expectation, we have

$$\mathbb{E}[\|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^6] \leq \mathbb{E}[\|\mathbf{g}_n - \mathbf{g}^*\|_v^6](1 - 4\eta\gamma_{n+1}) + \gamma_{n+1}^4 \left( 240 \cdot \text{Diam}_2(\mathcal{C})^2 + \frac{32}{\eta^2} \cdot 10^3 \right) + \frac{2^{10}}{\eta}\gamma_{n+1}^5 + 64\gamma_{n+1}^6,$$

875 where the terms involving  $\text{Diam}_2(\mathcal{C})$  appears from the crude bound  $\|\mathbf{g}_n - \mathbf{g}\|_v \leq \text{Diam}_2(\mathcal{C})$ .

876 Applying Lemma F.3 in a similar way as in the proof of Theorem 4.4 gives  $\mathbb{E}[\|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^6] \lesssim \frac{\gamma_n^3}{\eta^3}$ ,  
 877 so by Jensen's inequality, we conclude

$$\mathbb{E}[\|\mathbf{g}_{n+1} - \mathbf{g}^*\|_v^{2p}] \lesssim \frac{\gamma_n^p}{\eta^p} \quad \text{for } p \in \{1, 2, 3\}.$$

878 □

## 879 C.2 Convergence of the averaged iterates.

880 **Theorem 4.4** (Averaged iterates) Under Assumptions (A1) or (B1–B3), and assuming that  $f_\mu$  is  
 881  $\alpha$ -Hölder with  $\alpha \in (0, 1]$ , for any decay schedule of the form  $\gamma_n = \gamma_1/n^b$  with  $\gamma_1 > 0$  and  
 882  $b \in \left(\frac{1}{1+\alpha}, 1\right)$ , we have the convergence rate

$$\mathbb{E}[\|\bar{\mathbf{g}}_n - \mathbf{g}^*\|_v^2] = \mathcal{O}(1/n).$$

883 Without assuming  $f_\mu$  to be  $\alpha$ -Hölder, and for  $b \in (1/2, 1)$ , we still obtain

$$\mathbb{E}[\|\bar{\mathbf{g}}_n - \mathbf{g}^*\|_v^2] = \mathcal{O}(1/n^b).$$

884 *Proof.* For this proof, we introduce the additional following notation:

885 For any  $c > 0$  we define the function  $t \mapsto \Psi_c(t)$  such that

$$\sum_{t=1}^T t^{-c} \leq \Psi_c(T) := \begin{cases} 1 + \ln(T+1) & \text{if } c = 1, \\ \frac{2c-1}{c-1} & \text{if } c > 1, \\ 1 + \frac{1}{1-c}(T+1)^{1-c} & \text{if } c < 1. \end{cases} \quad (15)$$

886 We start by a decomposition of the gradient step, already present in [18]. We define the differences

$$\begin{aligned} p_k &:= \text{Proj}_{\mathcal{C}}(\mathbf{g}_k - \gamma_{k+1}\nabla_{\mathbf{g}}h(\mathbf{g}_k, X_{k+1})) - (\mathbf{g}_k - \gamma_{k+1}\nabla_{\mathbf{g}}h(\mathbf{g}_k, X_{k+1})), \\ \xi_{k+1} &:= \nabla H(\mathbf{g}_k) - \nabla_{\mathbf{g}}h(\mathbf{g}_k, X_{k+1}), \\ \delta_k &:= \nabla H(\mathbf{g}_k) - \nabla^2 H(\mathbf{g}^*)(\mathbf{g}_k - \mathbf{g}^*). \end{aligned}$$

887 Noting  $I_M$  the identity matrix in  $\mathcal{M}_M(\mathbb{R})$ , we observe that, by incorporating each introduced term  
 888 sequentially, for any  $k \in \mathbb{N}$ , we have

$$\begin{aligned} \mathbf{g}_{k+1} - \mathbf{g}^* &= \text{Proj}_{\mathcal{C}}(\mathbf{g}_k - \gamma_{k+1}\nabla_{\mathbf{g}}h(\mathbf{g}_k, X_{k+1})) - \mathbf{g}^* \\ &= \mathbf{g}_k - \gamma_{k+1}\nabla_{\mathbf{g}}h(\mathbf{g}_k, X_{k+1}) - \mathbf{g}^* - p_k \\ &= \mathbf{g}_k - \gamma_{k+1}\nabla H(\mathbf{g}_k, X_{k+1}) - \mathbf{g}^* + \gamma_{k+1}\xi_{k+1} - p_k \\ &= (I_M - \gamma_{k+1}\nabla^2 H(\mathbf{g}^*))(\mathbf{g}_k - \mathbf{g}^*) - \gamma_{k+1}\delta_k + \gamma_{k+1}\xi_{k+1} + p_k. \end{aligned}$$

889 Thus, we have that

$$\nabla^2 H(\mathbf{g}^*)(\mathbf{g}_k - \mathbf{g}^*) = \frac{\mathbf{g}_k - \mathbf{g}_{k+1}}{\gamma_{k+1}} - \delta_k + \xi_{k+1} + \frac{p_k}{\gamma_{k+1}}.$$

890 Observing that  $\frac{1}{n+1} \sum_{k=0}^n (\mathbf{g}_k - \mathbf{g}^*) = \bar{\mathbf{g}}_n - \mathbf{g}^*$ , we have the following decomposition of the  
 891 averaged iterates

$$\nabla^2 H(\mathbf{g}^*)(\bar{\mathbf{g}}_n - \mathbf{g}^*) = \frac{1}{n+1} \sum_{k=0}^n \frac{\mathbf{g}_k - \mathbf{g}_{k+1}}{\gamma_{k+1}} - \frac{1}{n+1} \sum_{k=0}^n \delta_k + \frac{1}{n+1} \sum_{k=0}^n \xi_{k+1} + \frac{1}{n+1} \sum_{k=0}^n \frac{p_k}{\gamma_{k+1}}.$$

892 We will now give the convergence rate of each sum.

893 • **Convergence rate for**  $\frac{1}{n+1} \sum_{k=0}^n \frac{\mathbf{g}_k - \mathbf{g}_{k+1}}{\gamma_{k+1}}$ .

$$\begin{aligned} \sum_{k=0}^n \frac{\mathbf{g}_k - \mathbf{g}_{k+1}}{\gamma_{k+1}} &= \sum_{k=0}^n \frac{(\mathbf{g}_k - \mathbf{g}^*) - (\mathbf{g}_{k+1} - \mathbf{g}^*)}{\gamma_{k+1}} \\ &= \sum_{k=0}^n \frac{\mathbf{g}_k - \mathbf{g}^*}{\gamma_{k+1}} - \sum_{k=0}^n \frac{\mathbf{g}_{k+1} - \mathbf{g}^*}{\gamma_{k+1}} \\ &= \sum_{k=1}^n \left( \frac{1}{\gamma_{k+1}} - \frac{1}{\gamma_k} \right) (\mathbf{g}_k - \mathbf{g}^*) + \frac{\mathbf{g}_0 - \mathbf{g}^*}{\gamma_1} - \frac{\mathbf{g}_{n+1} - \mathbf{g}^*}{\gamma_{n+1}}. \end{aligned}$$

894 Remark that  $\gamma_{n+1}^{-1} - \gamma_n^{-1} \leq 2\gamma_1^{-1}n^{b-1}$ . Using Minkowski's inequality and that, by Theorem 4.4  
895 (non-averaged iterates),  $\mathbb{E} [\|\mathbf{g}_n - \mathbf{g}^*\|_v^2] \lesssim \frac{\gamma_1}{\eta} (n+1)^{-b}$ ,

$$\mathbb{E} \left[ \left\| \sum_{k=0}^n \frac{\mathbf{g}_k - \mathbf{g}_{k+1}}{\gamma_{k+1}} \right\|_v^2 \right]^{\frac{1}{2}} \lesssim \frac{1}{\eta} \Psi_{1-b/2}(n+1) + \text{Diam}_2(\mathcal{C})\gamma_1^{-1} + \frac{1}{\sqrt{\gamma_1\eta}}(n+1)^{b/2}.$$

896 We thus have the convergence rate

$$\frac{1}{n+1} \mathbb{E} \left[ \left\| \sum_{k=0}^n \frac{\mathbf{g}_k - \mathbf{g}_{k+1}}{\gamma_{k+1}} \right\|_v^2 \right]^{\frac{1}{2}} \lesssim \frac{1}{\eta(n+1)^{1-b/2}}.$$

897 • **Convergence rate for**  $\frac{1}{n+1} \sum_{k=0}^n \delta_k$ .

898 We recall that  $\delta_k = \nabla H(\mathbf{g}_k) - \nabla^2 H(\mathbf{g}^*)(\mathbf{g}_k - \mathbf{g}^*)$  and that the Hessian using either Theorem B.3  
899 or Proposition 4.1, depending on our setting, there exists a ball  $B(\mathbf{g}^*, d_1)$  with  $d_1 > 0$  where  $H$  is  
900  $\alpha$ -Hölder. Therefore, applying Lemma F.4, if  $\mathbf{g}_k \in B(\mathbf{g}^*, d_1)$ , we have

$$\|\delta_k\| \lesssim \|\mathbf{g}_k - \mathbf{g}^*\|_v^{1+\alpha}.$$

901 Otherwise, since the Hessian, whose expression is provided in Proposition B.5, is uniformly bounded  
902 by an application of Lemma F.2, there exists a constant  $C_\delta$  such that for any  $\mathbf{g} \in \mathcal{C}$ ,  $\|\nabla H(\mathbf{g}) -$   
903  $\nabla^2 H(\mathbf{g}^*)(\mathbf{g} - \mathbf{g}^*)\| \leq C_\delta$ .

904 Since  $\mathbb{P}(\mathbf{g}_k \notin B(\mathbf{g}^*, d_1)) = \mathbb{P}(\|\mathbf{g}_k - \mathbf{g}^*\| > d_1)$ , using Markov's inequality gives

$$\begin{aligned} \mathbb{E}[\|\delta_k\|_v^2] &= \mathbb{E}[\|\delta_k\|_v^2 \mathbf{1}_{\mathbf{g}_k \in B(\mathbf{g}^*, d_1)}] + \mathbb{E}[\|\delta_k\|_v^2 \mathbf{1}_{\mathbf{g}_k \notin B(\mathbf{g}^*, d_1)}] \\ &\lesssim \mathbb{E}[\|\mathbf{g}_k - \mathbf{g}^*\|_v^{2+2\alpha}] + \frac{C_\delta^2}{d_1^{2+2\alpha}} \mathbb{E}[\|\mathbf{g}_k - \mathbf{g}^*\|_v^{2+2\alpha}] \\ &\lesssim \mathbb{E}[\|\mathbf{g}_k - \mathbf{g}^*\|_v^{2+2\alpha}]. \end{aligned}$$

905 Therefore, using Minkowski's inequality, we have

$$\begin{aligned} \frac{1}{n+1} \mathbb{E} \left[ \left\| \sum_{k=0}^n \delta_k \right\|_v^2 \right]^{\frac{1}{2}} &\lesssim \frac{1}{n+1} \sum_{k=0}^n \frac{1}{\eta^{\frac{1+\alpha}{2}} \gamma_{k+1}^{\frac{1+\alpha}{2}}} \\ &\leq \frac{1}{\eta^{\frac{1+\alpha}{2}} (n+1)} \Psi_{\frac{b+\alpha b}{2}} \\ &\lesssim \frac{1}{\eta^{\frac{1+\alpha}{2}} (n+1)^{\frac{b+\alpha b}{2}}}. \end{aligned}$$

906 • **Convergence rate for**  $\frac{1}{n+1} \sum_{k=0}^n \xi_{k+1}$ .

907 We recall that  $\xi_{k+1} = \nabla H(\mathbf{g}_k) - \nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1})$  and thus  $\mathbb{E}[\xi_{k+1}] = 0$ .

908 Observe that

$$\mathbb{E} \left[ \left\| \sum_{k=0}^n \xi_{k+1} \right\|_v^2 \right] = \mathbb{E} \left[ \left\| \sum_{k=0}^{n-1} \xi_{k+1} \right\|_v^2 + 2 \left\langle \sum_{k=0}^{n-1} \xi_{k+1}, \xi_{n+1} \right\rangle_v + \|\xi_{n+1}\|_v^2 \right]$$

909 with

$$\mathbb{E} \left[ \left\langle \sum_{k=0}^{n-1} \xi_{k+1}, \xi_{n+1} \right\rangle_v \right] = \mathbb{E} \left[ \left\langle \sum_{k=0}^{n-1} \xi_{k+1}, \mathbb{E}[\xi_{n+1} | \mathcal{F}_n] \right\rangle_v \right] = 0.$$

910 Thus, since for all  $k$ ,  $\mathbb{E}[\|\xi_k\|^2] \leq 4$ , we have the convergence rate

$$\frac{1}{n+1} \mathbb{E} \left[ \left\| \sum_{k=0}^n \xi_{k+1} \right\|_v^2 \right]^{\frac{1}{2}} \leq \frac{2}{\sqrt{n+1}}.$$

911 • **Convergence rate for**  $\frac{1}{n+1} \sum_{k=0}^n \frac{p_k}{\gamma_{k+1}}$ .

912 Take  $d_0$  such that  $B(\mathbf{g}^*, d_0) \subset \mathcal{C}$ . Using the notation  $\nabla_k := \nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1})$  for conciseness, we  
913 obtain

$$\begin{aligned} \mathbb{E}[\|p_k\|_v^2] &= \mathbb{E}[\|\text{Proj}_{\mathcal{C}}(\mathbf{g}_k - \gamma_{k+1} \nabla_k) - (\mathbf{g}_k - \gamma_{k+1} \nabla_k)\|_v^2] \\ &= \mathbb{E}[\|\text{Proj}_{\mathcal{C}}(\mathbf{g}_k - \gamma_{k+1} \nabla_k) - (\mathbf{g}_k - \gamma_{k+1} \nabla_k)\|_v^2 \mathbf{1}_{\mathbf{g}_k - \gamma_{k+1} \nabla_k \notin \mathcal{C}}] \end{aligned}$$

914 Since for any  $y \in \mathcal{C}$ , one has  $\|x - \text{Proj}_{\mathcal{C}}(x)\|_v \leq \|x - y\|_v$ , taking  $y = \mathbf{g}_k$ , and since  $\mathbf{g}_k - \gamma_{k+1} \nabla_k \notin \mathcal{C}$   
915 is satisfied only if  $\|\mathbf{g}_k - \gamma_{k+1} \nabla_k - \mathbf{g}^*\|_v > d_0$ , we have

$$\begin{aligned} \mathbb{E}[\|p_k\|_v^2] &\leq \mathbb{E}[\|\gamma_{k+1} \nabla_k\|_v^2 \mathbf{1}_{\|\mathbf{g}_k - \gamma_{k+1} \nabla_k - \mathbf{g}^*\|_v > d_0}] \\ &\leq 4\gamma_{k+1}^2 \frac{\mathbb{E}[\|\mathbf{g}_k - \gamma_{k+1} \nabla_k - \mathbf{g}^*\|_v^4]}{d_0^4} \\ &\leq \frac{\gamma_{k+1}^2}{d_0^4} \left( 2^5 \mathbb{E}[\|\mathbf{g}_k - \mathbf{g}^*\|_v^4] + 2^9 \gamma_{k+1}^4 \right) \\ &\lesssim \frac{1}{\eta^2} \gamma_{k+1}^4 \end{aligned}$$

916 Where we used Markov's inequality and the inequality  $(A + B)^4 \leq 2^3(A^4 + B^4)$ , for all  $A, B \in \mathbb{R}$

917 and that by Proposition C.1,  $\mathbb{E}[\|\mathbf{g}_n - \mathbf{g}^*\|^4] \lesssim \frac{\gamma_1^2}{\eta^2} (n+1)^{-2b}$ .

918 We thus have

$$\begin{aligned} \frac{1}{n+1} \mathbb{E} \left[ \left\| \sum_{k=0}^n \frac{p_k}{\gamma_{k+1}} \right\|_v^2 \right]^{\frac{1}{2}} &\lesssim \frac{1}{n+1} \sum_{k=0}^n \frac{\gamma_{k+1}}{\eta} \\ &\lesssim \frac{\gamma_1}{\eta(n+1)^b}. \end{aligned}$$

919 • **Conclusion.**

920 Using the convergence rate of all our terms, Cauchy-Schwarz inequality and that  $(A + B)^2 \leq$   
921  $2(A^2 + B^2)$  for all  $A, B \in \mathbb{R}$  we conclude

$$\begin{aligned} \mathbb{E}[\|\nabla^2 H(\mathbf{g}^*)(\bar{\mathbf{g}}_n - \mathbf{g}^*)\|_v^2] &= \mathbb{E} \left[ \left\| \frac{1}{n+1} \sum_{k=0}^n \frac{\mathbf{g}_k - \mathbf{g}_{k+1}}{\gamma_{k+1}} - \frac{1}{n+1} \sum_{k=0}^n \delta_k + \frac{1}{n+1} \sum_{k=0}^n \xi_{k+1} + \frac{1}{n+1} \sum_{k=0}^n \frac{p_k}{\gamma_{k+1}} \right\|_v^2 \right] \\ &\lesssim \frac{1}{\eta^2(n+1)^{2-b}} + \frac{1}{\eta^{\frac{1+\alpha}{\alpha}}(n+1)^{b+\alpha b}} + \frac{1}{n+1} + \frac{\gamma_1^4}{\eta^2(n+1)^{2b}}. \end{aligned}$$

922 Since  $b \in (\frac{1}{2}, 1)$  and  $b \in (\frac{1}{1+\alpha}, 1)$  the limiting term is  $\frac{1}{n+1}$  and we have

$$\mathbb{E} \left[ \left\| \nabla^2 H(\mathbf{g}^*) (\bar{\mathbf{g}}_n - \mathbf{g}^*) \right\|_v^2 \right] \lesssim \frac{1}{n+1}.$$

923 Finally, observe that there is an orthogonal matrix  $U$  such that  $\nabla^2 H(\mathbf{g}^*) =$   
 924  $U \text{diag}(\lambda_1, \dots, \lambda_{M-1}, 0) U^\top$ . Therefore, noting by abuse of notation

$$(\nabla^2 H(\mathbf{g}^*))^{-1} = U \text{diag}(\lambda_1^{-1}, \dots, \lambda_{M-1}^{-1}, 0) U^\top$$

925 the inverse of  $\nabla_k^2$  in the space  $\text{Vect}(\mathbf{1}_M)^\perp$  we finally have

$$\mathbb{E}[\|\bar{\mathbf{g}}_n - \mathbf{g}^*\|_v^2] \lesssim \frac{1}{\lambda^2(n+1)}$$

926 where  $\lambda = \min_{j \in \llbracket 1, M-1 \rrbracket} \lambda_j > 0$  by either Theorem B.3 or Proposition B.12.  $\square$

### 927 C.3 Convergence rate of PSGD in the general setting

928 **Theorem 3.2** (PSGD in the general setting) Assuming that the semi-dual problem (4) admits at least  
 929 one solution  $\mathbf{g}^*$  and that there exists a compact set  $K$  such that  $\mu(K) \geq 1 - \frac{1}{2}w_{\min}$ , choosing the  
 930 learning rate  $\gamma_n = \gamma_1/n^b$  with  $\gamma_1 = \frac{\text{Diam}(\mathcal{C})}{2\sqrt{2}}$  and  $b = \frac{1}{2}$ , we obtain

$$\mathbb{E}[H(\bar{\mathbf{g}}_n) - H(\mathbf{g}^*)] \leq \frac{4\sqrt{2} \text{Diam}(\mathcal{C})}{\sqrt{n}}.$$

931 *Proof.* Define  $\gamma_k = \frac{\gamma_1}{\sqrt{k}}$  for  $\gamma_1 > 0, k \geq 1$  and denote by  $\mathbf{g}^* \in \mathcal{C}$  a minimizer of the functional  $H$ .  
 932 Thanks to Jensen's inequality coupled with the fact that no matter  $\mathbf{g} \in \mathbb{R}^M : H(\mathbf{g}) - H(\mathbf{g}^*) \leq$   
 933  $\nabla H(\mathbf{g})^\top (\mathbf{g} - \mathbf{g}^*)$ , it comes

$$\mathbb{E}[H(\bar{\mathbf{g}}_n) - H(\mathbf{g}^*)] \leq \mathbb{E} \left[ \frac{1}{n+1} \sum_{k=0}^n H(\mathbf{g}_k) - H(\mathbf{g}^*) \right] \leq \frac{1}{n+1} \mathbb{E} \left[ \sum_{k=0}^n \nabla H(\mathbf{g}_k)^\top (\mathbf{g}_k - \mathbf{g}^*) \right].$$

934 Then, since no matter  $k$ ,  $X_{k+1} \sim \mu$ , we have  $\mathbb{E}[\nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1}) \mid \mathcal{F}_k] = \nabla H(\mathbf{g}_k)$ , we have

$$\mathbb{E}[H(\bar{\mathbf{g}}_n) - H(\mathbf{g}^*)] \leq \frac{1}{n+1} \mathbb{E} \left[ \sum_{k=0}^n \nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1})^\top (\mathbf{g}_k - \mathbf{g}^*) \right]. \quad (16)$$

935 We will proceed to bound the right hand side of this inequality.

936 By definition of  $\mathbf{g}_{k+1}$ , and since  $\mathbf{g}^* \in \mathcal{C}$  and  $\text{Proj}_{\mathcal{C}}$  is 1-Lipschitz, it comes

$$\begin{aligned} \mathbb{E}[\|\mathbf{g}_{k+1} - \mathbf{g}^*\|^2] &\leq \mathbb{E}[\|\text{Proj}_{\mathcal{C}}(\mathbf{g}_k - \gamma_{k+1} \nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1})) - \mathbf{g}^*\|^2] \\ &\leq \mathbb{E}[\|\mathbf{g}_k - \gamma_{k+1} \nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1}) - \mathbf{g}^*\|^2] \\ &\leq \mathbb{E}[\|\mathbf{g}_k - \mathbf{g}^*\|^2 + \gamma_{k+1}^2 \|\nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1})\|^2 - 2\gamma_{k+1} \nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1})^\top (\mathbf{g}_k - \mathbf{g}^*)]. \end{aligned}$$

937 In addition, since  $\|\nabla_{\mathbf{g}} h(\mathbf{g}, X)\| \leq 2$  a.s. no matter  $\mathbf{g} \in \mathbb{R}^M$  and  $X \in \mathbb{R}^d$ ,

$$\mathbb{E}[2\gamma_{k+1} \nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1})^\top (\mathbf{g}_k - \mathbf{g}^*)] \leq \mathbb{E}[\|\mathbf{g}_k - \mathbf{g}^*\|^2 - \|\mathbf{g}_{k+1} - \mathbf{g}^*\|^2 + 4\gamma_{k+1}^2].$$

938 Then, with the help of Abel's summation formula,

$$\begin{aligned} 2\mathbb{E} \left[ \sum_{k=0}^n \nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1})^\top (\mathbf{g}_k - \mathbf{g}^*) \right] &\leq \mathbb{E} \left[ \sum_{k=0}^n \frac{\|\mathbf{g}_k - \mathbf{g}^*\|^2 - \|\mathbf{g}_{k+1} - \mathbf{g}^*\|^2}{\gamma_{k+1}} \right] + 4 \sum_{k=0}^n \gamma_{k+1} \\ &\leq \mathbb{E} \left[ \sum_{k=1}^n \|\mathbf{g}_k - \mathbf{g}^*\|^2 \left( \frac{1}{\gamma_{k+1}} - \frac{1}{\gamma_k} \right) \right] - \frac{\|\mathbf{g}_{k+1} - \mathbf{g}^*\|^2}{\gamma_{k+1}} + \frac{\|\mathbf{g}_0 - \mathbf{g}^*\|^2}{\gamma_1} + 4 \sum_{k=0}^n \gamma_{k+1} \end{aligned}$$

939 Then, since for all  $k$ ,  $\|\mathbf{g}_k - \mathbf{g}^*\| \leq \text{Diam}_2(\mathcal{C})$ , it comes

$$\begin{aligned} 2\mathbb{E} \left[ \sum_{k=0}^n \nabla_{\mathbf{g}} h(\mathbf{g}_k, X_{k+1})^\top (\mathbf{g}_k - \mathbf{g}^*) \right] &\leq \text{Diam}_2(\mathcal{C})^2 \sum_{k=1}^n \left( \frac{1}{\gamma_{k+1}} - \frac{1}{\gamma_k} \right) + \frac{D^2}{\gamma_1} + 4 \sum_{k=1}^n \gamma_k \\ &\leq \frac{\text{Diam}_2(\mathcal{C})^2}{\gamma_{n+1}} + 4 \sum_{k=1}^n \gamma_k \\ &\leq \frac{\text{Diam}_2(\mathcal{C})^2}{\gamma_1} \sqrt{n+1} + 8\gamma_1 \sqrt{n+1} \end{aligned}$$

940 Using the inequality (16) we obtain

$$\mathbb{E} [H(\bar{g}_n) - H(\mathbf{g}^*)] \leq \frac{1}{\sqrt{n+1}} \left( \frac{\text{Diam}_2(\mathcal{C})^2}{\gamma_1} + 8\gamma_1 \right).$$

941 The best constant  $\gamma_1$  is then  $\gamma_1 = \frac{\text{Diam}_2(\mathcal{C})}{2\sqrt{2}}$ , but no matter  $\gamma_1 > 0$ , we have the desired convergence  
942 rate

$$\mathbb{E} [H(\bar{g}_n) - H(\mathbf{g}^*)] = \mathcal{O}(1/\sqrt{n}).$$

943

□

## 944 D Proof of Lemma 3.1: Localisation of a projection set

945 **Lemma 3.1** (Existence of a projection set) As soon as the semi-dual problem is well-posed, there  
946 exists a minimizer  $\mathbf{g}^*$  in the set

$$\mathcal{C} := \{\mathbf{g} \in \{0\} \times \mathbb{R}^{M-1} \mid |g_j| \leq \|c\|_{K,\infty}\}$$

947 where  $\|c\|_{K,\infty} := \sup_{x \in K, j \in \llbracket 1, M \rrbracket} |c(x, y_j)|$  for any compact set  $K$  satisfying  $\mu(K) \geq 1 - \frac{1}{2} \min w_j$ .

948 *Proof.* By the first order condition, the minimizer of  $H$  satisfies  $\mu(\mathbb{L}_j(\mathbf{g})) = w_j$  for all  $j \in \llbracket 1, M \rrbracket$ .  
949 In particular, one can restrict the search set for an optimal potential to the set of potentials defined by

$$\mathcal{L} := \left\{ \mathbf{g} \in \mathbb{R}^M : \forall j \in \llbracket 1, M \rrbracket, \mu(\mathbb{L}_j(\mathbf{g})) \geq \frac{2}{3} w_{\min} \right\}.$$

950 Let us show that this set is contained in an  $L^\infty$  ball with an explicit radius. Consider any compact set  
951  $K$  such that  $\mu(K) \geq 1 - \frac{1}{2} w_{\min}$ . For  $\mathbf{g} \in \mathcal{L}$  and any  $j \in \llbracket 1, M \rrbracket$  we get that

$$\begin{aligned} \mu(\mathbb{L}_j(\mathbf{g}) \cap K) &= 1 - \mu((\mathbb{R}^M \setminus \mathbb{L}_j(\mathbf{g})) \cup (\mathbb{R}^M \setminus K)) \\ &\geq 1 - \mu((\mathbb{R}^M \setminus \mathbb{L}_j(\mathbf{g}))) - \mu((\mathbb{R}^M \setminus K)) \\ &= 1 - (1 - \mu(\mathbb{L}_j(\mathbf{g}))) - (1 - \mu(K)) \\ &\geq \frac{w_{\min}}{6} > 0, \end{aligned}$$

952 and so  $\mathbb{L}_j(\mathbf{g}) \cap K \neq \emptyset$ . In particular, for every  $j \in \llbracket 1, M \rrbracket$ , there exists  $x_j \in \mathbb{L}_j(\mathbf{g}) \cap K$  and so for  
953 all  $i \in \llbracket 1, M \rrbracket$

$$c(x_j, y_j) - g_j \leq c(x_j, y_i) - g_i.$$

954 Therefore, using the fact that the cost is non-negative, we have

$$\begin{aligned} \max_i g_i - \min_j g_j &\leq \max_i \max_j \{c(x_j, y_i) - c(x_j, y_j)\} \\ &\leq \max_{i,j \in \llbracket 1, M \rrbracket} \sup_{x \in K} c(x, y_i) = \|c\|_{K,\infty}. \end{aligned}$$

955 Moreover, since  $H(\mathbf{g} + \lambda \mathbb{1}_M) = H(\mathbf{g})$  one can fix  $g_1 = 0$ , which concludes the proof. □

## E Minimax estimation of OT quantities

Consider  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  with densities  $f_P$  and  $f_Q$ . We recall that the Hellinger distance is defined by

$$d_H(P, Q) := \left( \int_{\mathbb{R}^d} \left( \sqrt{f_P(x)} - \sqrt{f_Q(x)} \right)^2 d\lambda_{\mathbb{R}^d}(x) \right)^{\frac{1}{2}}. \quad (17)$$

We also recall the formulation of Le Cam's Lemma. We refer to [32], Chapter 15, for more details.

**Lemma E.1.** (*Le Cam's Lemma.*) Let  $\mathcal{P}$  be a set of probability distributions on a measurable space, and consider the problem of estimating a parameter  $\theta \in \Theta$  with a loss function  $\ell$  defined, for all  $\hat{\theta}, \theta \in \Theta$ , as

$$\ell(\hat{\theta}, \theta) = d(\hat{\theta}, \theta)^p,$$

where  $p \geq 1$  is an integer, and  $d$  is a distance on  $\Theta$ . Then, for all  $\theta_1, \theta_2 \in \Theta$ ,

$$R_M = \inf_{T^{(n)}} \sup_{\theta \in \Theta} \mathbb{E}_\theta[\ell(\theta, T^{(n)}(X))] \geq \frac{1}{2^p} (1 - \sqrt{nd_H(P_{\theta_1}, P_{\theta_2})}) d(\theta_1, \theta_2)^p,$$

where  $R_M$  is the minimax risk, and  $T^{(n)}$  is an estimator based on  $n$  i.i.d samples from  $P_\theta$ .

In the sequel, we note  $\mathcal{P}_H$  the set of absolute continuous measures that verify Assumption A.

### E.1 Kantorovich potential

*Proof.* Consider  $\nu = \frac{1}{2}\delta_{\{0\}} + \frac{1}{2}\delta_{\{1\}}$  and fix the cost to transfer mass to be the usual quadratic cost  $c(x, y) = \frac{1}{2}\|x - y\|^2$ .

For  $\delta \geq 0$ , we define  $\mu_{\theta^\delta} = \mathcal{N}(\delta, 1)$ . Since  $d = 1$ , the optimal transport map is monotone non-decreasing (see, for instance, Chapter 2 in [30]). Thus, we must have the identity

$$T_{\mu_{\theta^\delta}, \nu}(x) = \begin{cases} 0, & \text{if } x \leq \delta, \\ 1, & \text{otherwise.} \end{cases}$$

Therefore the vector  $\theta^\delta \in \mathbb{R}^2$  solves the semi-dual problem if and only if it satisfies the following inequalities

$$\begin{aligned} c(x, 0) - \theta_1^\delta &\leq c(x, 1) - \theta_2^\delta, \quad \forall x \leq \delta, \\ c(x, 1) - \theta_2^\delta &\leq c(x, 0) - \theta_1^\delta, \quad \forall x > \delta. \end{aligned}$$

Since  $c(x, y) = \frac{1}{2}\|x - y\|^2$  we can fix  $\theta_1^\delta = 0$  and compute

$$\theta^\delta = \left( 0, \frac{1}{2} - \delta \right).$$

Therefore, we parameterized the family of probability measures  $\mu_\theta \in \mathcal{P}_H$  so that the couple  $(\theta^c, \theta)$  is the unique solution in  $\Theta \subset \{\theta = (\theta_1, \theta_2) \in \mathbb{R}^2; \theta_1 = 0\}$  of the dual of  $\text{OT}(\mu_\theta, \nu)$ . In this class of probabilities, the minimax estimation of the optimal transport potential  $\theta$ , given  $n > 0$  i.i.d samples of the source measure, can be written as

$$R_n^\Theta := \inf_{\hat{\theta}^{(n)}} \sup_{\theta \in \Theta} \mathbb{E}_{\mu_\theta} [\|\hat{\theta}^{(n)} - \theta\|^2],$$

where  $\hat{\theta}^{(n)}$  is based on the  $n$  i.i.d samples from the source measure  $\mu$ . Note that

$$R_n^\Theta \leq \inf_{\mathbf{g}^{(n)}} \sup_{\mu \in \mathcal{P}_H} \mathbb{E}_\mu [\|\mathbf{g}^{(n)} - \mathbf{g}^*\|^2], \quad (18)$$

where the infimum is taken over all vectors  $\mathbf{g}^{(n)}$  based on  $n$  i.i.d samples of  $\mu$ .



Using the closed form of the Hellinger distance between Gaussian distributions, we have

$$d_H(\mu_{\theta^0}, \mu_{\theta^\delta}) = \sqrt{1 - \exp\left(-\frac{\delta^2}{8}\right)}.$$

979 For  $\delta \rightarrow 0$ , the Taylor expansion gives  $d_H(\mu_{\theta^0}, \mu_{\theta^\delta}) = \delta/\sqrt{8} + o(\delta)$ . Applying Le Cam's Lemma  
980 with  $\delta = 1/\sqrt{n}$  gives

$$\begin{aligned} R_n^\Theta &\geq \frac{1}{4} (1 - \sqrt{n} d_H(\mu_{\theta^0}, \mu_{\theta^\delta})) \|\theta^0 - \theta^\delta\|_2^2 \\ &\geq \frac{1}{4} \left(1 - 1/\sqrt{8} + o(1)\right) \frac{1}{n} \\ &\geq \frac{1}{10n} + o\left(\frac{1}{n}\right) \end{aligned}$$

981 Using the inequality (18) concludes the proof.  $\square$

## 982 E.2 OT map

983 *Proof.* Define for  $\delta \in [0, 1]$  the set of probability measures  $\mu_\delta$ , with density:

$$f_{\mu_\delta}(x) = \mathbf{1}_{x \in [0,1]} (1 + \delta g(x)), \quad x \in [0, 1],$$

984 where

$$g(x) = \begin{cases} 2(1 - 2x), & x \in [0, 1/2], \\ -2(2x - 1), & x \in [1/2, 1]. \end{cases}$$

985 The squared Hellinger distance between  $\mu_0$  (uniform) and  $\mu_\delta$  is:

$$\begin{aligned} d_H(\mu_0, \mu_\delta)^2 &= \frac{1}{2} \int_0^1 (\sqrt{1} - \sqrt{1 + \delta g(x)})^2 dx \\ &= \frac{1}{6\delta} \left( (1 - 2\delta)^{3/2} + 6\delta - (1 + 2\delta)^{3/2} \right) \\ &= \frac{1}{2} \delta^2 + o(\delta^2) \end{aligned}$$

986 when  $\delta \rightarrow 0$ . Therefore, we obtain

$$d_H(\mu_0, \mu_\delta) = \frac{1}{\sqrt{2}} \delta + o(\delta).$$

987 Since  $\mu_\delta \in \mathcal{P}_H$ ,  $\delta \in [0, 1]$ , we have

$$\inf_{T^{(n)}} \sup_{\mu \in \mathcal{P}_H} \mathbb{E}_\mu \left[ \left\| T^{(n)} - T_{\mu, \nu} \right\|_{L^p(\mu)}^p \right] \geq \inf_{T^{(n)}} \sup_{\delta \in [0,1]} \mathbb{E}_{\mu_\delta} \left[ \left\| T^{(n)} - T_{\mu_\delta, \nu} \right\|_{L^p(\mu_\delta)}^p \right]. \quad (19)$$

988 From the relation  $1 - 2\delta \leq f_{\mu_\delta}(x) \leq 1 + 2\delta, \forall x \in [0, 1]$ , we infer that no matter  $T^{(n)}$  and  $\delta$  small  
989 enough

$$\left\| T^{(n)} - T_{\mu_\delta, \nu} \right\|_{L^p(\mu_\delta)}^p \geq \frac{1}{2} \left\| T^{(n)} - T_{\mu_\delta, \nu} \right\|_{L^p(\mathcal{U}(0,1))}^p. \quad (20)$$

990 Consider  $\nu = \frac{1}{2} \delta_{\{0\}} + \frac{1}{2} \delta_{\{1\}}$  as above. Recall that the optimal transport map is monotone non-  
991 decreasing in dimension 1. Moreover by definition  $T_{\mu_\delta, \nu}(x) \in \{0, 1\}, x \in [0, 1]$ , and  $\mu_\delta(T_{\mu_\delta, \nu}^{-1}(0)) =$   
992  $1/2$ . Therefore one identifies  $T_{\mu_\delta, \nu} = \mathbf{1}_{x \geq M_\delta}$  where  $M_\delta$  is the median of  $\mu_\delta, \delta > 0$ , satisfying

$$\int_0^{M_\delta} (1 + \delta g(x)) dx = \frac{1}{2}.$$

993 Noticing that  $M_\delta \in [0, 1/2]$ , we solve

$$\int_0^{M_\delta} (1 + 2\delta(1 - 2x)) dx = \frac{1}{2},$$

994 which gives the solution

$$M_\delta = \frac{1 + 2\delta - \sqrt{1 + 4\delta^2}}{4\delta} = \frac{1}{2} - \frac{1}{2}\delta + o(\delta).$$

995 Observe that

$$\|T_{\mu_0, \nu} - T_{\mu_\delta, \nu}\|_{L^p(\mathcal{U}(0,1))}^p = \int_0^1 |\mathbf{1}_{x \geq M_0} - \mathbf{1}_{x \geq M_\delta}|^p dx = |M_\delta - M_0|.$$

996 We proved the relation  $\|T_{\mu_0, \nu} - T_{\mu_\delta, \nu}\|_{L^p(\mathcal{U}(0,1))}^p = \delta + o(\delta)$ . Applying Le Cam's Lemma with

997  $\delta = \frac{1}{\sqrt{n}}$  for  $n$  sufficiently large and any  $p \in [1, \infty)$ , we obtain

$$\begin{aligned} \inf_{T^{(n)}} \sup_{\delta \in [0,1]} \mathbb{E}_{\mu_\delta} \left[ \|T^{(n)} - T_{\mu_\delta, \nu}\|_{L^p(\mathcal{U}(0,1))}^p \right] &\geq \frac{1}{2^p} (1 - \sqrt{n} d_H(\mu_0, \mu_\delta)) \|T_{\mu_0, \nu} - T_{\mu_\delta, \nu}\|_{L^p(\mathcal{U}(0,1))}^p \\ &\geq \frac{1}{2^p} \left( 1 - \frac{1}{\sqrt{2}} + o(1) \right) \left( \frac{1}{\sqrt{n}} + o\left(\frac{1}{\sqrt{n}}\right) \right) \\ &\gtrsim \frac{1}{\sqrt{n}}. \end{aligned}$$

998 Therefore, combining inequalities (19) and (20) we conclude

$$\inf_{T^{(n)}} \sup_{\mu \in \mathcal{P}_H} \mathbb{E}_\mu \left[ \|T^{(n)} - T_{\mu, \nu}\|_{L^p(\mu)}^p \right] \gtrsim \frac{1}{\sqrt{n}}.$$

999

□

## 1000 F Technical Lemmas

### 1001 F.1 Technical Lemmas for Appendix B

1002 **Lemma F.1.** *Perturbation of Laplacian Matrices. Let  $A$  and  $B$  be symmetric Laplacian matrices of*  
 1003 *the same size such that:*

$$A_{ij} \leq 0, \quad B_{ij} \leq A_{ij} \quad \text{for all } i \neq j.$$

1004 Suppose  $\lambda_2(A) > 0$ , where  $\lambda_2(A)$  denotes the second smallest eigenvalue of  $A$ . Then:

$$\lambda_2(B) \geq \lambda_2(A)$$

1005 where  $\lambda_2(B)$  is the second smallest eigenvalue of  $B$ .

1006 *Proof.* We recall the variational characterization of the second smallest eigenvalue of a Laplacian  
 1007 matrix  $M$  is

$$\lambda_2(M) = \min_{x \perp \mathbf{1}} \frac{x^T M x}{x^T x}.$$

1008 Define the matrix  $C = B - A$ . For  $i \neq j$ :

$$C_{ij} = B_{ij} - A_{ij} \leq 0,$$

1009 and for diagonal elements

$$C_{ii} = B_{ii} - A_{ii} = - \sum_{j \neq i} C_{ij} \geq 0.$$

1010 Thus,  $C$  is a Laplacian matrix and so it is positive, semi-definite.

1011 Let  $y_2$  be the eigenvector corresponding to  $\lambda_2(B)$ , the second smallest eigenvalue of  $B$ , that is

$$By_2 = \lambda_2(B)y_2 \quad \text{with } y_2 \perp \mathbf{1}.$$

1012 Since  $A = B - C$ , we have

$$y_2^T A y_2 = y_2^T (B - C) y_2 = y_2^T B y_2 - y_2^T C y_2.$$

1013 Thus,

$$\frac{y_2^T A y_2}{y_2^T y_2} = \lambda_2(B) - \frac{y_2^T C y_2}{y_2^T y_2}.$$

1014 By the variational principle:

$$\lambda_2(A) = \min_{x \perp \mathbf{1}} \frac{x^T A x}{x^T x} \leq \frac{y_2^T A y_2}{y_2^T y_2} \leq \lambda_2(B) - \frac{y_2^T C y_2}{y_2^T y_2}.$$

1015 Since  $C$  is a Laplacian matrix, we have  $y_2^T C y_2 \geq 0$ , no matter  $y_2$  and thus  $\lambda_2(A) \leq \lambda_2(B)$  which  
1016 completes the proof.

1017 □

1018 **Lemma F.2.** Let  $f_\mu$  satisfy, for all  $x \in \mathbb{R}^d$ , the decay condition

$$\sum_{r \geq 1} r^{d-1} \sup_{x \in \mathbb{R}^d \setminus B(0, r-1)} f_\mu(x) < \infty.$$

1019 Then, there exists a constant  $C > 0$  such that, for any hyperplane  $H \subset \mathbb{R}^d$ , we have

$$\int_H f_\mu(x) d\mathcal{H}^{d-1}(x) \leq C.$$

1020 *Proof.* Defining for any  $r \in \mathbb{N}^*$ ,  $R(r) := B(0, r) \setminus B(0, r-1)$ , we have

$$\begin{aligned} \int_H f_\mu(x) d\mathcal{H}^{d-1}(x) &= \sum_{r \geq 1} \int_{H \cap R(r)} f_\mu(x) d\mathcal{H}^{d-1}(x) \\ &\leq \sum_{r \geq 1} \int_{H \cap R(r)} \sup_{x \in R(r)} f_\mu(x) d\mathcal{H}^{d-1}(x) \\ &\leq \sum_{r \geq 1} \int_{H \cap R(r)} \sup_{x \in \mathbb{R}^d \setminus B(0, r-1)} f_\mu(x) d\mathcal{H}^{d-1}(x) \end{aligned}$$

1021 Then, using the fact that  $\mathcal{H}^{d-1}(H \cap R(r)) \leq \mathcal{H}^{d-1}(H \cap B(0, r))$ , and noting that  $H \cap B(0, r)$  is  
1022 a  $(d-1)$ -dimensional ball of radius  $r$ , we have

$$\mathcal{H}^{d-1}(H \cap R(r)) \leq \frac{\pi^{(d-1)/2}}{\Gamma\left(\frac{d+1}{2}\right)} r^{d-1}.$$

1023 Therefore, incorporating this bound, we obtain:

$$\int_H f_\mu(x) d\mathcal{H}^{d-1}(x) \leq \sum_{r \geq 1} \frac{\pi^{(d-1)/2}}{\Gamma\left(\frac{d+1}{2}\right)} r^{d-1} \sup_{x \in \mathbb{R}^d \setminus B(0, r)} f_\mu(x)$$

1024 which is finite by our decay assumption on  $f_\mu$ .

1025 □

## 1026 F.2 Technical Lemmas for Appendix C

1027 **Lemma F.3.** *Let  $(\gamma_n)_{n \geq 0}$  and  $(m_n)_{n \geq 0}$  be some positive and decreasing sequences and let  $(\delta_n)_{n \geq 0}$ ,*  
 1028 *satisfying the following:*

1029 • *The sequence  $\delta_n$  follows the recursive relation:*

$$\delta_{n+1} \leq (1 - \omega\gamma_{n+1})\delta_n + m_{n+1}\gamma_{n+1}, \quad (21)$$

1030 *with  $\delta_0 \geq 0$  and  $\omega > 0$ .*

1031 •  *$\gamma_n$  converges to 0.*

1032 • *Let  $n_0 = \inf \{n \geq 1 : \omega\gamma_{n+1} \leq 1\}$ ,  $\delta_{n_0}$  is non-negative.*

1033 *Then, for all  $n \geq n_0$ , we have the upper bound:*

$$\delta_n \leq \exp\left(-\omega \sum_{i=n_0+1}^n \gamma_i\right) \left(\sum_{k=n_0}^n \gamma_k m_k + \delta_{n_0}\right) + \frac{1}{\omega} m_{\lceil \frac{n}{2} \rceil - 1}$$

1034 *Proof.* For all  $n \geq n_0$ , one has

$$\delta_n \leq \underbrace{\prod_{i=n_0+1}^n (1 - \omega\gamma_i) \delta_{n_0}}_{=: \tilde{U}_{1,n}} + \underbrace{\sum_{k=n_0+1}^n \prod_{i=k+1}^n (1 - \omega\gamma_i) \gamma_k m_k}_{=: \tilde{U}_{2,n}}$$

1035 One can consider two cases:  $\lceil n/2 \rceil - 1 \leq n_0$  and  $\lceil n/2 \rceil - 1 > n_0$ .

1036 **Case where  $\lceil n/2 \rceil - 1 \leq n_0 < n$ :** Since  $m_k$  is decreasing,

$$\begin{aligned} U_{2,n} &\leq m_{n_0+1} \sum_{k=n_0+1}^n \prod_{i=k+1}^n (1 - \omega\gamma_i) \gamma_k \\ &= \frac{1}{\omega} m_{n_0+1} \sum_{k=n_0+1}^n \prod_{i=k+1}^n (1 - \omega\gamma_i) - \prod_{i=k}^n (1 - \omega\gamma_i) \\ &= \frac{1}{\omega} m_{n_0+1} \left(1 - \prod_{i=n_0+1}^n (1 - \omega\gamma_i)\right) \\ &\leq \frac{1}{\omega} m_{n_0+1} \end{aligned}$$

1037 Since  $m_k$  is decreasing, it comes  $U_{2,n} \leq \frac{1}{\omega} m_{\lceil n/2 \rceil}$ .

1038 **Case where  $\lceil n/2 \rceil - 1 > n_0$ :** As in [4], for all  $m = n_0 + 1, \dots, n$ , one has

$$U_{2,n} \leq \exp\left(-\omega \sum_{k=m+1}^n \gamma_k\right) \sum_{k=n_0+1}^m \gamma_k m_k + \frac{1}{\omega} m_m$$

1039 Then, taking  $m = \lceil n/2 \rceil - 1$ , it comes

$$U_{2,n} \leq \exp\left(-\omega \sum_{k=\lceil n/2 \rceil}^n \gamma_k\right) \sum_{k=n_0+1}^{\lceil n/2 \rceil - 1} \gamma_k m_k + \frac{1}{\omega} m_{\lceil n/2 \rceil - 1}$$

1040 □

1041 **Lemma F.4.** *Linearization of the Hessian If a function  $H : \mathbb{R}^M \rightarrow \mathbb{R}$  is such that its Hessian is*  
 1042  *$\alpha$ -Hölder on the ball  $B(0, r)$ , with  $r > 0, \alpha \in (0, 1)$  and constant  $L$ , then for any  $\mathbf{g}, \mathbf{g}^* \in B(0, r)$ ,*  
 1043 *we have*

$$\|\nabla H(\mathbf{g}) - \nabla^2 H(\mathbf{g}^*)(\mathbf{g} - \mathbf{g}^*)\| \leq C \|\mathbf{g} - \mathbf{g}^*\|^{1+\alpha},$$

1044 *where  $C = \frac{L}{\alpha+1}$ .*

1045 *Proof.* Consider the Taylor expansion of  $\nabla H(\mathbf{g})$  around  $\mathbf{g}^*$ :

$$\nabla H(\mathbf{g}) = \nabla H(\mathbf{g}^*) + \nabla^2 H(\mathbf{g}^*)(\mathbf{g} - \mathbf{g}^*) + R(\mathbf{g}),$$

1046 where the remainder term  $R(\mathbf{g})$  is given by:

$$R(\mathbf{g}) = \int_0^1 [\nabla^2 H(\mathbf{g}^* + t(\mathbf{g} - \mathbf{g}^*)) - \nabla^2 H(\mathbf{g}^*)] (\mathbf{g} - \mathbf{g}^*) dt.$$

1047 By the assumption of  $\alpha$ -Hölder continuity of the Hessian, we have

$$\|\nabla^2 H(\mathbf{g}^* + t(\mathbf{g} - \mathbf{g}^*)) - \nabla^2 H(\mathbf{g}^*)\| \leq L\|t(\mathbf{g} - \mathbf{g}^*)\|^\alpha = Lt^\alpha \|\mathbf{g} - \mathbf{g}^*\|^\alpha.$$

1048 Thus,

$$\|R(\mathbf{g})\| \leq \int_0^1 Lt^\alpha \|\mathbf{g} - \mathbf{g}^*\|^\alpha \|\mathbf{g} - \mathbf{g}^*\| dt = L\|\mathbf{g} - \mathbf{g}^*\|^{1+\alpha} \int_0^1 t^\alpha dt.$$

1049 Evaluating the integral:

$$\int_0^1 t^\alpha dt = \frac{1}{\alpha + 1}.$$

1050 Therefore,

$$\|R(\mathbf{g})\| \leq \frac{L}{\alpha + 1} \|\mathbf{g} - \mathbf{g}^*\|^{1+\alpha}.$$

1051 This implies the desired inequality:

$$\|\nabla H(\mathbf{g}) - \nabla^2 H(\mathbf{g}^*)(\mathbf{g} - \mathbf{g}^*)\| \leq C\|\mathbf{g} - \mathbf{g}^*\|^{1+\alpha},$$

1052 where  $C = \frac{L}{\alpha+1}$ . □

1053